Fahimeh Yazdanpanah, Carlos Alvarez-Martinez, Daniel Jimenez-Gonzalez, and Yoav Etsion, Member, IEEE

**Abstract**—General purpose hybrid dataflow/von-Neumann architectures are gaining attraction as effective parallel platforms. Although different implementations differ in the way they merge the conceptually different computational models, they all follow similar principles: they harness the parallelism and data synchronization inherent to the dataflow model, yet maintain the programmability of the von-Neumann model. In this paper, we classify hybrid dataflow/von-Neumann models according to two different taxonomies: one based on the execution model used for inter- and intrablock execution, and the other based on the integration level of both control and dataflow execution models. The paper reviews the basic concepts of von-Neumann and dataflow computing models, highlights their inherent advantages and limitations, and motivates the exploration of a synergistic hybrid computing model. Finally, we compare a representative set of recent general purpose hybrid dataflow/von-Neumann architectures, discuss their different approaches, and explore the evolution of these hybrid processors.

Index Terms—Dataflow architectures, von-Neumann model, parallel processors, hybrid systems, scheduling and task partitioning

# **1** INTRODUCTION

**D**OWER-EFFICIENCY is today one of the main challenges in computer architecture. One of the approaches for tackling this challenge is the use of homogeneous and heterogeneous multi-core architectures that help to: 1) use more power-efficient cores, and 2) exploit the existing parallelism on the applications. These multi-core architectures are conventionally based on the von-Neumann (traditional control flow) computing model, which is inherently sequential because of its use of a program counter and an updateable memory. Nevertheless, the von-Neumann computing model is able to exploit some limited instruction level parallelism (ILP), data level parallelism (DLP), and thread level parallelism (TLP). However, DLP and TLP should be explicitly expressed by the programmer and/or compiler, while ILP is limited by the sequential execution of the instructions.

The dataflow model is a recurrent alternative to the von-Neumann execution model. The dataflow computing model is known to overcome the limitations of the traditional control flow model by fully exploiting the parallelism inherent in programs. In the dataflow model, the operands trigger the execution of the operation to be performed on them. In other words, dataflow architectures use the availability of data to fetch instructions rather than the availability of instructions to fetch data. Unlike the

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TPDS.2013.125 von-Neumann model, the dataflow model is neither based on memory structures that require inherent state transitions, nor does it depend on history sensitivity and program counter to execute a program sequentially. These properties allow the use of the model to represent maximum concurrency to the finest granularity and to facilitate dependency analysis among computations.

In this sense, the dataflow model has the potential to be an elegant execution paradigm with the ability to exploit the inherent parallelism available in applications. Furthermore, this model is self-scheduled and more powerefficient than the control flow model, which shows inefficiencies [41], [51], [56]. However, although the benefits of the dataflow have been known for a long time, this model has not yet been fully exploited for commercial systems. In fact, implementations of the model have failed to deliver the promised performance because the dataflow model has some inefficiencies and limitations. One significant drawback of the dataflow model is its inability to effectively manage data structures, memory organizations and traditional programming languages.

Therefore, in order to increase the performance and power efficiency of multi-core systems, these systems can be designed as hybrid architectures that combine the dataflow and von-Neumann models of computation. The convergence of the dataflow and control flow execution models allows for the incorporation of conventional control flow execution into the dataflow approach, or for exploiting a dataflow approach in von-Neumann architectures. This alleviates the inefficiencies associated with both methods. Hybrid dataflow/von-Neumann models therefore bind the power of the dataflow model for exposing parallelism together with the execution efficiency of the von-Neumann model in order to overcome the limitations of both models. While different hybrid implementations differ in the way they merge the two conceptually different execution models, they all follow similar principles.

The objective of this paper is to provide a better understanding of the evolution of the hybrid models and

1045-9219 © 2013 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications\_standards/publications/rights/index.html for more information.

F. Yazdanpanah, C. Alvarez-Martinez, and D. Jimenez-Gonzalez are with the Universitat Politècnica de Catalunya (UPC) and also with the Barcelona Supercomputing Center (BSC), 08034 Barcelona, Spain. E-mail: {fahimeh, calvarez, djimenez}@ac.upc.edu.

Y. Etsion is with the Electrical Engineering and Computer Science Departments, Technion-Israel Institute of Technology, Haifa 32000, Israel. E-mail: yetsion@tce.technion.ac.il.

Manuscript received 10 May 2012; revised 13 Feb. 2013; accepted 9 Apr. 2013. Date of publication 18 Apr. 2013; date of current version 16 May 2014. Recommended for acceptance by M. Kandemir.

their main characteristics. We classify them according to two taxonomies: one based on the execution model used for inter- and intrablock execution, and the other based on the integration level of both the control flow and the dataflow models. Using these taxonomies, we classify a representative set of recent general purpose hybrid models (works from the year 2000 or later), absent to the best of our knowledge from other surveys [101], [109], [110], as well as summarizing their main features and compare their benefits and issues. However, in order to acquire a fully historical point of view, we also describe some of the previous main contributions on hybrid models. On the other hand, to keep the length of this survey within bounds, software frameworks and specific purpose dataflow accelerators are regarded as beyond the scope of the paper.

The rest of the article is organized as follows: Section 2 discusses the von-Neumann (control flow) computing model. Section 3 overviews the dataflow computing model as well as different dataflow architectures. Section 4 presents hybrid dataflow/von-Neumann models, and classifies them according to two taxonomies. In Section 5, we describe some recent general purpose hybrid dataflow/ von-Neumann architectures. A comparison and discussion on main features of recent hybrid architectures and their common trends are given in Section 6. Finally, we conclude in Section 7.

# 2 VON-NEUMANN COMPUTING MODEL

The von-Neumann computation model [128] is the most common and commercially successful model to date. The main characteristic of this model is a single separate storage structure (the memory) that holds both program and data. Another important characteristic is the transfer of control between addressable instructions, using a program counter (PC). The transfer is either implicit (auto-increment of PC) or through explicit control instructions (jumps and branches, assignment to PC). It is for this reason that the von-Neumann model is commonly referred to as a control flow model.

A key tenet of the model is the set of memory semantics it provides in which loads and stores occur in the order in which the PC fetched them. Enforcing this order is required to preserve true (read-after-write), output (write-afterwrite), and anti (write-after-read) dependences between instructions.

Furthermore, the serial execution of instructions is a hallmark of the von-Neumann architecture. However, this simplistic sequential execution, together with data, control and structural hazards during the execution of instructions, may be translated into an under-utilization of the hardware resources. In that sense, exploiting parallelism at different granularities—instruction level parallelism (*ILP*), data level parallelism (*DLP*), and thread level parallelism (*TLP*)—is a mechanism for increasing hardware resource utilization.

Pipelined (IBM Stretch 1959 [12]) and superscalar [2] processors that try to process several instructions at the same time are the most common examples of **ILP**. Arguably the most notable class of superscalar processors is that of

the dynamically scheduled Out-of-Order processors [92] that maintain a window of pending instructions dispatching them in dataflow manner. In all these processors, parallelism is further enhanced by using a set of techniques such as register renaming, branch prediction and speculative execution, which are used in addition to dynamically dispatching independent instructions in parallel to multiple functional units (see details in Section 5.1). Another way of exploiting ILP is by means of very long instruction word (VLIW) processors [36]. The explicitly parallel instruction sets for VLIW enable the compiler [31] to express instruction independence statically in the binary code, thereby reducing the necessary hardware support for dynamically managing data and control hazards in Out-of-Order processors.

Architectures with **DLP** apply a single operation to multiple, independent data elements. Probably the most common examples of DLP are the single instruction multiple data (SIMD) extensions. SIMD extensions are mechanisms that statically express parallelism in the form of a single instruction that operates on wide, multi-element registers (a method sometimes referred to as sub-word parallelism). These extensions appeared in supercomputers such as the Thinking Machines CM-1 [55] and CM-2 [19], and are now ubiquitous in all general purpose processors. A derivative of SIMD processors, known as the single instruction multiple thread (SIMT) architecture, is nowadays common in graphics processing units (GPUs) [87].

Finally, TLP (or multi-threading) is applied by executing parallel threads on separate processing units. Nevertheless, some architectures utilize this coarse-grain parallelism to hide memory latencies and improve the utilization of hardware resources by interleaving multiple threads on a single physical processor. This technique is known as simultaneous multi-threading (SMT) [124], [130] and has been implemented in large machines such as HEP [112] and Tera [4] (as well as many others [1], [75], [129]). SMT has even made it to consumer products, starting with the Pentium 4 [81] and Power 5 [17] processors. However, despite all these efforts, effective utilization of parallel von-Neumann machines is inherently thwarted by the need to synchronize data among concurrent threads. Thread synchronization and memory latencies were identified [7] as the fundamental limitations of multiprocessors.

The need for efficient data synchronization has grave programmability implications and has placed emphasis on the cache coherency and consistency in shared-memory machines, particularly as the number of processing units continuously increase [14]. Transactional memory architectures [53] aim to alleviate that problem somewhat by providing efficient and easy-to-use lock-free data synchronization. Alternatively, speculative multithreading architectures exploit TLP dynamically by scheduling the threads in parallel [114], as Out-of-Order architectures for instructions, masking the synchronization issues. Experience shows that multithreaded control flow machines are feasible, though some doubt their scalability due to two major issues that limit their parallel processing capabilities: memory latency and synchronization. In summary, improvements in the memory system, ILP, DLP and TLP significantly reduce the memory latency issue of von-Neumann architectures, but they are still limited by the execution in control flow manner. On the other hand, the dataflow architectures can overcome this limitation thanks to the exploitation of the implicit parallelism of programs [7], [23].

## **3 DATAFLOW COMPUTING MODEL**

The dataflow computing model represents a radical alternative to the von-Neumann computing model. This model offers many opportunities for parallel Processing, because it has neither a program counter nor a global updatable memory, i.e., the two characteristics of the von-Neumann model that inhibit parallelism. Thanks to these properties, it is extensively used as a concurrency model in software and as a high-level design model for hardware.

The principles of dataflow were originated by Karp and Miller [65]. They proposed a graph-theoretic model for the description and analysis of parallel computations. Shortly after, in the early 1970s, the first dataflow execution models were developed by Dennis [26] and Kahn [64]. Dennis originally applied the dataflow idea to the computer architecture design while Kahn used it in a theoretical context for modeling concurrent software.

The dataflow model is self-scheduled since instruction sequencing is constrained only by data dependencies. Moreover, the model is asynchronous because program execution is driven only by the availability of the operands at the inputs to the functional units. Specifically, the firing rule states that an instruction is enabled as soon as its corresponding operands are present, and executed when hardware resources are available. If several instructions become fireable at the same time, they can be executed in parallel. This simple principle provides the potential for massive parallel execution at the instruction level. Thus, dataflow architectures implicitly manage complex tasks such as processor load balancing, synchronization, and accesses to common resources.

A dataflow program is represented as a directed graph, called *dataflow graph* (DFG). This consists of named nodes and arcs that represent instructions and data dependencies among instructions, respectively [24], [66]. Data values propagate along the arcs in the form of packets, called *tokens*. A DFG can be created at different computing stages. For instance, it can be created for a specific algorithm used for designing special-purpose architectures (common for signal processing circuits). However, most dataflow-based systems convert a high-level code into DFG at compile time, decode time, or even during execution time, depending on the architecture organization. Unlike control flow programs, binaries compiled for a dataflow machine explicitly contain the data dependency information.

In practice, real implementation of the dataflow model can be classified as static (single-token-per-arc) and dynamic (multiple-tagged-token-per-arc) architectures. The first dataflow architecture [28] followed the static model. This approach allows at most one token to reside on any arc. This is accomplished by extending the basic firing rule as follows: A node is enabled as soon as tokens are



Fig. 1. DFG of a loop (a) the static and (b) the dynamic dataflow.

present on its input arcs and there is no token on any of its output arcs [29]. In order to implement the restriction of having at most one token per arc, and to guard against nondeterminacy, extra reverse arcs carry acknowledge signals from consuming to producing nodes [29].

The implementation of the static dataflow model is simple, but since the graph is static, every operation can be instantiated only once, and thus loop iterations and subprogram invocations can not proceed in parallel. Fig. 1a shows an example of static dataflow graph for computing a loop which is executed *N* times sequentially (note that in this figure, the graph for controlling iteration of the loop is not illustrated). Despite this drawback, some machines were designed based on this model, including the MIT Dataflow Architecture [28], [30], DDM1 [25], LAU [96], and HDFM [125].

The dynamic dataflow model tries to overcome some of the deficiencies of static dataflow by supporting the execution of multiple instances of the same instruction template, thereby supporting parallel invocations of loop iterations and subprogram. Fig. 1b shows the concurrent execution of different iterations of the loop. This is achieved by assigning a *tag* to each data token representing the dynamic instance of the target instruction (e.g.,  $a_1, a_2, \ldots$ ). Thus, an instruction is fired as soon as tokens with identical tags are present at each of its input arcs. This enabling rule also eliminates the need for acknowledge signals, increases parallelism, and reduces token traffic. Dynamic dataflow machines employ two types of control instructions: Data-steering instructions and Tag management instructions. Data-steering instructions explicitly guide data values to the correct path after a branch, which is a control flow instruction. Each live value requires its own datasteering instruction [26]. Tag management instructions are inserted into tagged-token dataflow programs to differentiate between multiple dynamic instances of named program values (e.g. variables in executing iterations of a loop simultaneously). Notable examples of this model are the Manchester Dataflow [50], the MIT Tagged-Token [6], DDDP [70] and PIM-D [62].

The dynamic dataflow can execute out-of-order, bypassing any token with complex execution and delays the remaining computation. Another noteworthy benefit of this model is that little care is required to ensure that tokens remain in order.

The main disadvantage of the dynamic model is the extra overhead required to match tags on tokens. In order to reduce the execution time overhead of matching tokens, dynamic dataflow machines require expensive associative memory implementations [50]. One notable attempt to eliminate the overheads associated with the token store is the *Explicit Token Store* (ETS) [22], [54]. The idea is to allocate a separate memory frame for every active loop iteration and subprogram invocation. Since frame slots are accessed using offsets relative to a frame pointer, the associative search is eliminated. To make that concept practical, the number of concurrently active loop iterations must be controlled. Hence, the condition constraint of k-bounded loops was proposed [9], which bounds the number of concurrently active loop iterations. The Monsoon architecture [90] is the main example of this model.

