

Advances in the dataflow computational model

Walid A. Najjar^{a,*}, Edward A. Lee^{b,2}, Guang R. Gao^c

^a Department of Computer Science, Colorado State University, FT Collins, CO 80623, USA

^b Department of Electrical Engineering and Computer Science, University of California, Berkeley, USA

^c Department of Electrical Engineering, University of Delaware, USA

Abstract

The dataflow program graph execution model, or dataflow for short, is an alternative to the stored-program (von Neumann) execution model. Because it relies on a graph representation of programs, the strengths of the dataflow model are very much the complements of those of the stored-program one. In the last thirty or so years since it was proposed, the dataflow model of computation has been used and developed in very many areas of computing research: from programming languages to processor design, and from signal processing to reconfigurable computing. This paper is a review of the current *state-of-the-art* in the applications of the dataflow model of computation. It focuses on three areas: multithreaded computing, signal processing and reconfigurable computing. © 1999 Elsevier Science B.V. All rights reserved.

Keywords: Computational models; Dataflow; Multithreaded computer architecture; von Neumann computer; Dataflow history; Memory models

1. Introduction

Dataflow is a sound, simple, and powerful model of parallel computation. In dataflow programming and architectures there is no notion of a single point or locus of control – nothing corresponding to the program counter of a conventional sequential computer. The dataflow model describes computation in terms of locally controlled events; each event corresponds to the “firing” of an actor. An actor can be

* Corresponding author.

¹ Supported in part by DARPA under US AirForce Research Laboratory contract F33615-98-C-1319.

² Supported in part by the Ptolemy project, which is supported by the DARPA, the State of California MICRO program, and the following companies: Cadence, Hewlett Packard, Hitachi, Hughes Space and Communications, Motorola, NEC, and Philips.

a single instruction, or a sequence of instructions (per se, the dataflow model does not imply a limit on the size or complexity of actors). An actor fires when all the inputs it requires are available. In a dataflow execution, many actors may be ready to fire simultaneously (locally controlled by their operand availability), and thus these actors represent many asynchronous concurrent computation events.

Work on dataflow computer architecture emerged in the early 1970s with the use of dataflow program graphs to represent and exploit the parallelism in programs [5,28,30]. In a Dennis dataflow graph [30,31], operations are specified by actors that are enabled just when all actors that produce required data have completed their execution. The dependence relationships between pairs of actors are defined by the arcs of a graph, which maybe thought of as conveying results of an actor to successor actors, and by the firing rules, which specify exactly what data are required for an actor to fire. Decision and control actors may be included to represent conditional expressions and iterations. They alter the routing of tokens to effect data-driven control. Data structures may be constructed and accessed by appropriate dataflow actors.

Also in the early 1970s, a rather different style of dataflow models emerged [70]. Called Kahn process networks, this model replaces actors with sequential processes. These processes communicate by sending messages along channels that conceptually consist of unbounded FIFO queues. Kahn dataflow can be viewed as a generalization of Dennis dataflow [83]. While Dennis dataflow was originally applied to computer architecture design, Kahn dataflow was used by concurrency theorists for modeling concurrent software. Multithreaded architectures, with dataflow roots, use a style of dataflow that can be viewed as having elements of both.

In computer architecture, dataflow program graphs were originally used as a machine-level program representation. Two forms of dataflow architecture have become known: In a static architecture, the arc connecting one instruction to another can contain only a single result value (a token) from the source instruction. In this scheme there can be only one instance of a dataflow actor in execution at any time. In a dynamic dataflow architecture, tags are conceptually or actually associated with tokens so that tokens associated with different activations of an actor may be distinguished. This enables arcs to simultaneously carry multiple tokens, thereby exposing more data parallelism.

Because of the inherently parallel nature of dataflow execution, dataflow computers provide an efficient and elegant solution to the two fundamental problems of von Neumann computer: the memory latency and synchronization overhead, as described by Arvind and Iannucci [6]. The ability of the dataflow model to tolerate latency, by switching dynamically between ready computation threads, and to support low overhead distributed synchronization in hardware, has made it the candidate of choice for what has later been called “latency tolerant” architectures.

A history of the evolution of the dataflow computation model(s) and the variety of architecture and computational models that have been inspired by it is beyond the scope of this paper. The reader is referred to, among many others, the following texts [47,63,73].

As a computation model, the dataflow approach has had influence on many areas of computer science and engineering research. Examples include programming languages, processor design, multithreaded architectures, parallel compilation, high-level logic design, signal processing, distributed computing, and programming of systolic and reconfigurable processors. This paper reviews the recent developments in three of these areas: multithreaded execution, signal processing, and reconfigurable computing. The first two areas are relatively mature and have generated a substantial body of research results and commercial application. The last area is relatively young and is still in an emerging state. The first two sections review the state of the art of the dataflow computational model in multithreaded architectures and signal processing, respectively. The third section proposes the dataflow model as a programming model for reconfigurable computing.

2. Dataflow and multithreaded execution

Issues and challenges for multithreaded computer architecture suitable for use in parallel computers for general purpose computation are the subject of intensive debate. These issues and challenges depend heavily on the choice of an appropriate program execution model that will affect the programming model, the organization of the underlying system as well as the development support of complex applications onto the proposed architectures – compilers, runtime systems, and tools. The choice of program execution models may also have a profound impact on programming methodology [36,38,39].

In this section, we introduce basic concepts, characteristics and evolution of multithreaded computer architectures with dataflow origin. Our focus will be the review of issues and advances in program execution and architecture models upon which the design of these architectures is based.

In our terminology, a program execution model defines a basic low-level layer of programming abstraction of the underlying system architecture upon which the architecture model, programming model, compilation strategy, runtime system, and other software components are developed. It serves as an interface between the architecture and the software. Program execution model is in a broader (or higher) level than *instruction set architecture* (ISA) specification. An ISA usually provides a description of an instruction set for a particular machine that involves specific details of instructions such as instruction encoding and the organization of the machine registers set. In the context of this paper, a program execution model for multithreaded machines includes the following aspects: the thread model, the memory model and the synchronization.

The rest of this section is organized as follows: In Section 2.1, we briefly review the evolution of multithreaded architectures with dataflow origin. The subsequent sections will review the issues, challenges and advances in three areas: the thread models in Section 2.2, the memory models in Section 2.3 and the synchronization models in Section 2.4.

2.1. The evolution of the hybrid dataflow/von Neumann model and multithreading

The dataflow model and von Neumann serial control-flow model are generally viewed as two extremes of execution models on which a spectrum of architecture models can be based. The two models are in fact not orthogonal. Starting with the operational model of a pure dataflow graph, one can easily extend the model to support von Neumann style program execution. A region of actors within a dataflow graph can be grouped together as a thread to be executed sequentially under its own private program counter control, while the activation and synchronization of threads are data-driven. The new *hybrid* model is flexible in combining dataflow and control-flow evaluation, as well as in exposing parallelism at a desired level.

Such hybrid multithreaded architecture models have been proposed by a number of research groups with their origins in either static dataflow or dynamic dataflow. A number of articles have been published on multithreaded execution and architecture models with dataflow origin and can be found in a survey article by Dennis and Gao [33]. Principal projects and representative work before 1995 have been discussed in two monographs [47,63] as well as other survey articles [71]. Below we briefly mention a few research projects on each model.

Inspired by the static dataflow model, the *McGill dataflow architecture model* [43,45] has been proposed based on the *argument-fetching* principle [32]. The architecture departs from a direct implementation of dataflow graphs by having instructions fetch data from memory or registers instead of having instructions deposit operands (*tokens*) in “operand receivers” of successor instructions. The completion of an instruction will post an *event* (called a signal) to inform instructions that depend on the results of the instruction. This implements a modified model of dataflow computation called *dataflow signal graphs*. The architecture includes features to support efficient loop execution through *dataflow software pipelining* [42], and the support of threaded function activations. A good summary of the McGill dataflow architecture model (MDAM) and its extensions has been presented in Hum’s Ph.D. thesis [59].

Based on his experience with the MIT dynamic (tagged-token) dataflow architecture [4], Iannucci combined dataflow ideas with sequential thread execution to define a hybrid computation model described in his Ph.D. thesis [62]. The ideas later evolved into a multithreaded architecture project at IBM Yorktown Research Center as described elsewhere [63]. The architecture includes features such as a cache memory with synchronization controls, prioritized processor ready queues and features for efficient process migration to facilitate load balancing.