The dataflow model has the potential to be an elegant execution paradigm with the ability to exploit inherent parallelism available in applications. However, implementations of the model have failed to deliver the promised performance due to inherent inefficiencies and limitations. One reason for this is that the static dataflow is unable to effectively uncover large amount of parallelism in typical programs. On the other hand, dynamic dataflow architectures are limited by prohibitive costs linked to associative tag lookups, in terms of latency, silicon area, and power consumption.

Another significant problem is that dataflow architectures are notoriously difficult to program because they rely on specialized dataflow and functional languages. Dataflow languages are required in order to produce large dataflow graphs that expose as much parallelism as possible to the underlying architecture. However, these languages have no notion of explicit computation state, which limits the ability to manage data structures (e.g., arrays). To overcome these limitations, some dataflow systems include specialized storage mechanisms, such as the I-structure [8], which preserve the single assignment property. Nevertheless, these storage structures are far from generic and their dynamic management complicates the design.

In contrast, imperative languages such as C, C++, or Java explicitly manage machine state through load/store operations. This modus operandi decouples the data storage from its producers and consumers, thereby concealing the flow of data and making it virtually impossible to generate effective (large) dataflow graphs. Furthermore, the memory semantics of C and C++ support arithmetic operations on memory pointers, which result in memory aliasing, where different semantic names may refer to the same memory location. Memory aliasing cannot be resolved statically, thus further obfuscating the flow of data from between producers and consumers. Consequently, data-flow architectures do not effectively support imperative languages.

In summary, the dataflow model is effective in uncovering parallelism, due to the explicit expression of parallelism among dataflow paths and the decentralized execution model that obviates the need for a program counter to control instruction execution. Despite these advantages, programmability issues limit the usefulness of dataflow machines. Moreover, the lack of a total order on instruction execution makes it difficult to enforce the memory ordering that imperative languages require. While this section describes the key features, characteristics and limitations of the dataflow model, a complete survey of the model is beyond the scope of this paper. For further details, we refer the reader to more extensive literature on the subject [85], [115], [126].

# 4 HYBRID DATAFLOW/VON-NEUMANN MODELS

The inherent limitations of both dataflow and von-Neumann execution models motivate the exploration of a convergent model that can use synergies to leverage the benefits of both individual models. Therefore, the hybrid models try to harness the parallelism and data synchronization inherent to dataflow models, while maintaining existing programming methodology and abstractions that are largely based on von-Neumann models. While different hybrid implementations differ in the way they merge the two conceptually different models, they all follow similar principles.

Most notably, hybrid models alleviate the inefficiencies associated with dataflow model, either by increasing the basic operation granularity or by limiting the size of the DFG. Additionally, they incorporate control flow abstractions and shared data structures. As a result, different hybrid architectures employ a mix of control flow and dataflow instruction scheduling techniques using different partial scheduling methods. Furthermore, in the hybrid models, nodes of a DFG vary between a single instruction (fine-grain) to a set of instructions (coarse-grain).

A further significant benefit of hybrid models is clearly evident in their memory models. Hybrid models combine single assignment semantics, inherent to dataflow, with consistent memory models that support external sideeffects in the form of load/store operations. This relieves one of the biggest (if not the biggest) restriction of pure dataflow programming: the inability to support a shared state, and specifically shared data structures [85]. Therefore, hybrid models are capable of executing imperative languages. As a result, combining dataflow and von-Neumann models facilitates designing efficient architectures that benefit from both computing models, while the remaining issue concerns the best granularity-parallelism trade-off.

# 4.1 Evolution of Hybrid Architectures until 2000

The first idea of combining dataflow and control flow arose in the early 1980s [63], [99], [112], [123], and included data and memory structure management (e.g., Multithreaded Monsoon (MT. Monsoon) [91]), self-scheduling and asynchronous execution to simplify thread synchronization (e.g., HEP [63], [112]; Tera [4]; MT. Monsoon [91]), as well as the ability to execute both conventional and dataflow programs in the same machine [7], [14]. Some hybrid models [14], [60] even included a program counter to a dataflow architecture in order to execute sequential instructions in control flow manner. In this regard, other studies explored the threaded dataflow model [101], [109], in which partial data sub-graphs are processed as von-Neumann instruction streams. In particular, given a dataflow graph (program), each sub-graph that exhibits a low degree of parallelism is identified and transformed into a sequential thread of instructions. Such a thread is issued consecutively by the matching unit without matching further tokens, except for the first instruction of the thread. Data passed between instructions in the same thread is stored in registers instead of being written back to memory. These registers may be referenced by any succeeding instruction in the thread. This improves single-thread performance, because the total number of tokens needed to schedule program instructions is reduced, which in turn saves hardware resources. In addition, pipeline bubbles caused by runtime overhead associated with token matching are avoided for dyadic (two-operand) instructions within a thread. Two threaded dataflow execution techniques can be distinguished: 1) the direct token recycling technique, which allows cycle-by-cycle instruction interleaving of threads in a manner similar to multithreaded von-Neumann computers (e.g., MT. Monsoon architecture), and 2) consecutive execution of the instructions of a single thread technique (e.g., Epsilon [46], [47] and EM-4 [5] architectures). In the second technique, the matching unit is enhanced with a mechanism that, after firing the first instruction of a thread, delays matching of further tokens in favor of consecutive issuing of all instructions of the started thread. In addition, some architectures based on threaded dataflow use instruction pre-fetching and token pre-matching to reduce idle times caused by unsuccessful matches. EM-4 [5], EM-X [71] and RWC-1 [108] are examples of this kind of architectures, which are also referred to as macro-dataflow [78].

Until the late 80s and early 90s, the common wisdom was that fine-grain execution was much more suited to masking network and memory latencies than a coarsegrain execution, and would obviously provide a much better load leveling across processors and hence faster execution. However, it has been demonstrated that coarsegrain execution is equally suited to exploit parallelism as fine-grain [83], [86], [122]. On one hand, Gao's group [39], [58], [121] was the first to develop a coarse-grain data flow simulator and compiler from scratch, and to report on very extensive evaluations of very complex applications. On the other hand, Najjar's group [13], [34], [35], [85], [102] focused on modifying the Sisal compiler [77] in two ways: 1) by generating coarse-grained data flow code from a finegrained one, and 2) by generating coarse-grained data flow code from scratch using the Sisal compiler.

In addition to the coarsening of nodes in the DFG, another technique for reducing dataflow synchronization frequency (and overhead) is the use of complex machine instructions, such as vector instructions. With these instructions, structured data is referenced in block rather than element-wise, and can be supplied in bursts while also introducing the ability to exploit parallelism at the sub-instruction level. This technique introduces another major difference with conventional dataflow architectures; that is, tokens do not carry data (except for the values true or false). Data is only moved and transformed within the execution stage. Examples of such machines are Stollman [42], ASTOR [133], DGC [34], [35], and SIGMA-1 multiprocessor [132].

In parallel, the Out-of-Order model [59], [92], which emerged in the late 80s, incorporated the dataflow model to extract ILP from sequential code. This approach has been



Fig. 2. Inter- and intrablock scheduling of organizations of hybrid dataflow/ von-Neumann architectures. (a) Enhanced control flow, (b) control flow/ dataflow, (c) dataflow/control flow, and (d) enhanced dataflow. Blocks are squares and big circles.

further developed by Multiscalar [114] and thread level speculation (TLS) [98], [103], which can be viewed as coarse-grain versions of Out-of-Order.

Efforts have been made to survey hybrid models up to year 2000 [101], [109], [110], and also dataflow multithread models [27], [61], [68], [78]. However, to the best of our knowledge, since 2000 there has been no a comprehensive survey describing hybrid architecture. Hence, the main focus of this paper is on classifying recent hybrid data-flow/von-Neumann architectures, which have mainly attempted to improve the conventional architectures exploiting several aspects of dataflow concepts [32], [67], [76], [105], [118], or to utilize the dataflow approach as accelerators [44], [84], [127]. Most of those recent works are classified and compared in the following sections.

# 4.2 Taxonomy Based on Block Execution Semantics

The inherent differences between dataflow and von-Neumann execution models appear to place them at two ends of a spectrum that covers a wide variety of hybrid models. However, with our approach, the coarsening of the basic operation granularity, from a single instruction to a block of instructions, together with the inter- and intrablock execution semantics, enable us to partition the spectrum into four different classes of hybrid dataflow/ von-Neumann: Enhanced Control Flow, Control Flow/Dataflow, Dataflow/Control Flow and Enhanced Dataflow class. This taxonomy is based on whether they employ dataflow scheduling inside and/or between code blocks. Block is defined on the basis of the boundary between where the two scheduling models (inter- and intrablock scheduling) are mainly applied. In this way, the number of instructions in a block (block granularity) depends on the specific model. Fig. 2 illustrates inter- and intrablock scheduling of conventional organizations of hybrid dataflow/von-Neumann architectures.

# 4.2.1 Enhanced Control Flow Class

Models in this class schedule blocks in control flow manner, whereas the instructions within a block are scheduled in a mixed approach of control flow and dataflow manner. Fig. 2a illustrates the organization of this class.

The main example of this class is the Out-of-Order (restricted dataflow) model [59], [92]. The Out-of-Order

model, as an extension of superscalar processors, incorporates the dataflow model only in the issue and dispatch stages to extract ILP from sequential code. It is also referred to as local dataflow or micro dataflow architecture [101], [109], [110].

## 4.2.2 Control Flow/Dataflow Class

Models in this class schedule the instructions within a block in dataflow manner, whereas blocks are scheduled in control flow manner (Fig. 2b). This method is used in RISC dataflow architectures, which support the execution of existing software written for conventional processors.

Main examples of this class are TRIPS [105], [106], Tartan [84], Conservation cores (C-Cores) [127], DySER [44] and other architectures that rely on domain specific dataflow accelerators. TRIPS was a new design that tried to overcome the foreseen limitations of large core architectures by adding new layers of flexibility to the hardware. Explicit dataflow execution within blocks was a necessary way to improve fine-grain ILP while keeping hardware complexity within bounds. TRIPS unifies dataflow and von-Neumann into a single execution model. However, other architectures in this class essentially use dataflow to accelerate parts of the code (hyperblocks in Tartan; kernels in C-Cores; phases in DySER). Their decision on which parts of the code to accelerate is mostly static while TRIPS uses dynamic scheduling decisions to map hyperblocks to dataflow cores. Tartan, C-Cores and DySER use profiling to determine the parts of the code to be accelerated in a dataflow unit (or units), mapped on reconfigurable hardware coupled to a classical von-Neumann processor. Unlike Tartan and C-Cores, DySER also supports reconfigurations at runtime. This behavior allows the DySER architecture to capture a significant percentage of computation from a single application as multiple accelerated phases can be mapped to the same accelerator.

## 4.2.3 Dataflow/Control Flow Class

Models in this class employ dataflow rules between blocks and control flow scheduling inside the blocks (Fig. 2c). A block is a set of sequential instructions, where data is passed between instructions using register or memory (coarse-grain dataflow models [101], [109], [110]). Under these restrictions, blocks are issued by the matching unit, and token matching needs only to be performed on a block basis. Thus, the total number of tokens needed to schedule program instructions is reduced, which in turn saves hardware resources. Main examples of models of this class are: Star-T (\*T) [89], TAM [21], ADARC [117], EARTH [57], [121], P-RISC [88], MT. Monsoon [91], Pebbles [102], SDF<sup>1</sup> [67], DDM [76], and Task Superscalar [32].

Fig. 3 shows a further decomposition of this class based on the number of cores and number of instructions in a block (i.e., size of block) targeted by every specific model, as well as the year in which it was first published. Fig. 3a depicts the relationship between core granularity and the publication year of the proposed architectures. First hybrid Fig. 3. Different architectures of Dataflow/Control Flow class. (a) Number of cores and year. (b) Number of cores and size of blocks.

(b)

designs tend to have a small number of cores, while recently proposed architectures tend to use a larger number of cores. Fig. 3b shows the variance in core granularity in hybrid design. Architectures with a larger number of cores typically use fewer numbers of instructions per block, and designs with a fewer number of cores tend to use larger blocks (with more than 1000 instructions per block).

## 4.2.4 Enhanced Dataflow Class

Models in this class use dataflow firing rules for instructions inside the blocks and for the blocks themselves. In effect, this class consists of two-level dataflow models (Fig. 2d) utilizing some concepts of the von-Neumann model (e.g., storage management) to add the abilities of running imperative languages and managing data structures. SIGMA-1 [132], Cedar [74] and WaveScalar [118] are the main examples in this class.

# 4.2.5 Comparison of Hybrid Classes

Every one of the four classes presents advantages and drawbacks. *Enhanced Control Flow class* machines can very naturally execute control flow codes and uncover more ILP than the strict von-Neumann models. However, as the technology only allows them to address small to medium block sizes, the amount of parallelism they can expose is typically limited (some architectures such as Kilo-instruction Processors [20] try to overcome this problem by targeting much larger block sizes).

*Control Flow/Dataflow class* machines try to overcome the limitations of the previous class by forcing the pure

EARTH Pebbles Few TAM ADARC (<20) vea 1985 1990 1995 2000 2005 2010 (a) # cores M.Monsoon SDF TSS Large DDM (>100)\*T Medium EARTH P-RISC TAM ADARC Few Pebbles (<20)#inst. /block Basic block(<100) Block Hyper block(>1000)

TŚS

DDM

SDF

M.Monsoon

# cores

Large

(>100)

Medium

<sup>1.</sup> Please note that here SDF is the acronym for scheduled dataflow, as opposed to synchronous dataflow (SDF) [79]. The latter is a dataflow based execution model for signal processing algorithms and does not include any von-Neumann properties.

dataflow execution of the instructions inside a block. These models attempt to expose ILP statically at the block level, deferring memory operations to inter-block synchronization. Indeed, the *Control Flow/Dataflow* general strategy has shown a great potential in both performance and power savings [44], [84], although it poses the same problems as the previous class (e.g., smaller block sizes than desirable for fully exploiting dataflow advantages at ILP level).

For their part, *Dataflow/Control Flow class* models have taken advantage of the recent growth in the number of parallel hardware structures in cores, chips, machines and systems. As models in this class address parallelism at a coarse grain, they are able to exploit all these resources more effectively than conventional (von-Neumann) models while retaining the programming model inside the blocks.

Finally, *Enhanced Dataflow class* models constitute a complete rethinking of the execution problem. Since they do not use a program counter, they face several difficulties when executing conventional codes and managing memory organizations, and therefore need more hardware resources to be used effectively. On the other hand, *Enhanced Dataflow class* models may be regarded as an addition to both *Dataflow/Control Flow* and *Control Flow/Dataflow* classes, and in this sense they posses great potential.

### 4.3 Taxonomy Based on Execution Model

Hybrid models can also be classified from an execution model point of view; *unified-hybrid* models versus *dataflow accelerator* models. In a unified-hybrid architecture, a program must be executed using both dataflow and control flow scheduling since both models are intimately bound in the architecture. Although the majority of the models presented belong to this group, it does present some drawbacks. The additional hardware needed by the interconnection and synchronization mechanisms (e.g., hardware of Out-of-Order architectures) leads to more complexity and power consumption. Furthermore, as all programs should be executed with the same hybrid scheduling schema, they are not able to adapt to specific cases in which a pure dataflow or von-Neumann model would be better.

On the other hand, in dataflow accelerator based architectures, the decision about which parts of the code to accelerate is mostly static (made by the programmer or compiler, and sometimes based on profiling). In addition, a whole program may be executed without the use of the dataflow accelerator. As mentioned above, Tartan, C-Cores and DySER are architectures that use dataflow to accelerate kernels (or hyperblocks) and thus belong to this group.

# 5 EXAMPLES OF RECENT HYBRID DATAFLOW/ VON-NEUMANN ARCHITECTURES

In this section, we describe recent examples of hybrid dataflow/von-Neumann architectures for each of the above mentioned taxonomy classes, in chronological order. Out-of-Order architectures [59], [92], [113] are presented for *Enhanced Control Flow class*. Although Out-of-Order appeared before 2000, we included it here because of its popularity and its significant contribution to the class, as well as to highlight how the introduction of dataflow

execution into an otherwise control flow model can dynamically extract parallelism. TRIPS [105], [106], WaveScalar [118], [120] and Task Superscalar [32], [33] are presented for *Control Flow/Dataflow*, *Enhanced Dataflow*, and *Dataflow/Control flow classes* respectively. Although DySER [44] belongs to the *Control Flow/Data flow class*, it has been included as a recent representation of a wide range of pure control flow processors that use dataflow accelerators.

Other relevant architectures exist which, due to limitations of space, included in the supplementary file which is available in the Computer Society Digital Library at http:// doi.ieeecomputersociety.org/10.1109/125; for instance, MT. Monsoon and the SIGMA-1 multiprocessor, while non-recent, are highly relevant, representations of *Dataflow/ Control flow class* and *Enhanced dataflow class* architectures, respectively. Moreover, DDM and SDF architectures are relevent and recent *Data flow/Control Flow* architectures that are also included in the supplementary file available online.

Main characteristics of all foregoing architectures are described in Table 1 and discussed in Section 6.

# 5.1 Out-of-Order Execution Model

The *Out-of-Order* architecture (Restricted Dataflow) [59], [92], [113] is a fine-grain hybrid architecture belonging to the *Enhanced Control Flow class*. The Out-of-Order architecture is also referred to as local dataflow or micro dataflow architectures [101], [109].

# 5.1.1 Execution Model

Out-of-Order processors employ dataflow principles to extract instruction level parallelism (ILP) and optimize the utilization of the processor's resources. The processor relies on hardware mechanisms that determine dynamically data dependencies among the instructions in the instruction window. In other words, in this paradigm, a processor executes instructions in an order governed by the availability of input data, rather than by their original order in a program.<sup>2</sup> In doing so, the processor can both extract ILP and hide short data fetch latencies by processing subsequent instructions that are ready to run. Each instruction window of Out-of-Order processor is a block granularity for the intrablock scheduling.

### 5.1.2 Architecture Organization

Fig. 4 shows the general scheme of the Out-of-Order execution pipeline. Instructions are fetched in order, then after register renaming are decoded and placed into a pool of pending instructions (the instruction window) and the reorder buffer. The reorder buffer saves the program order and the execution states of the instructions. For increasing the effective instruction window size, these architectures rely on branch prediction and speculation. Therefore, they require complex check-pointing mechanisms to recover from branch mis-predictions and mis-speculated executions (not shown in the Figure).

Dispatch and Issue determine the out-of-order and dataflow execution of the microprocessor. The matching of the executable instructions in the microprocessor is restricted to the pending instructions of the instruction

<sup>2.</sup> The memory accesses are done in order.

	SIGMA-1	Out-of- Order	MT. Monsoon	DDM	SDF	TRIPS	Wavescalar	Task Superscalar	DySER (w/ inorder GPP)
Year	1982	1985	1991	2000	2001	2003	2003	2010	2011
ISA	DFG w/ complex machine ops.	RISC / CISC	DFG w/ thread extensions	RISC / CISC	RISC (preload/ store + computation)	EDGE	WaveScalar	RISC / CISC	RISC /CISC + Ultra-wide insts.
Main Features	Structure-flow processing, Chained hashing hardware	Out-of- order instruction execution	Explicit token store, Multithreaded	Decoupled non-blocking, CacheFlow policy (prefetching)	Decoupled non-blocking multithreading	Polymorphous (multiple parallelism modes), memory ordering, branch prediction	Wave-ordered and unordered memory, hierarchically interconnection	Out-of-order task execution	Profiling based detection of ultra-wide insts. compiler support to execute them on DySER blocks
Core Granularity	SIGMA-1 PE and SE	PE agnostic	Monsoon PE	PE agnostic	Simple processor	4x4 ALU (ET) array	5 stage dynamically scheduled pipeline	PE agnostic	DySER block: 8x8 FU / GPP: PE agnostic
Scalability	> 100 PEs+SEs	~ 10 FUs	> 1000PEs	~ 100 PEs	~ 100 PEs	~ 100 PEs	> 1000 FUs	>> 100 PEs	< 10 DySER Blocks
Parallelism level	ILP	ILP, DLP, TLP (dual threaded)	TLP	TLP	TLP	ILP, TLP, DLP	TLP, ILP	TLP	ILP, DLP, TLP (dual threaded)
Block Granularity	Vector length	Inst. window size	Thread (BB)	BB size (codeblock, > one thread, in TSU graph mem.)	BB size, <128- inst. blocks (27(15 in EP) <51 (39 in EP))	128-inst. block (EDGE)	<64 insts. per PE cache, any num. of insts. in a wave	Task size (any size) > 10K	Equivalent to few hundreds of ISA insts.
Inter-block Scheduling	DF	CF	DF	Dynamic DF (dependencies specified in programs)	Static DF (programmer/ compiler)	Static CF (compile time)	Dynamic DF (dependencies detected at execution time)	Dynamic DF (dependencies specified in programs)	CF
Intra-block Scheduling	Hybrid DF/CF	Hybrid CF/DF	CF (direct token recycling)	CF	CF (scheduled DF)	Static DF (compile time)	Dynamic DF (execution time)	CF	Static DF (compile time)
Inter-block Communication	Memory/ direct inter- connection	Register /cache /memory	Token/ memory	Cache	Frame memory and registers	Registers	Memory /direct Inter-connection	Memory	Register /memory / FIFO
Intra-block Communication	Memory/ direct inter- connection	Register /cache /memory	Register / memory	Register / memory	Register	Memory / direct Inter-connection	Memory / direct Inter-connection	Register / memory	Direct interconnection
Examples	SIGMA-1	Many	MT. Monsoon	D2NOW, Flux, DDM-VMc	SDA	TRIPS, TFlex	WaveCache	Task Superscalar	DySER

 
 TABLE I

 Comparison of the Recent Hybrid Dataflow/von-Neumann Architectures. DF, CF, and DFG Stand for Dataflow, Control Flow, and Dataflow Graph, Respectively

window. The matching hardware can therefore be restricted to a small number of instructions slots. In addition, because of the sequential program order, the instructions in this window are likely to be executable soon.

Once the instructions are executed, they are retired to permanent state machine (memory) in source program order (*commit* in the Figure). Another advantage of these architectures is their sequential execution of the instruc-



Fig. 4. Out-of-Order execution pipeline.

tions, exploiting the spatial locality of the program. This locality enables a memory hierarchy to be employed for storing the instructions and data, potentially executed in the following cycles and close to the executing processor.

## 5.1.3 Implementation Examples

Arguably the first Out-of-Order execution machine was the CDC 6600 (1964), which used a scoreboard to resolve conflicts. The IBM 360/91 (1966) introduced Tomasulo's algorithm, supporting full Out-of-Order execution. In 1990, the first Out-of-Order microprocessor appeared, the POWER1, but its Out-of-Order execution was limited to floating point instructions.

As mentioned above, Out-of-Order microprocessors have an instruction window that is restricted to a sequence of instructions. Thread-level speculation (TLS) processors may be regarded as an extension of Out-of-Order hybrid dataflow/von-Neumann architecture that increases the instruction window and potentially uncover more ILP. TLS is a technique which empowers the compiler to identify potential parallel threads, despite uncertainty as to whether those threads are actually independent [116]. TLS allows these threads to be speculatively executed in parallel, while squashing and re-executing any thread that suffers dependence violations. The instruction window is thus the addition of the sequence of instructions of all non-speculative and speculative threads executing in parallel, potentially larger than the Out-of-Order instruction window.

Thread creation, and the mechanism for buffering speculative state and tracking data dependences between speculative threads, are important features of the different TLS approaches. Some of them are implemented entirely in software [49], [98], [104], others in hardware [3], [80], [103], while others are a combination of software and hardware [16], [37], [43], [52], [72], [73], [97], [134]. Two relevant works are the LRPD test [98] (software-only support) and the Multiscalar architecture [114]. LRPD test allows the compiler to parallelize loops without fully disambiguating all memory references, and applies only to array-based codes. Disambiguation is performed with the use of shadow arrays to detect any cross-iteration Read-after-Write dependence. Multiscalar architecture was the first complete evaluation of an architecture designed specifically for supporting TLS. The compiler statically performs the distribution of the instructions among tasks (potential speculative threads). Address resolution buffer (ARB) [38], forward and release bits, and CFG (control flow graph) information are the mechanisms used for tracking control and data dependences between speculative threads.

## 5.2 TRIPS

*TRIPS* (Tera-op, Reliable, Intelligently adaptive Processing System) [105], [106] was designed at the University of Texas in Austin as a grid architecture for implementing the EDGE (Explicit Data Graph Execution) ISA [15], [105], [111]. It is an example of *Control Flow/Dataflow class* models.

#### 5.2.1 Execution Model

TRIPS combines control flow execution across hyperblocks of code consisting of up to 128 instructions with a dataflow execution inside the blocks. In TRIPS, a hyperblock is equivalent to the block granularity. This scheme enforces conventional memory semantics across hyperblocks and so enables imperative code to be executed without major modifications.

The TRIPS architecture is fundamentally block oriented. The compiler is responsible for statically scheduling each block of instructions onto the computational engine such that inter-instruction dependences are explicit. Therefore, the compiler role is vital for the final performance of the application [18]. Each block has a static set of state inputs, and a potentially variable set of state outputs that depends upon the exit point from the block. At runtime, the basic operational flow of the processor includes fetching a block from memory, loading it into the computational engine, executing it to completion, committing its results to the persistent architectural state if necessary, and then proceeding to the next block.

TRIPS has a block-atomic execution mode which means a block of instructions must be fetched and executed as though it were a single unit providing interruptions at block level. It provides direct communication for instructions within a block therefore instructions within a block can directly send values to dependent instructions within the same block. This behavior allows the architecture to have very large windows (up to 1024 instructions) that execute in dataflow order.



Fig. 5. TRIPS architecture (figure based on [105], [106]).

TRIPS provides three modes of execution that enable *polymorphous parallelism*: Desktop-morph (D-morph) is aimed at exploiting instruction level parallelism (ILP), Thread-level morph (T-morph) for thread level parallelism (TLP), and Stream-level morph (S-morph) for data level parallelism (DLP). How TRIPS works in these different modes is explained in the following subsection.

#### 5.2.2 Architecture Organization

Fig. 5 shows the TRIPS architecture that is a tiled and distributed. TRIPS processor consists of four cores and a tiled secondary memory (M tiles in Fig. 5a), surrounded by a tiled network (N tiles in the Figure) that acts as translation agents for determining where to route memory system requests.

Each of the TRIPS cores is implemented using five unique tiles: one global control tile (GT),  $4 \times 4$  execution tiles (ET), four register tiles (RT), four data tiles (DT), and five instruction tiles (IT), as shown Fig. 5b. Each tile only interacts with its immediate neighbors through microarchitectural networks (micronets). Micronets have roles such as transmitting operands between instructions, distributing instructions from the instruction tiles to the execution tiles, or communicating control messages from the program sequencer [107]. The major processor core micronet is the operand network (OPN), which handles transport of all data operands along ETs, DTs, GTs, and RTs.

The GT contains the PC running of the blocks in the TRIPS core; the instruction cache tag arrays; the I-TLB, and the next-block predictor. The GT handles TRIPS block management, including prediction, fetch, dispatch, completion detection, flush (on mis-predictions and interrupts), and commit. In addition, GT is used to set up the control register that configures the processor into different

speculation, execution, and threading modes. GT also maintains the state of all in-flight blocks (maximum 8) running in the ETs of the TRIPS core. When a block finishes, the block predictor (tournament local/gshare predictor-based) provides the predicted address of the next target block. The block is fetched and loaded into the reservation stations of ET.

Each ET consists of a fairly standard single-issue pipeline; a bank of 128 reservation stations (two-operand instructions); an integer unit; a floating point unit, and an operand router (shown in Fig. 5c). When a reservation station contains a valid instruction and a pair of valid operands, the node can select the instruction for execution. After execution, the node can forward the result to any of the operand slots in local or remote reservation stations within the ALU array ( $4 \times 4$  ETs' ALUs). Instructions are statically placed into the locations of the ET, and executed in dataflow manner using the direct instruction communication between intrablock producers and consumers, specified by the TRIPS ISA.

The Instruction Cache is tiled into five banks (i.e., IT) to increase the memory bandwidth. Each IT acts as a slave of the GT which holds the single tag array.

The register file is divided into four 32-register banks (tiles) that are nodes of the OPN micronet, allowing the compiler to place critical instructions that read and write from/to a given bank close to that bank. The registers file holds a portion of the architectural state, so values passed between hyperblocks, where direct instruction communication is not possible, are transmitted through the register file.

The primary memory is divided into four data tiles (DT). Each DT holds one L1 data cache bank. It can be accessed by any ALU through the local grid routing network.