The P-RISC [99] is an interesting hybrid model exploring the possibilities of constructing a multithreaded architecture around an RISC processor. The Star-T project – a successor of the Monsoon project – has defined a multiprocessor architecture using an extension of an off-the-shelf processor architecture to support fine-grain communication and scheduling of user microthreads [10]. The architecture is intended to retain the latency-hiding feature of the Monsoon split-phase global memory operations. Later development of the Start-T project can be found in [3,25].

Several other interesting projects have been carried out elsewhere in the US [64–66,95,96,109] and in the world such as in Japan [75,100]. For example, the RWC-1 project [106] of the Real World Computing Partnership in Japan was initiated based on the experience of working with dataflow architectures and their multi-threaded extensions.

It is possible to define a thread model combining the advantages from both the static and dynamic dataflow models. An example is the EARTH model [61] which provides the full generality of threaded function invocation as in the dynamic dataflow model, while maintaining the simplicity of the static dataflow at the finer level of threads through software pipelining.

2.2. The thread models

Multithreaded execution models with dataflow origin provide support for fine-grain threads at two levels. For example, under the EARTH model, the first level of thread is called *threaded function invocation*: parallel function invocation forks a thread to execute the function in parallel. Note that the caller continues its own execution without waiting for the return of the forked threaded function. At a lower (finer) level, the body of a threaded function can be further partitioned (by a user or a compiler) into fibers [112]: a collection of operations that can be forked as a separate thread.

The support of such finer level of threads is a distinguishing characteristic of multithreaded architecture models with dataflow origin. Such a thread may be generated from a region of a dataflow graph derived from functional programming such as the so-called “strongly connected regions” in the EM-4 architecture and its successors [74,75], or the superactors in the Super-actor machine [59] based on the McGill dataflow architecture model (MDFFA [45]). With the recent advances in compiler technology, new algorithms and methods have been developed to perform thread partitioning for programs written in conventional programming paradigms [56,110,111]. A more in-depth discussion on such multithreaded compilation technology can be found in Tang’s Ph.D. thesis [105]. Such compilers may also require the use of a cost model for the underlying multithreaded machines. This is a non-trivial task due to the features of dynamic thread scheduling and latency tolerance – which have motivated the study to develop such models [98]. The evolution of a combination of compiler advances and architecture concepts has led to the proposal of the *superstrand* execution and architecture model [90–92].

Theobald in his Ph.D. thesis [112] has illustrated that such finer-level threads – also called *fibers* in his thesis – have played an important role for efficiently exploiting producer–consumer style fine-grain parallelism in scientific applications. He has demonstrated how to exploit a style of *dataflow software pipelining* between producers and consumers through fine-grain synchronization between fibers (discussed below in Section 2.4), and how such “pipes” can be established smoothly through threaded function invocations.

2.3. The memory models

Before we discuss the synchronization model, we need first to discuss the issues and advances in memory models. This is, perhaps, one of the most important areas where the evolution of such multithreaded execution and architecture models has departed from the dataflow models. The pure dataflow models do not use the conventional notion of an updatable memory. The programming models (e.g. functional or dataflow programming languages) are based on a semantic model where there are no side-effects due to execution of any operations. As can be seen from the following discussion, researchers in multithreaded execution models have advanced far beyond the pure dataflow models in developing memory models and cache memory schemes that take advantage of situations where single assignment rule can apply, while allowing normal load/store operations (with side-effects) when required without imposing unnecessary constraints for parallelism and efficiency [60].

Following the convention in [104], a memory model in our terminology contains two essential aspects: (1) the support of a shared space addressing and (2) the so-called *memory consistency model*. Advances have been made in addressing both issues.

Global addressing capability. It is reasonable to assume that part of the storage space of the system is globally addressable. Experience has shown that a shared addressing capability makes naming logically shared data much easier for the programmer because any thread can directly reference any data in the shared space and the naming of the model is similar to that on a uniprocessor. This should greatly enhance the programmability for certain applications (although it is still a subject of debate). For a multithreaded execution model (no matter with dataflow origin or not) globally addressable memory provides a seamless extension of the memory model viewed by threads assigned to the same processing node. For applications with irregular and unpredictable data needs, the lack of global naming capability can hurt the programmability and efficiency [108,116]. The ability to name memory objects globally will also facilitate the support of dynamic load balancing of threads. As a result, shared-address space is adopted by most multithreaded execution models with dataflow origin.

Consistency and replication. A memory consistency model represents a binding “contract” between software and hardware in a shared-memory multiprocessor system. It is important to provide a memory consistency model that is easy to understand and that also allows for an efficient hardware implementation. Issues arise as to how the replication of non-local data are managed. In a classical shared address space, since non-local data may be accessed through ordinary processor reads and writes, opportunities exist for the underlying system to replicate data in a way that is transparent from the users.

The design of most shared-memory multiprocessors has been based on a multiprocessor memory consistency model, usually one derived from Lamport’s concept of sequential consistency [78] which requires the execution of a parallel program to appear as some interleaving of the memory operations on a sequential machine. Substantial research has been performed on how to develop relaxed memory

consistency models based on the sequential consistency [50,51,85]. However, there have been observations that the performance gained by using relaxed models does not justify their complexity [57], with the implication that perhaps multiprocessors should only support simple consistency models.

Some multithreaded program execution models with dataflow origin are taking an alternative point of view that the function of a memory system is to support correct execution of programs by the computer system of which it is a part [48]. A computer system, with its memory subsystem, must be a correct implementation of the program execution model. While the formulation of a memory consistency model may be part of the process of building a multiprocessor computer system, the assumed sequential consistency “requirement” is not essential to the construction of a useful and efficient parallel processor, and may even be a distraction from the primary task of correctly implementing a desirable model of parallel computing. The challenge is to investigate and develop a sound memory model fully integrated into the underlying multithreaded execution framework and design an efficient cache management scheme for its implementation.

Most scientific applications are determinate, and such programs written under a suitable execution model will not incur any scenario where memory races may occur. Therefore, memory coherence is not an issue here. Under a multithreaded execution model with dataflow origin, the so-called *single assignment* style is extended and enriched in such a way that writes in the same threads can be considered as ordinary “updates” to memory. However, the dataflow style synchronization between threads ensures that producers and consumers of values are guaranteed to be properly ordered.

In particular, under the single assignment property of the I-structure memory [7] once a data element is defined, it will never be updated again. The copies of the data elements in the local cache will never be updated. Therefore, cache coherence should not be an issue in I-structure memory systems. It makes the design of I-structure cache much simpler without having to take care of the cache coherence problem. Such a memory model and its cache management have been proposed and described in [35,37]. Implementation of I-structure caches in multithreaded architectures with dataflow origin has also been pursued in a number of research groups [2,7,75,86,100].

For applications where non-determinacy is required – such as transaction processing applications – several schemes based on an extension of functional programming and I-structures have been reported. One scheme is based on the M-structures proposal [9]. Another scheme is to propose a program execution model where sound non-determinacy operators are introduced so the programs can be written in such a way that non-determinacy can be isolated, and memory operations are properly ordered, such that its effect of the ordering of memory system can be ensured to have expected behavior as anticipated by the users without resorting to traditional memory coherence requirement. A proposal along this line is based on the non-deterministic merge operator as described in [34].

Another proposal is to develop a memory consistency model, as an alternative to the SC-derived model that does not make the coherence assumption. One such model is the *Location consistency* [46,48,49]. Instead of assuming that all writes to the

same location are serialized according to some total order, the state of a memory location is modeled as a partially ordered multi-set (pomset) of write operations and synchronization operations. The partial order in the pomset naturally follows from the ordering constraints defined by a concurrent program. A cache management mechanism under the LC model has been proposed in [49] which eliminates the unproductive traffic for maintaining cache coherence.

The implementation of these proposals and performance evaluation based on shared memory programs are yet to be carried out so their advantages and trade-offs can be judged both qualitatively and quantitatively.