Some of those hardware resources can be configured by using the GT, to operate differently depending on the mode: D-morph, T-morph and S-morph. For instance, the reservation stations can be managed differently depending on the execution mode. A physical frame is formed by the reservation stations with the same index across all the execution tiles (e.g., combining the first slot for all nodes in the grid forms frame 0). Frames containing one hyperblock form an architectural frame (A-frame). Thus, direct instruction communication is only possible within an A-frame.

In D-morph, all the frame space of a TRIPS core can be used, since it is a large, distributed, instruction issue window, by only one thread, allowing it to achieve maximum ILP. In addition, in order to increase the potential ILP, the hardware fills empty A-frames with speculatively mapped hyperblocks, predicting which hyperblock will be executed next, mapping it to an empty A-frame, and so on. The A-frames are treated as a circular buffer where the first is non-speculative, and the rest are speculative. When the non-speculative A-frame finishes, the first speculative A-frame becomes the non-speculative first A-frame of the circular buffer. In T-morph, the frame space is statically partitioned, so each thread can have its own frame partition. Within each thread, speculation is also used but extra prediction registers (e.g., block control state for each of the hardware threads) are needed. In S-morph, only one thread can be run and no speculation is performed. Instead, inner loops of a streaming application are unrolled to fill the reservation stations within multiple A-frames fused in a *super* A-frame. In this case, to reduce the power and instruction-fetch bandwidth overhead incurred by repeated fetching of the same code block across inner-loop iterations, the S-morph employs mapping reuse, in which a block is kept in the reservation stations and used multiple times. In this case, the L2 cache memory can be configured to be used as a stream register file [106], so that direct data array access and DMA transfer capabilities are allowed. Otherwise, the secondary memory works as a non-uniform cache access (NUCA) on-chip memory system.

### 5.2.3 Implementation Examples

Some studies have been carried out on different aspects of TRIPS. Sankaralingam *et al.* [107] describe the TRIPS control protocols. They detail each of the five types of reused tiles that compose the processor; the control and data networks that connect them, and the distributed micro-architectural protocols that implement instruction fetch, execution, flush, and commit. They also describe the physical design issues of implementation the micro-architecture in a 170 M transistor, 130 nm ASIC prototype chip composed of two 16-wide issue distributed cores and a distributed 1 MB NUCA on-chip memory system.

Gratz *et al.* [48] presented the design, implementation and evaluation of the TRIPS on-chip network (OCN) which is a wormhole-routed,  $4 \times 10$  2D mesh network with four virtual channels. These authors discussed the tradeoffs made in the design of the OCN; in particular, why area and complexity were traded off against latency. A full evaluation of a real TRIPS ASIC prototype and an EDGE compiler [40] demonstrates that the TRIPS machine is feasible. This work also shows that TRIPS is competitive with a Pentium 4 system in the number of cycles needed to execute an application. It is an impressive outcome for a new machine fully developed in an academic environment.

TFlex is another architecture based on the EDGE ISA. It is an implementation of the composable lightweight processor (CLP) [69], which is proposed to eliminate the problem of fixed-granularity processors, and consists of multiple simple, narrow-issue processor cores that can be aggregated dynamically to form more powerful singlethreaded processors. TFlex adds four capabilities to TRIPS in a distributed fashion: I-cache management; next-block prediction; L1 D-cache management, and memory disambiguation hardware. Robatmili *et al.* [100] presented a hardware implementation of mapping blocks to a distributed substrate of composable cores for the TFlex.

### 5.3 WaveScalar

*WaveScalar* [118], [120], is an example of the *Enhanced Dataflow class*. It is a dynamic, general purpose, decentralized superscalar dataflow architecture that is product of research performed at the University of Washington. WaveScalar is also the name of the dataflow instruction set and the execution model.



Fig. 6. Dataflow graph and wave-ordered memory (figure based on [120]).

### 5.3.1 Execution Model

The WaveScalar execution model is basically a dataflow model enhanced to support imperative languages. The key tenet of the Wavescalar execution model is that programs execute in *waves*, which are sets of connected instructions of the program graph. The wave name may come from the way data flows from the initial instruction to subsequent instructions in parallel.

Formally, a wave is a connected, directed acyclic portion of the control flow graph with a single entrance. The WaveScalar compiler partitions an application into maximal waves and adds wave management instructions. In fact, waves are similar to hyperblocks, but they may contain control flow joins and are generated using loop unrolling to make them larger (all instructions within a wave are partially ordered, so waves can not contain loops). In order to allow instructions to operate on different dynamic waves, all data elements travel with their wave number, which increases as the data goes out of a wave and enters a new one (or the same) using a special *waveadvance* instruction.

Therefore, in order to execute an imperative program in Wavescalar, it is compiled into an special code that contains the dataflow graph (i.e., the wave) as well as the memory order. An example of the memory order problem is illustrated in Fig. 6. It is assumed that the Load instruction must execute after the Store instruction to ensure correct execution because the two memory addresses are identical. In a pure dataflow graph this implicit dependence between the two instructions (the dashed line in Fig. 6) can not be expressed.

However, Wavescalar supports a wave-ordered memory mode in which the compiler annotates memory access instructions within each wave to encode the ordering constraints between instructions thereby forming a chain of memory instructions. A memory request can therefore only be executed if the previous request in the chain and all memory requests from the previous wave have already been executed. In order to be successful, the compiler must ensure that there is a complete chain of memory operations along every path through a wave. So, if there are no memory operations in one of the paths of a branch, a MemNop instruction must be inserted into that path to maintain the chaining. Furthermore, in order to increase parallelism (i.e., in loops), nondependent memory accesses can also be annotated with an additional ripple number that allow loads to execute in parallel and even out of order if all previous stores have finished.

By implementing wave-ordered memory, a total ordering of memory instructions can be achieved with little dynamic overhead. This feature, alongside indirect jumps for object linking, allows traditional von-Neumann models of computation to execute just as fast—if not faster—on the dataflow architecture. Its main advantage is that it is a dataflow hardware that runs programs written in standard programming languages, by efficiently providing the sequential memory semantics required by imperative languages.

In addition to wave-ordered memory, a second memory scheme in WaveScalar (*standard data firing rule* mode) allows the programmer to omit any unnecessary ordering and intertwine memory operations into the program graph by using the standard data firing rule. The unordered memory scheme introduces a new store instruction, storeunordered-ack, which returns zero to signify when it has been completed. Using this value as an input arc to other instructions enforces memory ordering while providing greater flexibility to the programmer.

Both wave-ordered and unordered memory can be used interchangeably within the same program or even within the same wave, to take advantage of fine-grain (unordered) and coarse-grain (wave-ordered) threads, resulting in significant performance improvements [118], [120].

#### 5.3.2 Architecture Organization

To execute WaveScalar programs, a scalable, tile-based processor architecture called *WaveCache* has been designed. Fig. 7 shows the WaveScalar architecture. Each basic processing element (PE) is a five-stage (Input, Match, Dispatch, Execute, and Output), dynamically scheduled execution pipeline. In WaveScalar, pairs of PEs are coupled into pods sharing ALU results via a common bypass network. Four pods (8 PEs) are grouped into a domain that communicate over a set of pipelined buses. Four domains (32 PEs) form a cluster supported by conventional memory hierarchy. In order build larger machines, multiple



Fig. 7. WaveScalar architecture (figure based on [118]).

clusters can be connected by a 2D mesh interconnection network.

Wave-ordered memory lies in the WaveCache's store buffers (one per cluster), which are responsible for implementing the wave-ordered memory interface that guarantees correct memory ordering. To reduce communication costs, the PEs are connected through a hierarchical interconnection infrastructure. WaveScalar's hierarchical interconnect plays an important role in overall scalability. Swanson *et al.* [119] have studied the area-performance trade-offs for WaveScalar.

The placement scheme of the instructions of a program has a compile-time and a runtime component. The compiler is responsible for grouping instructions into segments. Those segments have up to 64 instructions. As a program executes, the WaveCache maps the program instructions onto its array of PEs, placing a whole segment of instructions at the same PE. The instructions remain at their PEs for many invocations, and as the working set of instructions changes, the WaveCache removes unused instructions and maps new ones in their place. The instructions communicate directly with one another over a scalable, hierarchical on-chip interconnect, obviating the need for long wires and broadcast communication.

# 5.3.3 Implementation Examples

The only implementation of WaveScalar is WaveCache. WaveScalar uses a regular native DEC compiler for converting source code to Alpha binary, and a binary translator is used for translating an Alpha binary to a WaveScalar binary. Petersen et al. [95] present and analyze three compiler optimizations for the Wavescalar C compiler that significantly reduce control overhead with minimal additional hardware. The basis of the solution lies in recognizing that overhead instructions are relatively simple to implement in hardware and can generally execute in parallel with computation. Hence, the micro-architecture can be tuned to execute overhead instructions in parallel with computation instructions. Merzulo et al. [82] proposed the transactional WaveCache to exploit speculative execution of memory operations. Pei et al. [93] exploited speculative multithreading (SpMT) based on WaveScalar.

### 5.4 Task Superscalar

The *Task Superscalar* [32], [33] is a task-based dataflow architecture which generalizes the operational flow of dynamically scheduled Out-of-Order processors. It was designed at the Barcelona Supercomputing Center (BSC) and belongs to the *Dataflow/Control Flow class*. The Task Superscalar combines the effectiveness of Out-of-Order processors in uncovering parallelism together with the task abstraction, thereby providing a unified management layer for CMPs, which effectively employs processors as functional units.

### 5.4.1 Execution Model

The Task Superscalar processor combines dataflow execution of tasks with control flow execution within the tasks (i.e., the block granularity is a task). As ILP pipelines uncover parallelism in a sequential instruction stream, similarly, the Task Superscalar uncovers task level



Fig. 8. Task Superscalar architecture (figure based on [33]).

parallelism among tasks generated by a sequential thread. By utilizing intuitive programmer annotations of task inputs and outputs, the Task Superscalar pipeline dynamically detects inter-task data dependencies, identifies task-level parallelism, and executes tasks out-of-order. This design enables programmers to exploit many-core systems effectively, while simultaneously simplifying the programming model.

### 5.4.2 Architecture Organization

The high-level operational flow of the Task Superscalar is illustrated in Fig. 8. A task generator thread resolves the inter-task control path and sends non-speculative tasks to the pipeline front-end for dependency decoding. The task window may consist of tens of thousands of tasks, which enables it to uncover large amounts of parallelism [33]. The pipeline asynchronously decodes the task dependencies, generates the task dependency graph (with tasks as nodes and dependencies between tasks as arcs), and schedules tasks as they become ready. Finally, ready tasks are sent to the execution backend, which consists of a task scheduler and a queuing system.

As shown in Fig. 8, the front-end employs a tiled design, and is managed by an asynchronous point-to-point protocol. The front-end is composed of four module types: the *pipeline gateway; task reservation stations* (TRS); *object renaming tables* (ORT) and *object versioning tables* (OVT).

The gateway is responsible for controlling the flow of tasks into the pipeline including by allocating TRS space for new tasks and distributing tasks to the different modules, as well as stalling the task generator thread whenever the pipeline fills. TRSs store the in-flight task information and track the readiness of task operands. Inter-TRSs communication is used to register consumers with producers, and notify consumers when data is ready. The TRSs store the meta-data of all in-flight tasks.

The ORTs map parameters onto the most recent task by accessing the same memory object, and thereby detect object dependencies. Storing all data users (either producer or consumer), rather than only storing real data for producers, facilitates TRS consumer chaining. The OVTs track live operand versions, which are created whenever a new data producer is decoded. Each OVT is associated with exactly one ORT. The functionality of the OVTs is similar to a physical register file, but only for maintaining operand meta-data. Effectively, the OVT manages data anti- and output-dependencies, either through operand renaming or by chaining different output operands and unblocking them in-order by sending a ready message when the previous version is released.

Fig. 8 also shows the Out-of-Order components equivalent to the Task Superscalar modules. In Out-of-Order processors, dynamic data dependencies operates by matching each input register of a newly fetched instruction (consumer) with the most recent instruction that writes data to that register (producer). The instruction is sent to a reservation station to wait until all its inputs become available. Hence, the reservation stations effectively store the instruction dependency graph, which consists of all inflight instructions. In the Task Superscalar, the mechanism of decoding tasks identifies all possible effects a task may have on the shared processor state, so producers and consumers are identified correctly. Moreover, tasks are decoded in-order to guarantee correct ordering of producers and consumers, and specifically, that the decoding of a task producing a datum updates the renaming table, before any task consuming the datum performs a lookup.

#### 5.4.3 Implementation Examples

Etsion *et al.* [32] presented a design for a distributed Task Superscalar pipeline front-end which can be embedded into virtually any many-core fabric and manages it as a Task Superscalar multiprocessor. The Task Superscalar architecture uses the StarSs programming model [10], [94]. This programming model supports Out-of-Order execution of tasks, by enabling programmers to explicitly expose task side-effects by using annotating operands of kernel functions as input, output, or inout. The model can thus decouple the execution of the thread generating the tasks from their decoding and execution. At runtime, whenever the task generator thread reaches a call site at one of the kernels, the task creation code (injected by a source-tosource compiler) packs the kernel code pointer and all the operands, and writes them to the task pipeline.

Yazdanpanah *et al.* [131] presented an FPGA-based prototype of the Task Superscalar architecture. The implemented hardware is based on a tiled design that can operate in parallel and is easily scalable to manage hundreds of cores in the same way that Out-of-Order architectures manage functional units. The prototype operates at near 150 Mhz and maintains up to 1024 inflight tasks, managing the data dependencies in few cycles.

# 5.5 DySER

The *DySER* (Dynamically Specialized Execution Resource) [44], [45] is an architecture based on dataflow accelerators belonging to the *Control Flow/Dataflow class*. It was designed at the University of Wisconsin-Madison as the hardware substrate of the dynamically specialized execution (DySE) model.



Fig. 9. Processor pipeline with DySER (figure based on [44]).

## 5.5.1 Execution Model

The DySER integrates dataflow accelerators (DySER block) into a control flow processor pipeline as functional units. To achieve this goal, the program is explicitly partitioned by the compiler (profile-guided to determine common path trees of control flow) into phases (i.e., program sections). After that, for each phase, the compiler determines its kernels and tries to accelerate them by using the DySER block. The DySER block is basically a large, reconfigurable, functional unit composed of different arithmetic units whose connections are reconfigured at runtime, creating specialized dataflow blocks that can be pipelined. The model hinges on the assumption that only a few dataflow blocks are active during a given phase of an application and they are invoked several times. Thus, setting up the static routes once amortizes the execution of the DySER unit over many invocations.