2.4. Synchronization models

This section will focus mainly on the synchronization among the finer-level threads (as described in Section 2.2) – a distinctive feature of multithreaded models with dataflow origin. The semantics of threaded function invocations themselves is well understood and has been discussed elsewhere [33].

Recall that when a program is divided into these finer-level threads, the data and control dependencies among the threads must be preserved. Under a program execution model with dataflow origin, such control and data dependencies are made explicit in the code. A synchronization event is posted from one thread to another to inform the recipient that a specific control or data dependence has been satisfied. For instance, the sending thread may have produced data required by the receiving thread according to some data dependence, and the producer thread must tell the consumer thread that the data are ready. If a thread depends on more than one data or control event, it has to make sure that all dependences have been satisfied before it becomes enabled. The underlying architecture (and runtime system) should provide atomic operations for sending data (possibly to a global memory location) and posting an event to guarantee that the data have been properly transferred before the receiving thread that depends on the data is enabled by the synchronization event. Some recent studies on the effectiveness of these operations have been reported in [112]. As a remark, these operations are very powerful: they can name any location in the shared address space for the transaction, they also name a thread which should be informed when the transaction is completed. Note that a designated thread may or may not become enabled after the transaction – because it may wait for more than one synchronization event. So these operations are related to both traditional shared-memory synchronization operations as well as message passing operations (including *active messages* [113]) but are different from both.

3. Dataflow in signal processing

Signal processing has its intellectual roots in circuit theory. Algorithms tend to be modeled as compositions of components that conceptually operate concurrently and communicate via signals that are functions of time. The components are typically

filters, which transform the input signals to construct output signals. Traditionally, the signals are continuous functions of the time continuum.

Contemporary signal processing is more likely to involve discrete-time signals than continuous-time signals. Here, instead of time continuum, a discrete clock globally regulates the computation. Signals have values at the discrete clock ticks, but not in between. Multirate systems involve multiple clocks, but with clear relationships among them, so that ticks in the multiple clocks can be unambiguously associated.

A fairly direct model of discrete-time systems is obtained using the synchronous/reactive principle [12,13]. In this model, a concurrent program executes in a sequence of discrete steps, which correspond to ticks of a global clock. At each clock tick, the value of each signal is obtained by solving for a fixed point of a system of equations. One of the possible outcomes of this fixed point solution is that a signal has no value at a particular clock tick. By this mechanism, multirate systems are easily modeled. Some synchronous/reactive languages, notably Lustre [67,68] and Signal [80], have a dataflow flavor, in that signals are viewed as streams of values, where each value is aligned with a clock tick. However, the concurrent semantics of these languages is very different from Dennis [30,36] or Kahn [70] styles of dataflow, which are distinctly asynchronous. Whereas synchronous languages are analogous to synchronous circuits in their treatment of concurrency, Dennis and Kahn-style dataflow are analogous to self-timed circuits.

The synchronous model globally orders the tokens according to a global clock. In Dennis and Kahn dataflow, by contrast, signals are streams of tokens, where the relative ordering of tokens within a stream matters, but there is not necessarily any ordering of tokens across streams. Partial ordering constraints on tokens in distinct streams are imposed by the data precedences in the actors or processes.

In Dennis dataflow, actors are enabled by the presence of tokens on the input streams, and once enabled, can fire to produce tokens on the output streams. This enabling/firing sequence implies a partial order in which output tokens are required to occur after the input tokens that trigger the firing.

In Kahn's model, sequential processes communicate via FIFO queues that are conceptually unbounded. A process can write to a queue at any time, but if it reads from a queue, then it will block if there are no data. This blocking is a kind of dual of the Dennis notion of firing, where an actor can be fired once input data become available. In Kahn's model, a process replaces a Dennis actor, and it is active unless there are no input data. A sequence of firings of a Dennis actor can be viewed as a Kahn process [83], so we will henceforth refer to these two models simply as "dataflow". For their detailed semantic connection, see [84].

The synchronous/reactive model provides an abstraction of discrete-time systems where the metric properties of time are eliminated. That is, there is no measured or even specified time interval between clock ticks. Rather, there is a sequence of clock ticks. But one key property of time is retained: its global ordering of events. Dataflow abstracts things still further by eliminating this global ordering.

Streams of tokens provide a suitable abstraction for signals. Dataflow actors or Kahn processes operating on these streams provide a suitable abstraction for the

components. The notion of time is lost, but the data precedences are preserved. Thus, the abstraction is more highly concurrent even than the circuit-theory roots of signals processing, since the tight coordination implied by a global ordering based on time is no longer required. Instead, only the data precedences need to be respected. It is arguable that this reduces stream-based processing to its computational essence.

3.1. Industrial practice

The signal processing community has embraced dataflow modeling. A major reason for this is the widespread use of block diagrams for documenting and explaining signal processing algorithms. Dataflow semantics work well with a block diagram syntax. Moreover, the use of block diagrams evokes a circuit metaphor, thus making a connection with the historical roots of signal processing. But it also proves to be a convenient and lucid way to specify reasonably complex algorithms. By contrast, attempts to use dataflow principles directly in the computer architecture of embedded signal processors have largely failed.

Commercial software environments that combine dataflow semantics with a block diagram syntax include SPW from Cadence, Cossap from Synopsys, and several programs from smaller players. These software environments are used for algorithm-level modeling of systems, for specification, and as a starting point for synthesis of customized hardware designs and embedded software.

Some closely related block diagram environments have semantics that more closely resemble synchronous/reactive languages, such as Simulink from The MathWorks. These are more explicit about time. Simulink in particular has its roots in continuous-time modeling of control systems, and therefore has a very physical model of time as an integral part of its semantics.

Block diagram languages are by no means the only tool used in signal processing. A very popular alternative is matrix-oriented imperative languages such as Matlab from The MathWorks. These languages speak to the mathematical roots of signal processing more than to its roots in circuit theory. Indeed, these languages do not directly model signals. Instead, signals must be decomposed into finite and classical data structures, particularly vectors and arrays.

Since many of the more sophisticated algorithms in signal processing are based on matrix formulations, Matlab and related languages prove very convenient for describing them. These algorithms, however, most often get embedded in systems where stream-based processing is much more natural. A fairly common compromise, therefore, is to use coarse-grain dataflow, where actors can represent large scale computations that are specifically defined using a matrix-oriented language. At least one environment specifically permits the functionality of dataflow actors to be given using Matlab [20].

3.2. Decidable dataflow

Many signal processing systems involve repeated (infinite) execution of a well-defined finite computation on an infinite stream of data. Implementations have real-

time constraints, and often take the form of embedded software (such as assembly code for programmable DSPs [79]). This raises a number of interesting issues. In particular, it is important that the schedule of actor firings be predictable in order to ensure that real time constraints are met. It is also critical that a program never deadlocks. Because of the embedded system context, it is also important that the total memory devoted to storing unprocessed tokens be bounded.

In dataflow, deadlock occurs when all actors are starved. Deadlock is equivalent to halting. For general dataflow models, it is undecidable whether a program will deadlock. It is also undecidable whether a program has an infinite execution that consumes bounded memory for storing pending tokens [18,19]. One approach, therefore, is to use a subset of dataflow where these questions are decidable.

The so-called synchronous dataflow [81] and its extensions [16,44], which is not synchronous in the same sense as synchronous/reactive languages, is one such decidable model. In SDF, an actor is characterized by the number of tokens that it consumes and produces on its inputs and outputs on each firing. The inputs and outputs of each actor are labeled with an integer constant. For an input, this is the number of tokens that are required on that input stream in order to fire the actor. For the outputs, it is the number of tokens that are produced by a firing.

In a block-diagram syntax, an arc connecting the two actors represents the flow of tokens between them. If we assume that the topology of such flows is fixed, and that the number of tokens produced and consumed by each actor is fixed, then both deadlock and bounded memory are decidable. To get an idea of how this is possible, we review the notion of balance equations [81].

Suppose an actor A produces N tokens on an output that is connected to actor B. Suppose further that actor B requires M tokens on its input to fire. Suppose that actor A fires f_A times, and actor B fires f_B times. Then the balance principle requires that

$$f_A N = f_B M.$$

One such equation can be written for each arc in an interconnection of actors.

The balance principle is trivially satisfied if either $f_A = f_B = 0$ or $f_A = f_B = \infty$. A more interesting situation arises, however, if there is a bounded positive solution for f_i for each actor i such that all balance equations are satisfied. In this case, there may be a finite but non-zero set of firings that achieves balance.

It has been shown that for a connected dataflow graph, if the balance equations have a non-trivial solution, then they have a unique smallest positive integer solution [81]. If they have no solution, then there is no bounded memory infinite execution of the dataflow graph.

If the balance equations have a non-trivial solution, then there is a simple finite algorithm that will determine whether the graph deadlocks. Thus, both the bounded memory question and the deadlock question are decidable for SDF. Hence, SDF graphs can reliably be implemented in embedded real-time systems.

Although the decidability of the SDF model is convenient for certain applications, it comes at a very high cost in expressiveness. In particular, the fact that the token production/consumption patterns are fixed means that an application cannot use the

flow of tokens to effect control. That is, conditional variations in the flow are not allowed. Even signal processing applications involve a certain amount of conditional computation, so this restriction becomes very limiting.

Consider for example an actor which, under the control of a boolean-valued input stream, routes tokens from another input stream to one of two outputs. Such an actor is often called a “switch”. The partner of the switch is the “select”, which uses a boolean input stream to guide its decisions about which of two other input streams to read for the next token. After reading an input token from either of these two inputs, it outputs that token on its single output. These two actors thus serve the same function as railway switches, splitting or merging streams of data under the control of boolean signals. However, the introduction of these two actors alone to the SDF formalism makes both deadlock and bounded memory undecidable.

One solution is to broaden the dataflow model that is used by accepting the undecidability of the key propositions but nonetheless attempt to decide. Undecidability merely states that no algorithm can decide in finite time for all programs. It does not prevent us from constructing a formalism where a compiler can decide for most programs.

This is the approach taken by Buck [19], who based his approach on the token flow model of Lee [82]. This model generalizes the balance equations to permit the number of tokens produced and consumed to vary dynamically. To do this, it attaches a symbolic value rather than a numerical value to each port of each actor. The balance equations therefore include symbolic coefficients, and have to be solved symbolically. Frequently, it is possible to solve them. The resulting solution can be used to construct a schedule of actor firings that is provably free from deadlocks and that requires only bounded memory for token storage.

Buck’s model came to be called boolean dataflow, although the basic concept has been extended to non-boolean control signals [21]. BDF, however, has two problems. First, because the key underlying implementation questions remain undecidable, a programming environment that depends on being able to decide is somewhat fragile. Small and seemingly innocuous changes to a program may cause the scheduling algorithms to fail. Also, experience constructing programs using this model indicates that it is not often the most intuitive way to specify a computation. Programs often resemble puzzles, where considerable cleverness is required to construct them and to understand what they do. Using the routing of tokens for control flow may not be such a good idea.

3.3. *Reactive systems*

Signal processing systems almost never do only signal processing. Consider for example a cellular telephone. It includes intensive numerical signal processing in both the speech coder and radio modem. It may also include advanced features such as speech recognition for hands-free dialing. These signal processing components are each quite sophisticated, and invariably involve algorithms that stress the limitations of SDF. For example, equalization of a fading radio channel may involve the use of distinct algorithms during the establishment of a connection vs. steady state. Also,

power conservation dictates the use of simpler algorithms when the channel is benign, suggesting mode changes that would be driven by channel estimators. SDF by itself does not match these requirements well. BDF is a possibility, but it is not a complete solution if one considers the entire design of the phone.

Signal processing applications fall into a broader class of so-called reactive systems [89], which are those that interact with their environment at the speed of the environment. Much of what characterizes many reactive systems is randomly structured control logic that bears little resemblance to the regular algorithmic structure of core signal processing algorithms. A cellular phone, for example, includes not only sophisticated signal processing, but also a substantial amount of embedded control logic for call processing and multiple-access protocols. BDF no longer matches well, since even relatively simple control logic becomes very difficult to understand.

One possible solution is to embrace a heterogeneous approach, where modeling techniques can be combined systematically to construct more complex systems. Decidable dataflow models would be used for the hard-real-time signal processing in the cellular phone, but other formalisms would be used for the control logic. Some researchers advocate melding state machines with dataflow [52], much the way StateCharts melds a synchronous/reactive model with state machines [69]. State machines are good at representing intricate control logic. Moreover, they often yield to formal analysis, and with careful language construction, remain decidable.

4. Dataflow as a programming model for reconfigurable computing

The dataflow graph program representation model has been proposed as a platform for the programming of reconfigurable computing machines (RCMs). Reconfigurable computing consists of the dynamic mapping, and re-mapping, of computations directly onto hardware. Even though the reconfigurable computing model is not new [41], it has received renewed attention in recent years due to the development and rapid evolution of FPGA devices in particular [14,22] and of VLSI technology in general.

The design of computing systems is a trade-off between generality and efficiency. One of the great attractions of the reconfigurable computing model is its ability to increase efficiency with a smaller sacrifice in generality than vector processors for example. As a computation model, reconfigurable computing avoids one of the major limitations of the von Neumann stored-program model, namely the fetching of instructions from a storage, it can therefore achieve a much higher operation density than general-purpose processors [29,11,12,54,119,120].

Reconfigurable computing is a parallel execution model that combines aspects of the SIMD, systolic arrays and MIMD models. The granularity of the computation in this model can vary, within a same program, from bit-level manipulations to loop-level parallelism. Because of the wide variation in granularity and the flexibility of the execution model, in particular on machines that support partial dynamic reconfiguration, the programmability of RCMs is a major obstacle to their wider

acceptance. The development of an accessible yet efficient programming paradigm for reconfigurable computing is an important research challenge.

4.1. Reconfigurable computing machines

In its most general form, a reconfigurable computing machine (RCM) consists of a processor and an array of programmable blocks (RC array). One of the main distinguishing features of RCMs is the granularity of the programmable blocks. At one extreme are the FPGA-based machines (FRCMs) where the programmable block consists of a logical collection of gates. Their programming model is that of a gate-level net-list [87,88]. On the other extreme are the coarse-grain reconfigurable machines (CRCMs) where the logic block consists of an ALU and possibly a register file. In either of these categories, the architecture of the RCM can be a specialized processor (such as the Ganglion, built for neural network computations [26]), a coprocessor (such as the PRISM [8] and the Thinking machine CM-2X [27]) or an attached processor (such as the DEC PAM and PeRLe-0 and PeRLe-1 [14,15,114] and the SRC Splash 1 and 2 [22,53]).

Fine-grained RCMs. FRCMs rely on FPGAs as the reconfigurable array. The computation to be mapped onto hardware is written in a hardware description language (HDL) such as VHDL or Verilog. The HDL program is then compiled using a synthesis tool and the resulting net list is placed and routed for the specific FPGA device used. The resulting configuration codes are then down loaded to a board that comprises one or more FPGA devices. The programming model for these machines is, therefore, essentially that of hardware design. It is a lengthy process and requires expert knowledge of HDLs as well as the particular FPGA device family being targeted. Its advantages are its low cost and its very high flexibility: the configurable circuit is designed to exactly fit the computation. It was shown that the efficiency of computations on FRCMs, measured in bit operations per λ^2 second, is one to two orders of magnitude that of conventional computers [29]. This model, therefore, tends towards the efficiency levels of ASICs while preserving a much higher flexibility. Examples of this model are the DEC PAM and PeRLe machines, the SRC Splash 1 and 2 and the Annapolis Microsystems Wildforce board [1]. FRCMs are generally designed as attached processors communicating with a host processor on an I/O bus.

Coarse-grained RCMs. At the other extreme are the RCMs where the programmable blocks consist of complex ALUs with local registers and possible support for floating-point operations. In these machines the programming model is similar to a conventional processor or multiprocessor. Examples of this model include the UC Irvine Morphosys chip [77], the UC Berkeley Garp [55], the MIT RAW machine [115], the CMU Pipe-Rench [23,107] and the University of Washington RaPiD [40]. The advantage of this model is its simpler programming paradigm. It also benefits better from traditional compiler optimizations, both sequential and parallel. Its hardware structure, however, is more rigid and hence less flexible than that of its FRCMs. CRMs are generally designed as coprocessors or integrated in the processor design itself [23,24,77]. These two RC architecture models are by no means exclusive.