In order to execute each kernel within a phase with a pure dataflow accelerator, kernels are divided into a load-back slice and a computation slice. A load-back slice includes all the memory accesses, while a computation slice consists of computation operations without memory accesses that are grouped and executed in the dataflow accelerator. With this separation between slices, the usual processor memory disambiguation optimizations can proceed unhindered. Therefore, in DySER, these computation slices are the block granularities for the intrablock scheduling.

To make better use of the reconfigurable unit, when the control flow execution reaches a program phase, the DySER block is dynamically configured (specialized for the phase). Furthermore the execution model allows the use of multiple DySER blocks, where each block is configured differently. With multiple DySER blocks, the next block can be predicted and configured before its inputs are produced by the processor. The large granularity of the phases allows easy predictability.

### 5.5.2 Architecture Organization

Fig. 9 illustrates the DySER attached to a processor pipeline. The DySER block consists of a circuit-switched network of heterogeneous functional units. The functional units (FUs) form the basic computation fabric. Each FU is connected to four neighbor switches, from where it gets data and control input values and injects outputs. Each FU also includes a configuration register that specifies which function to perform, as well as one data register and one status register for each input switch. The status registers indicate the validity of values in the data registers. The data registers match the word-size of the machine. The switches (Ss) contain data and status registers, and include a configuration register which specifies the input to output port mappings.

DySER blocks are configured by writing into configuration registers at each functional unit and switch. After configuration, the switches in the DySER block form a circuit-switched network that creates explicit hardware paths from inputs to the functional units between functional units, and from functional units to outputs. The functional units are configured to perform the operation required to execute the desired dataflow graph. The idea is that for a given application phase, DySER blocks are configured once and re-used many times.

The basic execution inside a DySER block is dataflow driven by values arriving at a functional unit. When the valid bits for both left and right operands are set, the functional unit consumes those inputs, and a fixed number of cycles later produces the output by writing into the data and status register of the output switch.

All the inputs to a DySER block are fed through a logical FIFO, which delivers register inputs and memory values. Each entry specifies a switch and a port. Since a DySER block uses circuit-switched routing, this effectively decides where the value will be delivered in the block. Outputs follow a similar procedure; each port in the output switches corresponds to one possible DySER block output. Since for each output port, the DySER produces outputs in order, no FIFOs are required on the output side. When values arrive at the output ports, an output interface writes them to the corresponding register or memory.

DySER can be easily integrated into conventional inorder and Out-of-Order pipelines as an accelerator. Integration with an in-order pipeline is simple, and the DySER block interfaces with the instruction fetch stage for obtaining the configuration bits, the register file stage and the memory stage of the pipeline. A state machine must be added to the instruction cache to read configurations bits for a DySER block and send them to the input interface of that DySER block.

DySER integration with an Out-of-Order pipeline requires more careful design. The processor views DySER as a functional unit, although the input ports should be exposed to the issue logic to ensure two send operations are not executed out-of-order. Since loads can cause cache misses, when a load executes in the processor, the corresponding input port is marked as *busy* in the input buffers. When the data arrives from the cache, the input port is marked as *ready*, which prevents subsequent loads values from entering the DySER block earlier.

#### 5.5.3 Implementation Examples

Govindaraju *et al.* [44] implemented the DySER block in Verilog and synthesized it using Synopsys compiler with a 55 nm standard cell library. These authors developed pathtree, a program representation for application phases, in order to find the most frequently executed basic blocks for mapping on DySER. For evaluating DySER, they developed extensions to the GCC tool-chain which operates on the SPARC backend and performs path-profiling and DySER mapping. Benson *et al.* [11] described the integration of DySER into a commercial processor by designing an FPGA prototype based on the OpenSPARC T1 processor called the *OpenSPlySER*. Govindaraju *et al.* [45] studied challenges for DySER on data parallel workloads.

# 6 DISCUSSION OF RECENT HYBRID MODELS

This section highlights the main features of recent hybrid architectures, described in Section 5 and in the supplementary file available online, compares and discusses them, and shows their common trends. Table 1 introduces the main features of the architectures sorted according to the *year* in which the architecture appeared.

# 6.1 Main Features

Out-of-Order, DDM, Task Superscalar and DySER are based on RISC/CISC ISA. In addition, DySER has the ultrawide instructions, which are used to run part of the program in the reconfigurable DySER blocks. SDF is based on a RISC ISA defined for the execution and synchronization processors. SIGMA-1, MT.Monsoon, TRIPS and WaveScalar are based on dataflow ISAs. SIGMA-1 has complex machine operations to reduce loop operations and structure-flow processing overheads.

The main feature of SIGMA-1 is the structure-flow processing that allows it to perform vector processing efficiently. In the case of Out-of-Order it is the dataflow execution of a sequential instruction stream. The Explicit Token Store (ETS), which eliminates the associative search in the matching unit, and multithreading are the main features of MT. Monsoon architecture. The main feature of DDM is the introduction of the CacheFlow policy, which implies the execution of a DDM thread (basic block of instructions—BB) only if its data is already placed in the cache. Decoupling computation and synchronization, and non-blocking threads are also the main features of SDF and DDM. However, computation in the DDM is carried out by an off-the-shelf processor, while in the SDF it is carried out by a custom designed processor. Another difference is that in SDF data is preloaded in registers, while in DDM data is pre-fetched in the cache. Polymorphism is one of the main features of TRIPS and provides three modes of execution for exploiting one of the three types of parallelism ILP, TLP and DLP. The main feature of the WaveScalar is waveordering execution. In the wave-ordered memory of WaveScalar, memory instructions are annotated with extra information that orders them relative to other instructions of a block. The main feature of Task Superscalar is Out-of-Order task execution. DySER architecture introduces the idea of generic dataflow accelerators integrated within a general purpose processor through ultra-wide instructions. Those generic dataflow accelerators are dynamically configured at execution time.

The computational core granularity varies from any processing element (PE) or core size in the case of DDM, Task Superscalar and Out-of-Order processors to a small SDF core. MT. Monsoon uses the original dataflow Monsoon PE to sequentially execute the thread instructions using the direct token recycling technique. SIGMA-1 uses two types of operational units: the structure element (SE) and the processing element (PE). The SEs perform structure-flow (vector) processing so that data is referenced in block rather than element-wise. SEs are tightly coupled to the PEs, that process elements in a dataflow manner. Each TRIPS processor consists of a tiled  $4 \times 4$  ALU (ET) array (the basic core of TRIPS) with 128 reservation stations per ET. Hyperblocks are statically placed, at compile time, on the reservation stations of the ALU array, and dynamically executed by the tiled architecture. DySER blocks consists on 8 × 8 FUs circuitswitched networks, that can be reconfigured to use different accelerators on runtime. The scalability also varies from more than 1000 PEs in the case of MT. Monsoon down to less than 10 DySER blocks in the DySER architecture.

### 6.2 Comparison and Discussion

### 6.2.1 Enhanced Control Flow Class

Out-of-Order architectures (restricted dataflow architectures) are the main representation of the *Enhanced Control Flow class*. Out-of-Order processors support ILP, DLP, and TLP in the form of dual threaded cores. The number of instructions of a block (block granularity) is that of the instruction window, created at runtime. Out-of-Order processors use cache, memory and registers to communicate data between blocks. Furthermore, Out-of-Order processors use hybrid control flow/dataflow intrablock scheduling and the same communication mechanisms as for inter-block communication.

The main difference between dynamic dataflow architectures and restricted dataflow pipelines is that the latter are designed to reconstruct the dataflow graph dynamically from a sequential instruction stream. The success of such a reconstruction relies on the ability to view a window of sequential code without control instructions and is largely attributed to accurate branch prediction and speculative execution. However, such processors are also susceptible to the prohibitive costs of branch mis-prediction, which require unrolling the execution of the wrongly predicted paths. This operation is particularly costly in deeply pipelined microprocessors. On the other hand, the restricted size of the instruction window limits the number of in-flight instructions and thus, to some extent, avoids the scalability issues associated with token stores in dynamic dataflow processors.

Nevertheless, the ILP achieved by Out-of-Order microprocessors is limited by the size of the instruction window and the amount of parallelism available in the instruction stream. In this sense, TLS may increase ILP by using speculative thread execution and a large instruction window. Indeed, thanks to the fact that each processor or processing unit only works with a limited part (i.e., instructions of thread) of the large instruction window, the complexity of concurrently monitoring the instruction issue of all the pending instructions, the data dependency cross-check complexity among the instructions, and the overall branch mis-prediction are reduced. Unlike dataflow models, TLS does not require a large waiting-matching store, although, it may suffer from costly check-pointing of memory accesses, squashing and re-executing threads.

### 6.2.2 Control Flow/Dataflow Class

The main representative of this class, TRIPS, constitutes a major effort in rethinking the computation of conventional codes, while also seeking to overcome the limitations of architectures based on big cores that is, large communication delays inside ever-growing control structures. The key is dataflow execution inside 128-instruction hyperblocks (intrablock dataflow scheduling), because it allows large instruction windows to be executed with reasonable hardware resources. This intrablock scheduling is static and defined at compile time. TRIPS uses both direct interconnection and memory for intrablock communication. For inter-block communication, TRIPS uses registers. Furthermore, when not enough ILP is available, TRIPS can use its polymorphous nature to enable different modes of execution. Therefore, it is also able to exploit DLP or TLP through loop-unrolling or parallel thread execution. On the other hand, commercial processors can obtain similar performance results by exploiting TLP through simultaneous multithreading, and DLP through SIMD instructions [40]. Indeed, small Out-of-Order instruction windows are sufficient to extract the available ILP efficiently in conventional codes. TRIPS may therefore be regarded as an efficient architecture that obtains similar results to classical processors with a different approach.

Another group of processors within this class use a dataflow accelerator statically defined by the compiler, inside a classical von-Neumann processor. DySER stands out from other dataflow accelerators because it is general purpose and presents some degree of runtime reconfiguration. Designed with power-efficiency in mind, the DySER execution model is based on the idea that a limited number of dataflow accelerators are enough to capture highly reused sections of the applications. TRIPS and DySER differ in that the former unifies dataflow and von-Neumann into a single execution model while the latter essentially uses dataflow to accelerate parts of the code.

DySER architecture can also support ILP, DLP, and TLP in the form of dual threaded cores. It supports DLP and TLP based on its incorporated general purpose processor (GPP) and the DySER blocks integrated with the GPP. For DySER, a block is a part of the program with as many as hundreds of GPP ISA instructions. Unlike TRIPS, DySER uses FIFOs to communicate input data with the DySER block. DySER provides static dataflow manner with direct interconnection for intrablock scheduling. Therefore, DySER requires profiling analysis of applications in order to pre-define the instructions that are going to be accelerated with the use of the DySER blocks. Once those sets of instructions are defined, the DySER execution model dynamically reconfigures the switched-network of functional units on the DySER block for each phase of the application. This dynamic reconfiguration provides area efficiency (rather than dynamically arbitrated networks) and programmability, although it requires compiler support and a phase predictor that tries to reconfigure-the DySER block before it is needed to hide the reconfiguration time. The need for profiling, together with the limited amount of runtime adaptability are the main disadvantages of this subclass. On the other hand, they are able to provide significant improvements in both performance and power efficiency over von-Neumann approaches, particularly for computation intensive kernels.

### 6.2.3 Dataflow/Control Flow Class

In these architectures, blocks are scheduled in a dataflow manner, while control flow scheduling is used within the blocks. Models in this class thus tend to provide specific support only to TLP. In particular, based on dependencies specified in the program, DDM and Task Superscalar perform dynamic dataflow inter-block scheduling by using cache and memory, respectively, for inter-block communication. SDF and MT. Monsoon perform static dataflow and both use memory and registers for inter-block communication. MT. Monsoon uses the V register, within the computation descriptor, that can be stored in a token of the pure dataflow Monsoon view. DDM, SDF and MT. Monsoon blocks are equivalent to a basic block, being up to 128 instructions in the case of a SDF/MT. Monsoon block. Task Superscalar may have blocks of any size.

This large class can be further divided into two groups according to the size of the blocks: either small or large. The sizes of the blocks of DDM, SDF, and MT. Monsoon model tend to be small, a decision that allows large amount of parallelism to be discovered and executed but also increases the cost of the synchronization. In the case of DDM, this characteristic makes the thread scheduling unit as important as the workstation duplicating the number of necessary processing elements. Another key point in this model is that in order to be efficient, it needs more information about the program than the classical control flow model. Programs should thus be annotated either by the compiler or by the programmer, which increases the complexity of the tool-chain needed to develop new applications. Unlike DDM, SDF executes the instructions within a block in-order, thereby obtaining less ILP but allowing the execute processor of its architecture to be simpler and smaller. Another characteristic of the SDF paradigm is that, although it can benefit from the annotated code, it can execute the original code as is, automatically extracting the available parallelism. MT. Monsoon, however, executes instructions within a block in-order using the direct token recycling technique. In addition, the thread extensions included in the MT. Monsoon facilitate the fork, join, and split phases of block executions.

The Task Superscalar is another instance of *Dataflow/ Control Flow class* architectures, but in this case the blocks are designed to be as large as desired. The Task Superscalar pipeline is designed as a generalization of Out-of-Order processors to the task-level. Nevertheless, its scalability goals, which target dynamically, managing very large graphs consisting of tens of thousands of nodes, require an alternative design to that of Out-of-Order processors. This redesign is the result of the Out-of-Order pipeline's use of reservation stations and bypass networks, whose operation is similar to that of associative token stores and are known not to scale.

The designers of the Task Superscalar pipeline thus opted for a distributed structure through careful protocol design that ubiquitously employs explicit data accesses, practically eliminates the need for associative lookups. The benefit of this distributed design is that it facilitates high levels of concurrency in the construction of the dataflow graph. These levels of concurrency trade off the basic latency associated with adding a new node to the graph with overall throughput. Consequently, the rate at which nodes are added to the graph enables high task dispatch throughput, which is essential for utilizing large many-core fabrics.

In addition, the dispatch throughput requirements imposed on the Task Superscalar pipeline are further relaxed by the use of tasks, or von-Neumann code segments, as the basic execution unit. The longer execution time of tasks compared to that of instructions means that every dispatch operation occupies an execution unit for a few dozen microseconds, and thereby further amplifies the scalability of the design.

The main disadvantage of most of models in this class is the need for annotating the original codes in order to extract a significant amount of parallelism from these codes. In this sense, a trend towards simplifying the annotations as much as possible can be observed in the designs of the programming models. Another common trend in this class is the increase in the number of processing elements and the size of the blocks.

## 6.2.4 Enhanced Dataflow Class

SIGMA-1 and WaveScalar are the main examples of the Enhanced Dataflow class. WaveScalar supports ILP and TLP, meanwhile SIGMA-1 supports ILP with extensions for vector processing. Unlike DySER and TRIPS, which need compiler support, WaveScalar performs dynamic dataflow intrablock scheduling, since the dependences are detected in execution time. For SIGMA-1, however, dependencies are explicit in the program to be executed. SIGMA-1 and WaveScalar use memory and direct interconnection for both intrablock and interblock communications. For interblock communication, SIGMA-1 and WaveScalar use memory and direct interconnection. A WaveScalar block is equivalent to a wave of instructions, although every PE caches up to 64 instructions, called segments. A SIGMA-1 block is considered as the set of instructions that are structure-flow processed.

The fact that WaveScalar is the only example of a mainly dataflow architecture able to execute imperative codes explains by itself the difficulties of such challenge. SIGMA-1 is able to execute an adapted Dataflow C language. SIGMA-1 was enhanced with control flow features in order to reduce dataflow synchronization frequency (and overhead) on vector processing. On the other hand, this uniqueness provides WaveScalar with a very interesting set of properties. Probably the main characteristic of this model is that it was designed with Moore's Law in mind to make of the most of the increase in transistor density and count. Therefore, the whole application would ideally be mapped to the PEs at the same time, and in this scenario (i.e., using kernels) it is expected to clearly outperform Out-of-Order processors. However, to achieve this goal the processor should have a larger number of PEs than what has been possible to date. In the

meantime, the need for "loading and discarding" instructions in the PEs throughout the program execution is one of its main bottlenecks. On the other hand, as technology evolves, this problem is expected to diminish, and WaveScalar architecture should be able to provide an approach capable of benefitting from the increasing transistor count, while keeping power consumption within bounds.

### 6.3 Common Trends

In addition to individual features of the classes discussed above, they share common properties and advantages. Moreover, they face similar challenges in their design. Recent hybrid architectures can handle imperative programming languages and data structures, as well as memory operations. This fact makes them stand out amongst other hybrid dataflow/von-Neumann architectures. Scheduling and memory management are key challenges in the design of hybrid architectures. One common theme among these architectures is their attempt to improve traditional processors by using dataflow principles at various levels in order to increase the capability of providing high levels of parallelism and performance. As the matter of fact, several features of the dataflow model, such as static single assignment, register renaming, dynamic scheduling and Out-of-Order instruction execution, I-structure-like synchronization and non-blocking threads are used in modern processor architectures and compiler technology. Moreover, many studies on hybrid models replace large, centralized processor cores with many simpler processing elements. In fact, all of these architectures, except WaveScalar, SIGMA-1 and MT. Monsoon, are von-Neumann machines, and rely on a program counter between blocks (inter-block) or inside blocks (intrablock), with some concepts of dataflow scheduling. WaveScalar eliminates the program counter and the register file and, relies completely on the dataflow program graph, thereby allowing the arcs between waves to define interactions between them.

The hybrid architectures discussed in this paper were developed as general purpose processors, although some of them may have not achieved their goals, as they failed to deliver the expected performance. Some of the hybrid architectures have limited scalability (e.g., Out-of-Order processors). In other cases, performance improvement was less than expected (e.g., TRIPS), while some of the hybrid architectures rely on new programming models (e.g., Task Superscalar and DDM). Most of them are not focused on power saving, although some dataflow based accelerators integrated with general purpose processors have been designed for energy efficiency such as C-Cores, Tartan and DySER.

# 7 CONCLUSIONS

This work surveys the recent general-purpose hybrid dataflow/von-Neumann architectures. To this end, we review the benefits and drawbacks of the von-Neumann and the dataflow computing models. We then present the common characteristics of the different hybrid models classifying them according to two different taxonomies that

enable a better understanding of their features. After that, we describe, compare and discuss a representative set of recent general-purpose hybrid dataflow/von-Neumann models. Finally, we present an insightful discussion with the aim of identifying trends in the next generation of hybrid architectures.

Nowadays, the majority of computer systems are based on the von-Neumann model. Such processors use a program counter to sequence the execution of instructions of a program and global updatable memory. Consequently, the von-Neumann machines have two fundamental limitations: memory latencies and thread synchronization. The dataflow model has no program counter or global updatable memory, so that dataflow architectures have the potential for exploiting all the parallelism available in programs. Since instructions in the dataflow models do not impose any constraints on sequencing except real data dependencies in programs, the dataflow model is asynchronous and self-scheduled.

However, although the dataflow model has been investigated since 1970s, no commercially viable global pure dataflow system has been implemented. The amount of parallelism discovered by the model becomes an implementation issue due to token matching and memory resource limitations. In theory, the dataflow model offers better performance and power efficiency than the von-Neumann model. The main reasons are the parallelism inherent to this model and the absence of overhead on pipeline control structures and temporary state (i.e., register file). Nevertheless, the efficient parallel programming of the dataflow architectures is difficult due to the fact that dataflow and functional languages do not easily support data structures, and they are not popular. On the other hand, imperative languages cannot be compiled to dataflow architectures, mainly because of issues associated with memory semantics.

Research on modern microprocessor architectures revealed the advantages of dataflow concepts in the use of instruction level parallelism. Indeed, in order to build efficient dataflow based machines, the dataflow model has to exploit some concepts belonging to the von-Neumann computing model. Similarly, most von-Neumann based architectures borrow concepts and mechanisms from the dataflow world to simplify thread synchronization and tolerate memory latency. As a result, the dataflow and von-Neumann models are not orthogonal, but are at two ends of a continuum. Combination or even unification of von-Neumann and dataflow models is possible, and is preferred to treating them as two unrelated, orthogonal computing paradigms. Recent dataflow research incorporates more explicit notions of state into the architecture, and von-Neumann models use many dataflow techniques to improve tolerance to long latency operations of modern multithreaded systems.

Hybrid architectures exploit the benefits of dataflow while preserving von-Neumann capabilities and imperative languages, in order to obtain high performance and low power architectures. Our findings are that most studies of hybrid designs exploit dataflow concepts in von-Neumann based architectures, particularly in superscalar and VLIW systems, in order to increase the capability of providing high levels of parallelism. On the other hand, some architects of the hybrid models have attempted to increase the efficiency of dataflow based architectures by using some ideas from control flow models. Although it is a promising area, most of the discussed hybrid architectures are not focused on power saving.

Designing a general-purpose architecture is a common goal, and all recent hybrid architectures discussed in this paper were developed as general-purpose processors. Moreover, it is also clear that modern hybrid architectures are designed to have the ability of handling imperative programming languages and data structures as well as memory organizations. Another observed trend is that architects of recent hybrid models have attempted to replace centralized processors by several simpler processing elements, since scheduling and memory management pose key challenges in their designs. An increase in the number of processing elements may be observed, as well as the fact that all the architectures try to use the dataflow principles at the level (ILP, DLP or TLP) envisioned by their designers with the most potential parallelism. At the same time, the von-Neumann scheduling is maintained at the other levels to keep the required resources within bounds.

### ACKNOWLEDGMENT

This work is supported by the Ministry of Science and Technology of Spain and the European Union (FEDER funds) under contract TIN2007-60625, by the Generalitat de Catalunya (contract 2009-SGR-980), and by the European FP7 project TERAFLUX id. 249013. The authors wish to thank Mark Oskin for his insightful comments on earlier drafts of this document, and the blind reviewers that, with their comments, have significantly improved the quality of the paper.

### REFERENCES

- A. Agarwal, B.-H. Lim, D. Kranz, and J. Kubiatowicz, "APRIL: A Processor Architecture for Multiprocessing," in *Proc. Int'l Symp. Comput. Architect.*, 1990, pp. 104-114.
- *Comput. Architect.*, 1990, pp. 104-114.
  [2] T. Agerwala and J. Cocke, "High Performance Reduced Instruction Set Processors," IBM T.J. Watson Res. Center, Yorktown Heights, NY, USA, Tech. Rep. RC12434, 1987.
- [3] H. Akkary and M.A. Driscoll, "A Dynamic Multithreading Processor," in Proc. ACM/IEEE Int'I Symp. Microarchitect., 1998, pp. 226-236.
- [4] R. Alverson, D. Callahan, D. Cummings, B. Koblenz, A. Porterfield, and B. Smith, "The Tera Computer System," in *Proc. Int'l Symp. Supercomput.*, 1990, pp. 1-6.
- [5] A. Vikraman, L. Bic, and T. Ungerer, "Evolution of Dataflow Computers," in *Advanced Topics in Data-Flow Computing*, J.-L. Gaudiot and L. Bic, Eds. Englewood Cliffs, NJ, USA: Prentice Hall, 1991.
- [6] A. Vikraman and D.E. Culler, "Dataflow Architectures," Annu. Rev. Comput. Sci., vol. 1, pp. 225-253, June 1986.
- [7] A. Vikraman and R.A. Iannucci, "Two Fundamental Issues in Multiprocessing," in Proc. 4th Int'l DFVLR Semin. Found. Eng. Sci. Parallel Comput. Sci. Eng., 1988, pp. 61-88.
- [8] A. Vikraman, R.S. Nikhil, and K.K. Pingali, "I-Structures: Data Structures for Parallel Computing," *ACM Trans. Programm. Lang. Syst.*, vol. 11, no. 4, pp. 598-632, Oct. 1989.
  [9] P. Barahona and J.R. Gurd, "Simulated Performance of the
- [9] P. Barahona and J.R. Gurd, "Simulated Performance of the Manchester Multi-Ring Dataflow Machine," in *Proc. Parallel Comput.*, 1985, pp. 419-424.
- [10] P. Bellens, J. Perez, R. Badia, and J. Labarta, "CellSs: A Programming Model for the Cell BE Architecture," in Proc. ACM/IEEE Supercomput., 2006, pp. 1-5.

- [11] J. Benson, R. Cofell, C. Frericks, C.-H. Ho, V. Govindaraju, T. Nowatzki, and K. Sankaralingam, "Design, Integration and Implementation of the DySER Hardware Accelerator into OpenSPARC," in *Proc. IEEE Int'l Symp. High Perform. Comput. Architect.*, 2012, pp. 1-12.
- [12] E. Bloch, "The Engineering Design of the Stretch Computer," in *Proc. IRE-AIEE-ACM (Eastern) Comput. Conf.*, 1959, pp. 48-58.
  [13] W. Bohm, W. Najjar, B. Shankar, and L. Roh, "An Evaluation of
- [13] W. Bohm, W. Najjar, B. Shankar, and L. Roh, "An Evaluation of Coarse Grain Dataflow Code Generation Strategies," in *Proc. Programm. Models Massively Parallel Comput.*, 1993, pp. 63-71.
- [14] R. Buehrer and K. Ekanadham, "Incorporating Data Flow Ideas into Von Neumann Processors for Parallel Execution," *IEEE Trans. Comput.*, vol. 36, no. 12, pp. 1515-1522, Dec. 1987.
  [15] D. Burger, S. Keckler, K. McKinley, M. Dahlin, L.K. John, C. Lin,
- [15] D. Burger, S. Keckler, K. McKinley, M. Dahlin, L.K. John, C. Lin, C.R. Moore, J. Burrill, R.G. McDonald, and W. Yoder, the TRIPS Team, "Scaling to the End of Silicon with EDGE Architectures," *Computer*, vol. 37, no. 7, pp. 44-55, July 2004.
- [16] M. Cintra, J.F. Martínez, and J. Torrellas, "Architectural Support for Scalable Speculative Parallelization in Shared-Memory Multiprocessors," in Proc. Int'l Symp. Comput. Architect., 2000, pp. 13-24.
- pp. 13-24.
  [17] J. Clabes, J. Friedrich, M. Sweet, J. DiLullo, S. Chu, D. Plass, J. Dawson, P. Muench, L. Powell, M. Floyd, B. Sinharoy, M. Lee, M. Goulet, J. Wagoner, N. Schwartz, S. Runyon, G. Gorman, P. Restle, R. Kalla, J. McGill, and S. Dodson, "Design and Implementation of the POWER5 TM Microprocessor," in *Proc. Annu. Des. Autom. Conf.*, 2004, pp. 670-672.
- [18] K. Coons, X. Chen, S.K. Kushwaha, D. Burger, and K. McKinley, "A Spatial Path Scheduling Algorithm for EDGE Architectures," SIGPLAN Notices, vol. 41, no. 11, pp. 129-140, Nov. 2006.
- [19] T. M. Corp., "Connection Machine Model CM-2 Technical Summary," Thinking Machines Corp., Cambridge, MA, USA, Tech. Rep. TR89-1, 1989.
- [20] A. Cristal, O.J. Santana, F. Cazorla, M. Galluzzi, T. Ramirez, M. Pericas, and M. Valero, "Kilo-Instruction Processors: Overcoming the Memory Wall," *IEEE Micro*, vol. 25, no. 3, pp. 48-57, May/June 2005.
- [21] D.E. Culler, S.C. Goldstein, K.E. Schauser, and T. Eicken, "TAM: A Compiler Controlled Threaded Abstract Machine," J. Parallel Distrib. Comput., vol. 18, no. 3, pp. 347-370, July 1993.
- [22] D.E. Culler and G.M. Papadopoulos, "The Explicit Token Store," J. Parallel Distrib. Comput., vol. 10, no. 4, pp. 289-308, Dec. 1990.
- [23] D.E. Culler, K.E. Schauser, and T. Eicken, "Two Fundamental Limits on Dataflow Multiprocessing," in *Proc. IFIP WG 10.3 Conf. Architect. Compilation Tech. Medium and Fine Grain Parallelism*, 1993, pp. 153-164.
- [24] A.L. Davis and R. Keller, "Data Flow Program Graphs," Computer, vol. 15, no. 2, pp. 26-41, Feb. 1982.
  [25] L. Davis, "The Architecture and System Method of DDM1: A
- [25] L. Davis, "The Architecture and System Method of DDM1: A Recursively Structured Data Driven Machine," in Proc. Int'l Symp. Comput. Architect., 1978, pp. 210-215.
- [26] J.B. Dennis, "First Version of a Data Flow Procedure Language," in Proc. Programm. Symp., vol. 19, ser. Lecture Notes in Computer Science, B. Robinet, Ed., 1974, pp. 362-376, Springer-Verlag: Berlin, Germany.
- [27] J.B. Dennis and G.R. Gao, "Multithreaded Architectures: Principles, Projects, and Numbers," School Comput. Sci., McGill Univ., Montreal, QC, Canada, Tech. Rep., 1994.
- [28] J.B. Dennis and D.P. Misunas, "A Preliminary Architecture for a Basic Data-Flow Processor," in *Proc. Int'l Symp. Comput. Architect.*, 1975, pp. 126-132.
- [29] J.B. Dennis, "Data Flow Supercomputers," *Computer*, vol. 13, no. 11, pp. 48-56, Nov. 1980.
- [30] J.B. Dennis, "The Varieties of Data Flow Computers," Adv. Comput. Architect., pp. 51-60, 1986.
- [31] J.R. Ellis, "Bulldog: A Compiler for VLIW Architectures (Parallel Computing, Reduced-Instruction-Set, Trace Scheduling, Scientific)," Ph.D. dissertation, Yale Univ., New Haven, CT, USA, 1985.
- [32] Y. Etsion, F. Cabarcas, A. Rico, A. Ramirez, R.M. Badia, E. Ayguade, J. Labarta, and M. Valero, "Task Superscalar: An Out-of-Order Task Pipeline," in *Proc. IEEE/ACM Int'l Symp. Microarchitect.*, 2010, pp. 89-100.
  [33] Y. Etsion, A. Ramirez, R.M. Badia, E. Ayguade, J. Labarta, and
- [33] Y. Etsion, A. Ramirez, R.M. Badia, E. Ayguade, J. Labarta, and M. Valero, "Task Superscalar: Using Processors as Functional Units," in *Proc. Hot Topics Parallelism*, 2010, p. 16.