As a matter of fact, the model proposed by Gerald Estrin in 1963 [41] has been called the “Fixed plus Variable” model: It combines, in one processor, a traditional CPU and a fine-grained reconfigurable engine. In a way, most reconfigurable architectures follow this same model of fixed plus variable. However, the relative sizes and the degree of integration of the fixed and variable parts vary. The CMU PipeRench, the UC Berkeley Garp and the UC Irvine M1 processors consist of a CPU and the reconfigurable component on the same processor chip.

4.2. The programmability of RCMs

The poor programmability of RCMs is probably the most serious obstacle to their wider adoption. It derives from two distinct and orthogonal problems: (1) The complexity of the model and (2) The use of HDL (VHDL or Verilog) as a programming language in RCMs.

The complexity of the model. Reconfigurable computing is a powerful model that can exploit program parallelism at most levels and in an *integrated* way. The parallelism in this model ranges from bit-level manipulations, to instruction, to loop and task levels (such as producer–consumer relationships). While in a von Neumann processor, these levels of parallelism are exploited in a somehow hierarchical manner (reflected for example in the storage hierarchy: registers, cache, memory, virtual memory), in a reconfigurable hardware these levels of parallelism are integrated in a flat “address” space (the address space becomes physical space). Dynamic reconfiguration adds another dimension to the programming complexity of this model. Dynamic reconfiguration is the ability to “re-program” the RCM, all or just part of it, dynamically and based on the outcome of a computation. In summary, this computational model combines the power, and complexity, of the SIMD, systolic arrays and the MIMD models. The challenge is therefore to migrate, at least in part, some of that complexity from the user to the compiler.

The use of HDLs. At the root of the resurgence of the reconfigurable computing model is the emergence of high-density and high-speed FPGAs. It is only natural that such machines have been programmed using HDLs. However, both Verilog and VHDL are languages designed, initially, for system level simulation. They are both based on a discrete-event simulation paradigm and are not easily accessible to conventionally applications programmers. The challenge is to interface between traditional high-level languages and HDLs [118].

4.3. Dataflow graphs and reconfigurable computing

A dataflow graph represents the essence of a computation. The two major characteristics of a dataflow graph program representation are:

- Data are represented by values (arcs) and not by memory location.
- Operations are functional: the output depends only on the inputs.

As a result, the advantages of this model for reconfigurable computing are:

- The ordering of operations is determined only by true data dependencies. It is not constrained by artificial dependencies introduced by the stored program model, such as memory aliases.
- It displays parallelism in a computation at all levels. Fine-grained parallelism, at the operation level, can be exploited within a processor. Coarser, loop or function level parallelism can be exploited across processors. In contrast: systolic arrays exploit only fine grain parallelism, SMP only coarse grain.
- At the semantic level, combinational logic circuits are naturally functional (their output depends only on their inputs) and are therefore easily synthesizable from dataflow graphs.
- With certain restrictions on actors, dataflow graphs do not have any implicit state. They allow easy and correct synthesis of sequential circuits where the storage elements are determined by the synthesis tool rather than by a compiler.

Various forms of dataflow graphs have been proposed and used for the programming of systolic and wavefront arrays [72,94,117] as well as for the synthesis, hardware/software co-design, of application specific hardware [17,58,93].

The Cameron Project is a research effort that aims at leveraging the potentials of a dataflow graph program representation to provide high-level programming for RCMs [97]. The application domain that is targeted by the Cameron Project is image processing, however, the approach is general enough and could be extended to other areas. Cameron relies on Khoros as a graphical user interface. Khoros is a widely popular program development environment designed for image processing (IP) applications [76,103].

Program modules that run on the RCM are written in a single-assignment subset of C called SA-C. SA-C programs are tightly integrated into Khoros. SA-C has been extended to support image processing constructs as well as compilation for a hardware implementation. The most salient features of SA-C are its support for true n -dimensional arrays, window-based operations on these, and variable precision fixed-point arithmetic. The first feature simplifies the expression of image processing algorithms while the second improves the efficiency of the mapping to hardware. Other C-like hardware definition languages have been developed (e.g., Handel-C [101,118]) with a large degree of success. However, they still require the manual translation of algorithms into logic circuit structures. The objective of Cameron is to automate that process while still providing the programmer with some control over the mapping to hardware. SA-C programs are compiled and optimized into a dataflow graph. This format is then compiled to VHDL for FRCMs or to native machine code for CRCMs.

The project described in [102] relies on the synchronous dataflow model, as implemented in Ptolemy, to develop mapping techniques for RCMs tailored to signal processing. It extends Ptolemy by developing a new reconfigurable computing domain that separates the interface specification from implementation for each signal processing functional block.

Reconfigurable computing is a very powerful paradigm with a great potential. It is, however, still in its infancy. The implementation of this model is enabled by the rapid advances in VLSI technology. Its greatest challenge lies in its programmability.

Unless that challenge is effectively and efficiently addressed, reconfigurable computing runs the risk of remaining a potential paradigm.

5. Conclusion

As a model of computation, dataflow has a long history. It has been used as a principle for computer architecture, as a model of concurrency in software, as a high-level design model for hardware. It has demonstrated its flexibility and efficiency at representing computation by the wide variety of areas in which it has been used.

The astoundingly rapid evolution of computing machines and their phenomenal penetration of modern society in the last half century has been attributed, to a large degree, to the simplicity, constance and uniformity of the stored program model. It has provided the stable conceptual framework within which processors, operating systems, algorithms, languages and compilers were designed and implemented. However, as applications evolve and exploit all the computing power of modern machines, the demand for higher performance machines is rapidly increasing. It is in the implementations of high performance computing systems that the stored program model shows its limitations. These limitations are directly and elegantly addressed by the dataflow model. We believe that the dataflow model is not only still relevant in many computing research areas, but also that many emerging and future areas will benefit from this model.

Acknowledgements

The authors acknowledge the support from DARPA, NSA, NSF, and NASA. The last author thanks members of his Computer Architecture and Parallel Systems Laboratory (CAPSL) group, Department of Electrical and Computer Engineering, University of Delaware, for their valuable input in several discussions. And finally, thanks to Dr. Ruppa Thulasiram for his careful reading of the manuscript and comments.

References

- [1] Annapolis Micro Systems, Wildforce Board Web Page, <http://www.annapmicro.com>.
- [2] B.S. Ang, Arvind, D. Chiou, StarT the next generation: integrating global caches and dataflow architecture, in: G.R. Gao, L. Bic, J.-L. Gaudiot (Eds.), *Advanced Topics in Dataflow Computing and Multithreading*, IEEE Computer Society Press, Silver Spring, MD, 1995, pp. 19–54.
- [3] B.S. Ang, D. Chiou, D. Rosenband, M. Ehrlich, L. Rudolph, Arvind, Start-voyager: a flexible platform for exploring scalable smp issues, CSG Memo 415, Computation Structures Group, MIT Lab. for Comput. Sci., December 1998.
- [4] Arvind, R.S. Nikhil, Executing a program on the MIT tagged-token dataflow architecture, *IEEE Trans. Comput.* 39 (3) (1990) 300–318.
- [5] Arvind, K.P. Gostelow, The U-interpreter, *Computer* 15 (2) (1982) 42–49.