- [34] P. Evripidou and J.L. Gaudiot, "A Decoupled Graph/ Computation Data-Driven Architecture with Variable-Resolution Actors," in Proc. Int'l Conf. Parallel Process., 1990, pp. 405-414.
- [35] P. Evripidou and J.L. Gaudiot, "The USC Decoupled Multilevel Dataflow Execution Model," in Advanced Topics in Data-Flow Computing, J.-L. Gaudiot and L. Bic, Eds. Englewood Cliffs, NJ, USA: Prentice-Hall, 1991, pp. 347-379.
- [36] J.A. Fisher, "Very Long Instruction Word Architectures and the ELI-512," SIGARCH Comput. Architect. News, vol. 11, no. 3, pp. 140-150, June 1983.
- [37] M. Frank, C.A. Moritz, B. Greenwald, S. Amarasinghe, and A. Agarwal, "SUDS: Primitive Mechanisms for Memory Dependence Speculation," MIT, Cambridge, MA, USA, Tech. Rep., 1999.
- [38] M. Franklin and G.S. Sohi, "ARB: A Hardware Mechanism for Dynamic Reordering of Memory References," IEEE Trans. Comput., vol. 45, no. 5, pp. 552-571, May 1996.
- [39] J.-L. Gaudiot, T. DeBoni, J. Feo, W. Böhm, W. Najjar, and P. Miller, "The SISAL Model of Functional Programming and its Implementation," in Proc. Int'l Symp. Parallel Algorithms/Architect. Synth., 1997, pp. 112-123.
- [40] M. Gebhart, B.A. Maher, K.E. Coons, J. Diamond, P. Gratz, M. Marino, N. Ranganathan, B. Robatmili, A. Smith, J. Burrill, S.W. Keckler, D. Burger, and K.S. McKinley, "An Evaluation of the TRIPS Computer System," in Proc. Int'l Conf. Architect. Support Programm. Lang. Oper. Syst., 2009, pp. 1-12.
- [41] M. Gebhart, D.R. Johnson, D. Tarjan, S.W. Keckler, W.J. Dally, E. Lindholm, and K. Skadron, "Energy-Efficient Mechanisms for Managing Thread Context in Throughput Processors," in Proc. Int'l Symp. Comput. Architect., 2011, pp. 235-246. [42] E. Gluck-Hiltrop, M. Ramlow, and U. Schurfeld, "The Stollman
- Dataflow Machine," in Proc. Lect. Notes Comput. Sci., 1989, pp. 433-457
- [43] S. Gopal, T.N.V. James, E. Smith, and G.S. Sohi, "Speculative Versioning Cache," in Proc. Int'l Symp. High Perform. Comput. Architect., 1998, pp. 195-205.
- [44] V. Govindaraju, C. Ho, and K. Sankaralingam, "Dynamically Specialized Datapaths for Energy Efficient Computing," in Proc. IEEE Int'l Symp. High Perform. Comput. Architect., 2011, pp. 503-514.
- [45] V. Govindaraju, C.-H. Ho, T. Nowatzki, J. Chhugani, N. Satish, K. Sankaralingam, and C. Kim, "DySER: Unifying Functionality and Parallelism Specialization for Energy-Efficient Computing," IEEE Micro, vol. 32, no. 5, pp. 38-51, Sept./Oct. 2012
- [46] V.G. Grafe, G.S. Davidson, J.E. Hoch, and V. Holmes, "The Epsilon Dataflow Processor," in Proc. Int'l Symp. Comput. Architect., 1989, pp. 36-45.
- [47] V.G. Grafe and J. Hoch, "The EPSILON-2 Multiprocessor System," J. Parallel Distrib. Comput., vol. 10, no. 4, pp. 309-318, Dec. 1990.
- [48] P. Gratz, C. Kim, R. McDonald, S.W. Keckler, and D. Burger, "Implementation and Evaluation of On-Chip Network Architectures," in Proc. Int'l Conf. Comput. Des., 2006, pp. 477-484.
- [49] M. Gupta and R. Nim, "Techniques for Speculative Run-Time Parallelization of Loops," in Proc. IEEE/ACM Supercomput., 1998, pp. 1-15.
- [50] J.R. Gurd, C.C. Kirkham, and I. Watson, "The Manchester Prototype Dataflow Computer," Commun. ACM, vol. 28, no. 1, pp. 34-52, Jan. 1985.
- [51] R. Hameed, W. Qadeer, M. Wachs, O. Azizi, A. Solomatnikov, B.C. Lee, S. Richardson, C. Kozyrakis, and M. Horowitz, "Understanding Sources of Inefficiency in General-Purpose Chips," in Proc. Int'l Symp. Comput. Architect., 2010, pp. 37-47.
- [52] L. Hammond, M. Willey, and K. Olukotun, "Data Speculation Support for a Chip Multiprocessor," in Proc. Int'l Conf. Architect. Support Programm. Lang. Oper. Syst., 1998, pp. 58-69. [53] M. Herlihy and J.E.B. Moss, "Transactional Memory: Architec-
- tural Support for Lock-Free Data Structures," in Proc. Int'l Symp. Comput. Architect., 1993, pp. 289-300. [54] J. Hicks, D. Chiou, B.S. Ang, and A. Vikraman, "Performance
- Studies of Id on the Monsoon Dataflow System," J. Parallel Distrib. Comput., vol. 18, no. 3, pp. 273-300, July 1993.
- [55] D. Hillis, "The Connection Machine," Ph.D. dissertation, Dept. Elect. Eng. Comput. Sci., MIT, Cambridge, MA, USA, 1988.
- [56] S. Hong and H. Kim, "An Integrated GPU Power and Performance Model," in Proc. Int'l Symp. Comput. Architect., 2010, pp. 280-289.

- [57] H.H.J. Hum, O. Maquelin, K.B. Theobald, X. Tian, X. Tang, G.R. Gao, P. Cupryky, N. Elmasri, L.J. Hendren, A. Jimenez, S. Krishnany, A. Marquez, S. Merali, S.S. Nemawarkarz, P. Panangaden, X. Xue, and Y. Zhu, "A Design Study of the EARTH Multiprocessor," in *Proc. Int'l Conf. Parallel Architect*. Compilation Tech., 1995, pp. 59-68.
- [58] H.H.J. Hum, O. Maquelin, K. Theobald, X. Tian, G. Gao, and L. Hendren, "A Study of the EARTH-MANNA Multithreaded System," Parallel Programm., vol. 24, no. 4, pp. 319-348, Aug. 1996.
- [59] W. Hwu and Y.N. Patt, "HPSm, A High Performance Restricted Data Flow Architecture Having Minimal Functionality," in Proc. Int'l Symp. Comput. Architect., 1986, pp. 297-306.
- [60] R.A. Iannucci, "Toward a Dataflow/Von Neumann Hybrid Architecture," in Proc. Int'l Symp. Comput. Architect., 1988, pp. 131-140.
- [61] R.A. Iannucci, G.R. Gao, R.H. Halstead Jr., and B. Smith, Multithreaded Computer Architecture: A Summary of the State of the Art. Boston, MA, USA: Kluwer, 1994.
- [62] N. Ito, M. Sato, E. Kuno, and K. Rokusawa, "The Architecture and Preliminary Evaluation Results of the Experimental Parallel Inference Machine PIM-D," in Proc. Int'l Symp. Comput. Architect., 1986, pp. 149-156.
- [63] H.F. Jordan, "Performance Measurements on HEP-A Pipelined MIMD Computer," in Proc. Int'l Symp. Comput. Architect., 1983, pp. 207-212.
- [64] G. Kahn, "The Semantics of a Simple Language for Parallel Programming," in Proc. IFIP Congr., 1974, vol. 74, pp. 471-475.
- [65] R.M. Karp and R.E. Miller, "Properties of a Model for Parallel Computations: Determinacy, Termination, Queueing," SIAM J. Appl. Math., vol. 14, no. 5, pp. 1390-1411, 1966.
- [66] K.M. Kavi, B. Buckles, and U. Bhat, "A Formal Definition of Data Flow Graph Models," IEEE Trans. Comput., vol. C-35, no. 11, pp. 940-948, Nov. 1986.
- [67] K.M. Kavi, R. Giorgi, and J. Arul, "Scheduled Dataflow: Execution Paradigm, Architecture, and Performance Evaluation," IEEE Trans. Comput., vol. 50, no. 8, pp. 834-846, Aug. 2001.
- [68] C. Kim and J.L. Gaudiot, Dataflow and Multithreaded Architectures. Hoboken, NJ, USA: Wiley, 1997.
- C. Kim, S. Sethumadhavan, M.S. Govindan, N. Ranganathan, [69] D. Gulati, D. Burger, and S.W. Keckler, "Composable Lightweight Processors," in Proc. IEEE/ACM Int'l Symp. Microarchitect., 2007, pp. 381-394.
- [70] M. Kishi, H. Yasuhara, and Y. Kawamura, "DDDP-A Distributed Data Driven Processor," in Proc. Int'l Symp. Comput. Architect., 1983, pp. 236-242.
- [71] Y. Kodama, H. Sakane, M. Sato, H. Yamana, S. Sakai, and Y. Yamaguchi, "The EM-X Parallel Computer: Architecture and Basic Performance," in Proc. Int'l Symp. Comput. Architect., 1995, pp. 14-23. [72] V. Krishnan and J. Torrellas, "A Chip-Multiprocessor Architec-
- ture with Speculative Multithreading," IEEE Trans. Comput., vol. 48, no. 9, pp. 866-880, Sept. 1999.
- [73] V. Krishnan and L.J. Torrellas, "The Need for Fast Communication in Hardware-Based Speculative Chip Multiprocessors," in Proc. Int'l Conf. Parallel Architect. Compilation Tech., 1999, pp. 24-33.
- [74] D. Kuck, E. Davidson, D. Lawrie, A. Sameh, C.Q. Zhu, A. Veidenbaum, J. Konicek, P. Yew, K. Gallivan, W. Jalby, H. Wijshoff, R. Bramley, U.M. Yang, P. Emrath, D. Padua, R. Eigenmann, J. Hoeflinger, G. Jaxon, Z. Li, T. Murphy, and J. Andrews, "The Cedar System and an Initial Performance Study," in Proc. Int'l Symp. Comput. Architect., 1993, pp. 213-223.
- [75] K. Kurihara, D. Chaiken, and A. Agarwal, "Latency Tolerance Through Multithreading in Large-Scale Multiprocessors," in Proc. Int'l Symp. Comput. Architect., 1991, pp. 91-101.
- [76] C. Kyriacou, P. Evripidou, and P. Trancoso, "Data-Driven Multithreading Using Conventional Microprocessors," IEEE Trans. Parallel Distrib. Syst., vol. 17, no. 10, pp. 1176-1188, Oct. 2006.
- [77] J. McGraw, S. Skedzielewski, S. Allan, R. Oldehoeft, J. Glauert, C. Kirkham, B. Noyce, and R. Thomas, SISAL: Streams and Iteration in a Single Assignment Language. Language Reference Manual. Livermore, CA, USA: Livermore National Laboratory, 1985.
- [78] B. Lee and A. Hurson, "Dataflow Architectures and Multithreading," *Computer*, vol. 27, no. 8, pp. 27-39, Aug. 1994. [79] E.A. Lee and D.G. Messerschmitt, "Synchronous Data Flow,"
- Proc. IEEE, vol. 75, no. 9, pp. 1235-1245, Sept. 1987.

- [80] P. Marcuello and A. González, "Clustered Speculative Multithreaded Processors," in Proc. Int'l Symp. Supercomput., 1999, pp. 365-372.
- [81] D.T. Marr, F. Binns, D.L. Hill, G. Hinton, D.A. Koufaty, J.A. Miller, and M. Upton, "Hyper-Threading Technology Architecture and Microarchitecture," Intel Technol. J., vol. 6, no. 1, pp. 1-12, Feb. 2002.
- [82] L.A.J. Marzulo, F.M.G. Franca, and V.S. Costa, "Transactional WaveCache: Towards Speculative and Out-of-Order Dataflow Execution of Memory Operations," in Proc. Int'l Symp. Comput. Architect. High Perform. Comput., 2008, pp. 183-190. [83] W.M. Miller, W.A. Najjar, and A.P.W. Bohm, "A Quantitative
- Analysis of Locality in Dataflow Programs," in Proc. Int'l Symp. Microarchitect., 1991, pp. 12-18.
- [84] M. Mishra, T.J. Callahan, T. Chelcea, G. Venkataramani, S.C. Goldstein, and M. Budiu, "Tartan: Evaluating Spatial Computation for Whole Program Execution," in Proc. Int'l Conf. Architect. Support Programm. Lang. Oper. Syst., 2006, pp. 163-174. [85] W.A. Najjar, E.A. Lee, and G.R. Gao, "Advances in the Dataflow
- Computational Model," Parallel Comput., vol. 25, no. 13/14, pp. 1907-1929, Dec. 1999.
- [86] S.S. Nemawarkar and G.R. Gao, "Measurement and Modeling of EARTH-MANNA Multithreaded Architecture," in Proc. Int'l Workshop Model., Anal., Simul. Comput. Telecommun. Syst., 1996, pp. 109-104.
- [87] J. Nickolls, I. Buck, M. Garland, and K. Skadron, "Scalable Parallel Programming with CUDA," ACM Queue, vol. 6, no. 2, pp. 40-53, Mar./Apr. 2008.
- [88] R.S. Nikhil, "Can Dataflow Subsume Von Neumann Computing?" in Proc. Int'l Symp. Comput. Architect., 1989, pp. 262-272.
- [89] R.S. Nikhil, G.M. Papadopoulos, and A. Vikraman, "\*T: A Multithreaded Massively Parallel Architecture," in Proc. Int'l Symp. Comput. Architect., 1992, pp. 156-167.
- [90] G.M. Papadopoulos and D.E. Culler, "Monsoon: An Explicit Token-Store Architecture," in Proc. Int'l Symp. Comput. Architect., 1990, pp. 82-91.
- [91] G.M. Papadopoulos and K.R. Traub, "Multithreading: A Revisionist View of Dataflow Architectures," in Proc. Int'l Symp. Comput. Architect., 1991, pp. 342-351.
- [92] Y.N. Patt, W.M. Hwu, and M. Shebanow, "HPS, A New Microarchitecture: Rationale and Introduction," in Proc. Int'l Symp. Microarchitect., 1985, pp. 103-108.
- [93] S. Pei, B. Wu, M. Du, G. Chen, L.A.J. Marzulo, and F.M.G. Franca, "SpMT WaveCache: Exploiting Thread-Level Parallelism in Wavescalar," in Proc. Congr. Comput. Sci. Inf. Eng., 2009, pp. 530-535.
- [94] J. Perez, R. Badia, and J. Labarta, "A Dependency-Aware Task-Based Programming Environment for Multi-Core Architectures," in Proc. IEEE Int'l Conf. Cluster Comput., 2008, pp. 142-151
- [95] A. Petersen, A. Putnam, M. Mercaldi, A. Schwerin, S. Eggers, S. Swanson, and M. Oskin, "Reducing Control Overhead in Dataflow Architectures," in Proc. Int'l Conf. Parallel Architect. Compilation Tech., 2006, pp. 182-191.
- [96] A. Plas, D. Comte, O. Gelly, and J. Syre, "LAU System Architecture: A Parallel Data-Driven Processor Based on Single Assignment," in Proc. Int'l Conf. Parallel Process., 1976, pp. 293-302.
- [97] M. Prvulovic, M.J. Garzarán, L. Rauchwerger, and J. Torrellas, "Removing Architectural Bottlenecks to the Scalability of Speculative Parallelization," in Proc. Int'l Symp. Comput. Architect., 2001, pp. 204-215.
- [98] L. Rauchwerger and D. Padua, "The LRPD Test: Speculative Run-Time Parallelization of Loops with Privatization and Reduction Parallelization," in Conf. on Programming Language Design and Implementation, 1995, pp. 218-232.
- [99] J.E. Řequa, "The Piecewise Data Flow Architecture Control Flow and Register Management," in Proc. Int'l Symp. Comput. Architect., 1983, pp. 84-89.
- [100] B. Robatmili, K.E. Coons, D. Burger, and K.S. McKinley, "Strategies for Mapping Dataflow Blocks to Distributed Hardware," in Proc. IEEE/ACM Int'l Symp. Microarchitect., 2008, pp. 23-34.
- [101] B. Robic, J. Silc, and T. Ungerer, "Beyond Dataflow," Comput. Inf.
- Technol., vol. 8, no. 2, pp. 89-101, 2000.
   [102] L. Roh and W. Najjar, "Design of Storage Hierarchy in Multi-threaded Architectures," in Proc. Int'l Symp. Microarchitect., 1995, pp. 271-278.

- [103] E. Rotenberg, Q. Jacobson, Y. Sazeides, and J.E. Smith, "Trace Processors," in Proc. IEEE/ACM Int'l Symp. Microarchitect., 1997, pp. 138-148. [104] P. Rundberg and P. Stenstrom, "Low-Cost Thread-Level Data
- Dependence Speculation on Multiprocessors," in Proc. 4th Workshop Multithreaded Execution, Architect. Compilation, 2000, pp. 1-9.
- [105] K. Sankaralingam, R. Nagarajan, H. Liu, C. Kim, J. Huh, D. Burger, S.W. Keckler, and C.R. Moore, "Exploiting ILP, TLP, and DLP with the Polymorphous TRIPS Architecture," in Proc. Int'l Symp. Comput. Architect., 2003, pp. 422-433.
- [106] K. Sankaralingam, R. Nagarajan, H. Liu, C. Kim, J. Huh, N. Ranganathan, D. Burger, S.W. Keckler, R.G. Mcdonald, and C.R. Moore, "TRIPS: A Polymorphous Architecture for Exploiting ILP, TLP, and DLP," ACM Trans. Architect. Code Optim., vol. 1, no. 1, pp. 62-93, Mar. 2004.
- [107] K. Sankaralingam, R. Nagarajan, R. McDonald, R. Desikan, S. Drolia, M.S. Govindan, P. Gratz, D. Gulati, H. Hanson, C. Kim, H. Liu, N. Ranganathan, S. Sethumadhavan, S. Sharif, P. Shivakumar, S.W. Keckler, and D. Burger, "Distributed Microarchitectural Protocols in the TRIPS Prototype Processor," in Proc. IEEE/ACM Int'l Symp. Microarchitect., 2006, pp. 480-491.
- [108] T. Sherwood, S. Sair, and B. Calder, "Predictor-Directed Stream Buffers," in Proc. IEEE/ACM Int'l Symp. Microarchitect., 2000, pp. 42-53.
- [109] J. Silc, B. Robic, and T. Ungerer, "Asynchrony in Parallel Computing: From Dataflow to Multithreading,' ' J. Parallel Distrib. Comput., vol. 1, no. 1, pp. 1-33, 1998.
- [110] J. Silc, B. Robic, and T. Ungerer, Processor Architecture: From Dataflow to Superscalar and Beyond. New York, NY, USA: Springer-Verlag, 1999.
- [111] A. Smith, J. Burrill, J. Gibson, B. Maher, N. Nethercote, B. Yoder, D.C. Burger, and K.S. McKinley, "Compiling for EDGE Architectures," in Proc. Int'l Symp. Code Gen. Optim., 2006, op. 185-195.
- [112] B.J. Smith, "Architecture and Applications of the HEP Multiprocessor Computer System," in Proc. SPIE Real Time Signal Process. IV, 1981, pp. 241-248.
- [113] J.E. Smith and A.R. Pleszkun, "Implementation of Precise Interrupts in Pipelined Processors," in Proc. Int'l Symp. Comput. Architect., 1998, pp. 291-299.
- [114] G.S. Sohi, S.E. Breach, and T.N. Vijaykumar, "Multiscalar Processors," in Proc. Int'l Symp. Comput. Architect., 1995, pp. 414-425. [115] V.P. Srini, "An Architectural Comparison of Dataflow Systems,"
- Computer, vol. 19, no. 3, pp. 68-88, Mar. 1986.
- [116] J.G. Steffan, "Hardware Support for Thread-Level Speculation," Ph.D. dissertation, Dept. Comput. Sci., Carnegie Mellon Univ., Pittsburgh, PA, USA, 2003.
- [117] J. Strohschneider, B. Klauer, S. Zickenheimer, and K. Waldschmidt, "ADARK: A Fine Grain Dataflow Architecture with Associative Communication Network," in Proc. EUROMICRO Conf., 1994, pp. 445-450.
- [118] S. Swanson, K. Michelson, A. Schwerin, and M. Oskin, "WaveScalar," in Proc. Int'l Symp. Microarchitect., 2003, pp. 291-302.
- [119] S. Swanson, A. Putnam, M.M. Mercaldi, K. Michelson, A. Petersen, A. Schwerin, M. Oskin, and S.J. Eggers, "Area-Performance Trade-Offs in Tiled Dataflow Architectures," in *Proc. Int'l Symp.* Comput. Architect., 2006, pp. 314-326.
- [120] S. Swanson, A. Schwerin, M. Mercaldi, A. Petersen, A. Putnam, K. Michelson, M. Oskin, and S. Eggers, "The WaveScalar Architecture," ACM Trans. Comput. Syst., vol. 25, no. 2, pp. 4:1-4:54, May 2007.
- [121] K.B. Theobald, "EARTH: An Efficient Architecture for Running Threads," Ph.D. dissertation, McGill Univ., Montreal, Quebec, CA, 1999.
- [122] X.-M. Tian, S. Nemawarkar, G.R. Gao, H. Hum, O. Maquelin, A. Sodan, and K. Theobald, "Quantitive Studies of Data-Locality Sensitivity on the EARTH Multithreaded Architecture: Preliminary Results," in Proc. Int'l Conf. High-Perform. Comput., 1996, pp. 362-367.
- [123] P. Treleaven, R. Hopkins, and P. Rautenbach, "Combining Data Flow and Control Flow Computing," Comput. J., vol. 25, no. 2, pp. 207-217, 1982.

- [124] D.M. Tullsen, S.J. Eggers, and H.M. Levy, "Simultaneous Multithreading: Maximizing On-Chip Parallelism," in Proc. Int'l Symp. Comput. Architect., 1995, pp. 392-403.
- [125] R. Vedder and D. Finn, "The Hughes Data Flow Multiprocessor: Architecture for Efficient Signal and Data Processing," in Proc. Int'l Symp. Comput. Architect., 1985, pp. 324-332.
- [126] A.H. Veen, "Dataflow Machine Architecture," ACM Comput. Surveys, vol. 18, no. 4, pp. 365-396, Dec. 1986.
- [127] G. Venkatesh, J. Sampson, N. Goulding, S. Garcia, V. Bryksin, J. Lugo-Martinez, S. Swanson, and M.B. Taylor, "Conservation Cores: Reducing the Energy of Mature Computations," in *Proc. Int'l Conf. Architect. Support Programm. Lang. Oper. Syst.*, 2010, pp. 205-218.
  [128] J. von Neumann, "First Draft of a Report on the EDVAC," U.S.
- [128] J. von Neumann, "First Draft of a Report on the EDVAC," U.S. Army Ordnance Dept. Univ. Pennsylvania Moore, School Elect. Eng., Philadelphia, PA, USA, Tech. Rep., 1945.
- [129] W.-D. Weber and A. Gupta, "Exploring the Benefits of Multiple Hardware Contexts in a Multiprocessor Architecture: Preliminary Results," in *Proc. Int'l Symp. Comput. Architect.*, 1989, pp. 273-280.
- [130] W. Yamamoto, M.J. Serrano, A.R. Talcott, R.C. Wood, and M. Nemirovsky, "Performance Estimation of Multistreamed, Supersealar Processors," in *Proc. Hawaii Int'l Conf. Syst. Sci.*, 1994, pp. 195-204.
- [131] F. Yazdanpanah, D. Jimenez-Gonzalez, C. Alvarez-Martinez, Y. Etsion, and R.M. Badia, "FPGA-Based Prototype of the Task Superscalar Architecture," in *Proc. HiPEAC Workshop Reconfi*gurable Comput., 2013, pp. 1-10.
- [132] T. Yuba, K. Hiraki, T. Shimada, S. Sekiguchi, and K. Nishida, "The SIGMA-1 Dataflow Computer," in Proc. Comput. Conf. Explor. Technol., Today Tomorrow, 1987, pp. 578-585.
- [133] E. Zehender and T. Ungerer, "The ASTOR Architecture," in Proc. Int'l Conf. Distrib. Comput. Syst., 1987, pp. 424-430.
- [134] Y. Zhang, L. Rauchwerger, and J. Torrellas, "Hardware for Speculative Parallelization of Partially-Parallel Loops in DSM Multiprocessors," in Proc. Int'l Symp. High Perform. Comput. Architect., 1999, pp. 135-139.



Fahimeh Yazdanpanah received the MSc degree from the University of Isfahan, Isfahan, Iran, in 2003. She is a PhD student at the Technical University of Catalunya (UPC), Catalonia, Spain, and researches at the Barcelona Supercomputing Center (BSC-CNS), Spain. Her research interests include computer architecture, digital and hardware design, VHDL, reconfigurable architecture, and FPGA-based design. She has been participating in the TERAFLUX European project.



**Carlos Alvarez-Martinez** received the MS and PhD degrees in computer science from the Technical University of Catalunya (UPC), Catalonia, Spain, in 1998 and 2007, respectively. Since 1998, he has been lecturing on computer design, computer architecture, and open source software. He has been a tenured Assistant Professor in the Computer Architecture Department at UPC since 2009. His research interests cover the areas of architectures, runtime systems, and reconfigurable solutions for high-performance

multiprocessor systems. Dr. Alvarez has coauthored more than 30 publications in international journals and conferences. Currently, he is advising two PhD students. He has been participating in the Hipeac Network of Excellence and in the TERAFLUX European project.



Daniel Jimenez-Gonzalez received the MS and PhD degrees in computer science from the Technical University of Catalunya (UPC), Catalonia, Spain, in 1997 and 2004, respectively. Since 2000, he has been lecturing on computer design, parallel, and architecture awareness programming. He has been a Collaborating Lecturer in the Computer Architecture Department at UPC since 2005. His research interests cover the areas of parallel and architecture awareness programming for bioinformatics ap-

plications, runtime systems, compilers, and reconfigurable solutions for high-performance multiprocessor systems. Dr. Jimenez-Gonzalez has coauthored more than 30 publications in international journals and conferences. Currently, he is co-advising two PhD students. He has been participating in the Hipeac Network of Excellence and in the SARC, ACOTES, TERAFLUX, and PRACE European projects.



Yoav Etsion received the MSc and PhD from the Hebrew University of Jerusalem, Jerusalem, Israel, in 2003 and 2009. He is an Assistant Professor at the Electrical Engineering and Computer Science departments in Technion-Israel Institute of Technology, Haifa, Israel. He was a Senior Researcher at the Barcelona Supercomputing Center (BSC-CNS) where he held a Juan de la Cierva Fellowship from the Ministry of Science and innovation of Spain. His research interests include computer architec-

ture, HW/SW interoperability, operating systems, and parallel programming models. He is a member of the IEEE.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.