- [6] Arvind, R. Iannucci, A critique of multiprocessing von Neumann style, in: International Symposium on Computer Architecture, Stockholm, Sweden, 1983.
- [7] Arvind, R.S. Nikhil, K.K. Pingali, I-structures: data structures for parallel computing, *ACM Trans. on Program. Languages Syst.* 11 (4) (1989) 598–632.
- [8] P.M. Athanas, H.F. Silverman, Processor reconfiguration through instruction set metamorphosis, *Computer* 26 (1993) 11–18.
- [9] P.S. Barth, R.S. Nikhil, Arvind, M-structures: extending a parallel, non-strict, functional language with state, CSG Memo 327, Computation Structures Group, MIT Lab. for Comput. Sci., March 1991.
- [10] M.J. Beckerle, Overview of the START(*T) multithreaded computer, in: Digest of Papers, COMPCON Spring '93, San Francisco, CA, February 1993, pp. 148–156.
- [11] A. Benedetti, P. Perona, Feature detection on a reconfigurable computer, in: Conference on Computer Vision and Pattern Recognition, Santa Barbara, CA, 1998.
- [12] A. Benveniste, G. Berry, The synchronous approach to reactive and real-time systems, in: *Proceedings of the IEEE*, vol. 79 (9), 1991, pp. 1270–1282.
- [13] A. Benveniste, P. Caspi, P. Le Guernic, N. Halbwachs, Data-flow synchronous languages, in: J.W. de Bakker, W.-P. de Roever, G. Rozenberg (Eds.), *A Decade of Concurrency – Reflections and Perspectives*, Lecture Notes in Computer Science, vol. 803, Springer, Berlin, 1994, pp. 1–45.
- [14] P. Bertin, D. Roncin, J. Vuillemin, Programmable active memories: a performance assesment, in: G. Borriello, C. Ebeling (Eds.), *Research on Integrated Systems*, MIT Press, Cambridge, MA, 1993, pp. 88–102.
- [15] P. Bertin, H. Touati, PAM programming environments: practice and experience, in: *Proceedings of the IEEE Workshop on FPGAs for Custom Computing Machines*, IEEE Computer Soc. Press, Los Alamitos, CA, 1994, pp. 133–139.
- [16] G. Bilsen, M. Engels, R. Lauwereins, J.A. Peperstraete, Static scheduling of multi-rate and cyclostatic DSP applications, in: *Proceedings of the 1994 Workshop on VLSI Signal Processing*, IEEE Press, New York, 1994.
- [17] M. Bolotski, A. DeHon, T.F. Knight, Unifying FPGAs and SIMD arrays, in: *Proceedings of the FPGA Workshop*, 1994.
- [18] J.T. Buck, E.A. Lee, Scheduling dynamic dataflow graphs with bounded memory using the token flow model, in: *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. I, Minneapolis, MN, April 1993, pp. 429–432.
- [18] J.T. Buck, Scheduling dynamic dataflow draphs with bounded memory using the token flow model, Technical Report UCB/ERL 93/69, Ph.D. Dissertation, Department of EECS, University of California, Berkeley, CA 94720, 1993.
- [20] J.T. Buck, S. Ha, E.A. Lee, D.G. Messerschmitt, Ptolemy: a framework for simulating and prototyping heterogeneous systems (special issue on Simulation Software Development) *Int. J. Comput. Simulation*, 4 (1994) 155–182 (<http://ptolemy.eecs.berkeley.edu/papers/JEurSim>).
- [21] J.T. Buck, Static scheduling and code generation from dynamic dataflow graphs with integer-valued control systems, Invited Paper, in: *Proceedings of the IEEE Asilomar Conference on Signals, Systems, and Computers*, 31 October–2 November, Pacific Grove, CA, 1994.
- [22] D.A. Buell, J.M. Arnold, W.J. Kleinfelder, *Splash 2: FPGAs in a Custom Computing Machine*, IEEE CS Press, Silver Spring, MD, 1996.
- [23] S. Cadambi, J. Weener, S.C. Goldstein, H. Schmit, D.E. Thomas, Managing pipeline-reconfigurable FPGAs, in: *Sixth International Symposium on Field Programmable Gate Arrays*, 1998.
- [24] R. Carley, S.C. Goldstein, T. Mukherjee, R. Rutenbar, H. Schmit, D. Thomas, PipeRench Web Page, Carnegie-Mellon University, <http://www.ece.cmu.edu/research/piperench/>.
- [25] D. Chiou, B.S. Ang, Arvind, M.J. Beckerle, A. Boughton, R. Greiner, J.E. Hicks, J.C. Hoe, StarT-NG: delivering seamless parallel computing, CSG Memo 371, Computation Structures Group, MIT Lab. for Comput. Sci., February 1995.
- [26] C.E. Cox, W. Ekkehard Blanz, Ganglion – a fast hardware implementation of a connectionist classifier, in: *Proceedings of the IEEE Custom Integrated Circuits Conference*, IEEE Press, New York, 1991, pp. 6.5.1–6.5.4.

- [27] S.A. Cuccaro, C.F. Reese, The CM-2X: a hybrid CM-2/Xilinx prototype, in: *Proceedings of the IEEE Workshop on FPGAs for Custom Computing Machines*, IEEE Computer Soc. Press, Los Alamitos, CA, 1993, pp. 121–131.
- [28] A.L. Davis, R.M. Keller, Data flow program graphs, *Computer* 15 (2) (1982) 26–41.
- [29] A. DeHon, Comparing computing machines, in: *Configurable Computing: Technology and Applications*, Proceedings of SPIE 3526, November 1998, p. 124.
- [30] J.B. Dennis, First version of a data-flow procedure language, in: *Proceedings of the Colloque sur la Programmation, Lecture Notes in Computer Science*, vol. 19, Paris, France, April 9–11, Springer, Berlin, 1974, pp. 362–376.
- [31] J.B. Dennis, Data flow ideas for supercomputers, in: *Digest of Papers, COMPCON Spring '84*, San Francisco, CA, February–March 1984, pp. 15–19.
- [32] J.B. Dennis, G.R. Gao, An efficient pipelined dataflow processor architecture, in: *Proceedings of the Supercomputing '88*, Orlando, FL, November 1988, pp. 368–373.
- [33] J.B. Dennis, G.R. Gao, Multithreaded architectures: principles, projects, and issues, ACAPS Technical Memo 29, School of Computer Science, McGill University, Montréal, Qué., February 1994, in: <ftp://ftp-acaps.cs.mcgill.ca/pub/doc/memos>.
- [34] J.B. Dennis, G.R. Gao, Multithreaded architectures: principles, projects, and issues, in: R.A. Iannucci, G.R. Gao, R.H. Halstead Jr., B. Smith (Eds.), *Multithreaded Computer Architecture: A Summary of the State of the Art*, chapter 1, Kluwer Academic Publishers, Norwell, Massachusetts, 1994.
- [35] J.B. Dennis, G.R. Gao, On memory models and cache management for shared-memory multiprocessors, ACAPS Technical Memo 90, School of Computer Science, McGill University, Montréal, Qué., December 1994, in: <ftp://ftp-acaps.cs.mcgill.ca/pub/doc/memos>.
- [36] J.B. Dennis, Streams data types for signal processing, in: J-L. Gaudiot, L. Bic (Eds.), *Advanced Topics in Dataflow Computing and Multithreading*, Prentice-Hall, Englewood Cliffs, NJ, 1995.
- [37] J.B. Dennis, G.R. Gao, On memory models and cache management for shared-memory multiprocessors, in: *Parallel and Distributed Processing*, IEEE Computer Soc., 1995.
- [38] J.B. Dennis, A parallel program execution model supporting modular software construction, in: *Proceedings of the Massively Parallel Programming Models (MPPM-97)*, IEEE Computer Soc., 1997, pp. 50–60.
- [39] J.B. Dennis, General parallel computation can be performed with a cycle-free heap, in: *Proceedings of the International Conference on Parallel Architectures and Compiler Technology*, Paris, France, October 1996, pp. 96–103.
- [40] C. Ebeling, D.C. Cronquist, P. Franklin, RaPiD – reconfigurable pipelined datapath, in *Proceedings of the Field Programmable Logic*, 1996.
- [41] G. Estrin, Parallel processing in a restructurable computer system, *IEEE Trans. Elect. Comput.* (1963).
- [42] G.R. Gao, *A Code Mapping Scheme for Dataflow Software Pipelining*, Kluwer Academic Publishers, Boston, Massachusetts, 1990.
- [43] G.R. Gao, H.H.J. Hum, J.-M. Monti, Towards an efficient hybrid dataflow architecture model, in: *Proceedings of the PARLE '91*, vol. I, Lecture Notes in Computer Science, vol. 505, Eindhoven, The Netherlands, June 1991, Springer, Berlin, pp. 355–371.
- [44] G.R. Gao, R. Govindarajan, P. Panangaden, Well-behaved dataflow for DSP computation, in: *Proceedings of the ICASSP-92*, San Francisco, March 1992.
- [45] G.R. Gao, An efficient hybrid dataflow architecture model, *J. Parallel Distrib. Comput.* 19 (4) (1993) 293–307.
- [46] G.R. Gao, V. Sarkar, Location consistency: stepping beyond the barriers of memory coherence and serializability, ACAPS Technical Memo 78, School of Computer Science, McGill University, Montréal, Qué., December 1994, in: <ftp://ftp-acaps.cs.mcgill.ca/pub/doc/memos>.
- [47] G.R. Gao, L. Bic, J-L. Gaudiot (Eds.), *Advanced Topics in Dataflow Computing and Multithreading*, IEEE Computer Soc. Press, New York, 1995, book contains papers presented at the Second International Workshop on Dataflow Computers, Hamilton Island, Australia, May 1992.

- [48] G.R. Gao, V. Sarkar, On the importance of an end-to-end view of memory consistency in future computer systems, in: *Proceedings of the International Symposium on High Performance Computing*, Fukuoka, Japan, 1997, pp. 30–41.
- [49] G.R. Gao, V. Sarkar, Location consistency – a new memory model and cache consistency protocol, CAPSL Technical Memo 16, Department of Elec. and Computer Engineering, University of Delaware, Newark, Delaware, February 1998, in: <ftp://ftp.capsl.udel.edu/pub/doc/memos>.
- [50] K. Gharachorloo, A. Gupta, J. Hennessy, Revision to ‘memory consistency and event ordering in scalable shared-memory multiprocessors’, Technical Report No. CSL-TR-93-568, Computer Systems Lab., Stanford University, Stanford, CA, April 1993.
- [51] K. Gharachorloo, D. Lenoski, J. Laudon, P. Gibbons, A. Gupta, J. Hennessy, Memory consistency and event ordering in scalable shared-memory multiprocessors, in: *Proceedings of the ISCA-17*, Seattle, Washington, May 1990, pp. 15–26.
- [52] A. Girault, B. Lee, E.A. Lee, Hierarchical finite state machines with multiple concurrency models, *IEEE Trans. on CAD*, 1999, 18 (6) (1999). (Revised from Memorandum UCB/ERL M97/57, Electronics Research Laboratory, University of California, Berkeley, CA 94720, August 1997) (<http://ptolemy.eecs.berkeley.edu/papers/98/starcharts>).
- [53] M. Gokhale, W. Holmes, A. Kopsler, S. Lucas, R. Minnich, D. Sweely, D. Lopresti, Building and using a highly parallel programmable logic array, *IEEE Comput.* 24 (1991) 81–89.
- [54] R. Hartenstein et al., A reconfigurable machine for applications in image and video compression, in: *Conference on Compression Technologies and Standards for Image and Video Compression*, Amsterdam, Netherlands, 1995.
- [55] J.R. Hauser, J. Wawrzynek, Garp: a MIPS processor with a reconfigurable co-processor, in: *Proceedings of the IEEE Symposium on Field-Programmable Custom Computing Machines*, April 1997.
- [56] L.J. Hendren, G.R. Gao, X. Tang, Y. Zhu, X. Xue, H. Cai, P. Ouellet, Compiling C for the EARTH multithreaded architecture, ACAPS Technical Memo 101, School of Computer Science, McGill University, Montréal, Qué., March 1996, in: <ftp://ftp-acaps.cs.mcgill.ca/pub/doc/memos>.
- [57] M.D. Hill, Multiprocessors should support simple memory-consistency models, *Computer* 31 (1998) 28–34.
- [58] J. Horstmannshoff, T. Grötke, H. Meyr, Mapping multirate dataflow to complex RT level hardware models, in: *11th International Conference on Application-specific Systems, Architectures and Processors*, 1997, (ASAP’97).
- [59] H.H.-J. Hum, The super-actor machine: a hybrid dataflow/von Neumann architecture, Ph.D. Thesis, McGill University, Montréal, Qué., May 1992.
- [60] H.H.J. Hum, G.R. Gao, A high-speed memory organization for hybrid dataflow/von Neumann computing, *Future Generation Comput. Syst.* 8 (4) (1992) 287–301.
- [61] H.H.J. Hum, O. Maquelin, K.B. Theobald, X. Tian, G.R. Gao, L.J. Hendren, A study of the earthman multithreaded system, *Int. J. Parallel Programming* 24 (4) (1996) 319–347.
- [62] R.A. Iannucci, Toward a dataflow/von Neumann hybrid architecture, in: *Proceedings of the ISCA-15*, Honolulu, Haw., May–June 1988, pp. 131–140.
- [63] R.A. Iannucci, G.R. Gao, R.H. Halstead Jr., B. Smith (Eds.), *Multithreaded Computer Architecture: A Summary of the State of the Art*, Kluwer Academic Publishers, Norwell, Massachusetts, 1994, book contains papers presented at the Workshop on Multithreaded Computers, Albuquerque, New Mexico, November 1991.
- [64] S. Jenks, J.-L. Gaudiot, Exploiting locality and tolerating remote memory access latency using thread migration, *Int. J. Parallel Programming*, 1996, in press.
- [65] S. Jenks, J.-L. Gaudiot, Nomadic threads – a migrating multithreaded approach to remote memory accesses in multiprocessors, in: *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques (PACT ’96)*, Boston, Massachusetts, October 1996.
- [66] S. Jenks, J.-L. Gaudiot, Nomadic threads: A migrating multithreaded approach to remote memory accesses in multiprocessors, in: *Proceedings of the 1996 Conference on Parallel Architectures and Compilation Techniques*, Boston, Massachusetts, October 1996.
- [67] N. Halbwachs, *Synchronous Programming of Reactive Systems*, Kluwer Academic Publishers, Dordrecht, 1993.

- [68] N. Halbwachs, P. Caspi, P. Raymond, D. Pilaud, The synchronous data flow programming language LUSTRE, in: *Proceedings of the IEEE*, vol. 79 (9), pp. 1305–1319.
- [69] D. Harel, Statecharts: a visual formalism for complex systems, *Sci. Comput. Programming* 8 (1987) 31–274.
- [70] G. Kahn, The semantics of a simple language for parallel programming, in: *Proceedings of the IFIP Congress 74*, North-Holland, Amsterdam, 1974.
- [71] C. Kim, J.-L. Gaudiot, Data-flow and multithreaded architectures, in: *Encyclopedia of Electrical and Electronics Engineering*, Wiley, New York, 1997.
- [72] I. Koren, B. Mendelson, I. Peled, G.M. Silberman, A data-driven VLSI array for arbitrary algorithms, *Computer* 21 (10) (1988) 30–43.
- [73] J.-L. Gaudiot, L. Bic (Eds.), *Advanced Topics in Data-Flow Computing*, Prentice-Hall, New York, 1991.
- [74] Y. Kodama, S. Sakai, Y. Yamaguchi, A prototype of a highly parallel dataflow machine EM-4 and its preliminary evaluation. in: *Proceedings of the InfoJapan 90*, October 1990, pp. 291–298.
- [75] Y. Kodama, H. Sakane, M. Sato, H. Yamana, S. Sakai, Y. Yamaguchi, The EM-X parallel computer: architecture and basic performance, in: *Proceedings of the ISCA-22*, Santa Margherita Ligure, Italy, June 1995, pp. 14–23.
- [76] KRI, Khoral Research Inc. Web Page, <http://www.kri.com>.
- [77] F. Kurdahi, N. Bagherzadeh, The Morphosys Project, University of California, Irvine, 1998, <http://www.eng.uci.edu/morphosys>.
- [78] L. Lamport, How to make a multiprocessor computer that correctly executes multiprocess programs, *IEEE Trans. Comput.* 28 (9) (1979) 690–691.
- [79] P. Lapsley, J. Bier, A. Shoham, E.A. Lee, *DSP Processor Fundamentals Architectures and Features*, IEEE Press, New York, 1997.
- [80] P. Le Guernic, T. Gauthier, M. Le Borgne, C. Le Maire, Programming real-time applications with SIGNAL, *Proc. IEEE* 79 (9) (1991).
- [81] E.A. Lee, D.G. Messerschmitt, Static scheduling of synchronous data flow programs for digital signal processing, *IEEE Trans. Comput.* (1987).
- [82] E.A. Lee, Consistency in dataflow graphs, *IEEE Trans. Parallel Distributed Syst.* 2 (2) (1991).
- [83] E.A. Lee, T.M. Parks, Dataflow process networks, *Proc. IEEE* 83 (5) (1995) 773–801.
- [84] E.A. Lee, A denotational semantics for dataflow with firing, Memorandum UCB/ERL M97/3, Electronics Research Laboratory, UC Berkeley, January 1997.
- [85] D. Lenoski, J. Laudon, T. Joe, D. Nakahira, L. Stevens, A. Gupta, J. Hennessy. The DASH prototype: implementation and performance, in: *Proceedings of the ISCA-19*, Gold Coast, Australia, May 1992, pp. 92–103.
- [86] W.-Y. Lin, J.-L. Gaudiot, I-structure software cache: a split-phase transaction runtime cache system, in: *Proceedings of the 1996 Conference on Parallel Architectures and Compilation Techniques*, Boston, Massachusetts, October 1996, pp. 122–126.
- [87] W.H. Mangione-Smith, Seeking solutions in configurable computing, *IEEE Comput.* 30 (1997) 38–43.
- [88] W.H. Mangione-Smith, Application design for configurable computing, *Computer* 30 (1997) 115–117.
- [89] Z. Manna, A. Pnueli, *The Temporal Logic of Reactive and Concurrent Systems*, Springer, Berlin, 1991.
- [90] O.C. Maquelin, H.H.J. Hum, G.R. Gao, Costs and benefits of multithreading with off-the-shelf RISC processors, in: *Proceedings of the EURO-PAR '95*, no. 966 in *Lecture Notes in Computer Science*, Stockholm, Sweden, Springer, August 1995, pp. 117–128.
- [91] A. Márquez, K.B. Theobald, X. Tang, T. Sterling, G.R. Gao, A superstrand architecture and its compilation, CAPSL Technical Memo 18, Department of Elec. and Computer Engineering, University of Delaware, Newark, Delaware, March 1998.
- [92] A. Márquez, K.B. Theobald, X. Tang, G.R. Gao, A superstrand architecture and its compilation, in: *Proceedings of the MTEAC99 Workshop held in conjunction with HPCA-5*, Orlando, FL, January 1999.

- [93] B. Mendelson, B. Patel, I. Koren, Designing special purpose co-processors using the data-flow paradigm, in: J.-L. Gaudiot, L. Bic (Eds.), *Advanced Topics in Data-Flow Computing*, Prentice-Hall, Englewood Cliffs, NJ, 1991, pp. 547–570.
- [94] B. Mendelson, G.M. Silberman, Mapping data flow programs on a VLSI array of processors, in: *Proceedings of the International Symposium on Computer Architecture*, Pittsburgh, Pennsylvania, June 1987, pp. 72–80.
- [95] W.A. Najjar, W.M. Miller, A.P.W. Böhm, An analysis of loop latency in dataflow execution, in: *Proceedings of the ISCA-19*, Gold Coast, Australia, May 1992, pp. 352–360.
- [96] W. Najjar, J.-L. Gaudiot, Multi-level execution in data-flow architectures, in: *Proceedings of the ICPP '87*, St. Charles, Ill., August 1987, pp. 32–39.
- [97] W. Najjar, A.P.W. Bohm, B. Draper, R. Beveridge, *The Cameron Project*, Colorado State University, Fort Collins, CO 1998, <http://www.cs.colostate.edu/cameron>.
- [98] S.S. Nemawarkar, Performance modeling and analysis of multithreaded architectures, Ph.D. Thesis, Montréal, Qué., August 1996.
- [99] R.S. Nikhil, Arvind, Can dataflow subsume von Neumann computing? in: *Proceedings of the ISCA-16*, Jerusalem, Israel, May–June 1989, pp. 262–272.
- [100] K. Okamoto, S. Sakai, H. Matsuoka, T. Yokota, H. Hirono, Multithread execution mechanisms on RICA-1 for massively parallel computation, in: *Proceedings of the 1996 Conference on Parallel Architectures and Compilation Techniques*, Boston, Massachusetts, October 1996, pp. 116–121.
- [101] Oxford Hardware Compilation Group, *The Handel Language*, 1997.
- [102] E.K. Pauer, C.S. Myers, P.D. Fiore, J.M. Smith, C.M. Crawford, E.A. Lee, J. Lundblad, C.X. Hylands, Algorithm analysis and mapping environment for adaptive computing systems, Second Annual Workshop on High Performance Embedded Computing, MIT Lincoln Labs, Lexington, MA, September 1998 (<http://ptolemy.eecs.berkeley.edu/papers/98/ACSmapping/>).
- [103] J. Rasure, S. Kubica, The KHOROS application development environment, in: J.L. Crowley (Ed.), *Experimental Environments for Computer Vision and Image Processing*, HICa, New Jersey, 1994.
- [104] S.K. Reinhardt, J.R. Larus, D.A. Wood, Tempest and typhoon: user-level shared memory, in: *Proceedings of the ISCA-21*, Chicago, Ill., April 1994, pp. 325–336.
- [105] X. Tang, Compiling for multithreaded architectures, Ph.D. Thesis, University of Delaware, Newark, DE, April 1999.
- [106] S. Sakai, K. Okamoto, H. Matsuoka, H. Hirono, Y. Kodama, M. Sato, Super-threading: architectural and software mechanisms for optimizing parallel computation, in: *Conference Proceedings of the 1993 International Conference on Supercomputing*, Tokyo, Japan, July 1993, pp. 251–260.
- [107] H. Schmit, Incremental reconfiguration for pipelined applications, *IEEE Symposium on FPGAs for Custom Computing Machines*, 1997.
- [108] J.P. Singh, J.L. Hennessy, A. Gupta, Scaling parallel programs for multiprocessors: methodology and examples, *IEEE Computer* 6 (26) (1993) 42–50.
- [109] A. Sohn, M. Sato, N. Yoo, J.-L. Gaudiot, Data and workload distribution in a multithreaded architecture, *J. Parallel Distributed Process.* 1997, pp. 256–264.
- [110] X. Tang, J. Wang, K.B. Theobald, G.R. Gao, Thread partitioning and scheduling based on cost model, in: *Proceedings of the SPAA '97*, Newport, Rhode Island, June 1997, pp. 272–281.
- [111] X. Tang, G.R. Gao, How hard is thread partitioning and how bad is a list scheduling based partitioning algorithm? in: *Proceedings of the Tenth Annual ACM Symposium on Parallel Algorithms and Architectures*, Puerto Vallarta, Mexico, June 1998.
- [112] K.B. Theobald, EARTH: an efficient architecture for running threads, Ph.D. Thesis, Montréal, Qué., January 1999.
- [113] T. von Eicken, D.E. Culler, S.C. Goldstein, K.E. Schauer, Active messages: a mechanism for integrated communication and computation, in: *Proceedings of the ISCA-19*, Gold Coast, Australia, May 1992, pp. 256–266.
- [114] J. Vuillemin et al., Programmable active memories: reconfigurable systems come of age, *IEEE Trans. on VLSI Syst.* (1996).

- [115] E. Waingold, M. Taylor, D. Srikrishna, V. Srakar, W. Lee, V. Lee, J. Kim, M. Frank, P. Finch, R. Barua, J. Babb, S. Amarsinghe, A. Agrawal, Baring it all to software: raw machines, *IEEE Comput.* 30 (1997) 86–93.
- [116] M. Warren, J.K. Salmon, A parallel hashed oct-tree n -body algorithm, in: *Proceedings of the Supercomputing '93*, Portland, Oregon, November 1993, pp. 12–21.
- [117] S. Weiss, I. Spillinger, G.M. Silberman, Architectural improvements for a data-driven vlsi processing array, *J. Parallel Distributed Comput.* 19 (4) (1993) 308–322.
- [118] N. Wirth, *The Hardware Description Language Lola*, June 1995.
- [119] M.J. Wirthlin, B.L. Hutchings, DISC : the dynamic instruction set computer, in: J. Schewel (Ed.), *Field Programmable Gate Arrays for Fast Board Development and Reconfigurable Computing*, vol. SPIE 2607, 1995, pp. 92–103.
- [120] M.J. Wirthlin, B.L. Hutchings, Improving functional density through run time constant propogation, in: *International Symposium on Field Programmable Gate Arrays*, April 1997, pp. 86–92.