

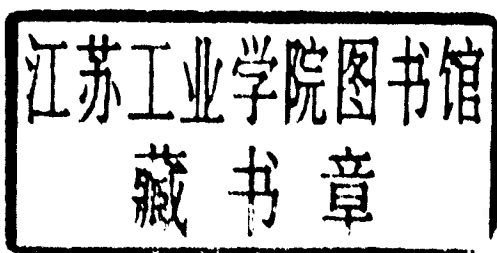
GREGORY M. PAPADOPOULOS

Implementation of a General-Purpose Dataflow Multiprocessor

Gregory M. Papadopoulos

MIT Laboratory for Computer Science

Implementation of a General-Purpose Dataflow Multiprocessor



Pitman, London

The MIT Press, Cambridge, Massachusetts

PITMAN PUBLISHING
128 Long Acre, London WC2E 9AN

© Gregory M. Papadopoulos 1991

First published 1991

Available in the Western Hemisphere and Israel from
The MIT Press
Cambridge, Massachusetts (and London, England)

ISSN 0953-7767

British Library Cataloguing in Publication Data

Papadopoulos, Gregory M.

Implementation of a general-purpose dataflow
multiprocessor.—(Research monographs in parallel and
distributed computing : ISSN 0953-7767

1. Computer systems. Multiprocessors

I. Title II. Series

004.35

ISBN 0-273-08835-1

Library of Congress Cataloging-in-Publication Data

Papadopoulos, Gregory Michael.

Implementation of a general purpose dataflow multiprocessor /
Gregory Michael Papadopoulos.—1st MIT Press ed.

p. cm.—(Research monographs in parallel and distributed
computing)

ISBN 0-262-66069-5 (pbk.)

1. Multiprocessors. 2. Computer architecture. 3. Parallel
processing (Electronic computers) I. Title. II. Series.

QA76.5.P29 1990

004'.35—dc20

All rights reserved; no part of this publication may be reproduced,
stored in a retrieval system, or transmitted in any form or by any
means, electronic, mechanical, photocopying, recording or
otherwise without the prior written permission of the publishers or
a licence permitting restricted copying in the United Kingdom
issued by the Copyright Licencing Agency Ltd, 33–34 Alfred Place,
London WC1E 7DP. This book may not be lent, resold, hired out or
otherwise disposed of by way of trade in any form of binding or cover
other than that in which it is published, without the prior consent of
the publishers.

Reproduced and printed by photolithography
in Great Britain by Biddles Ltd, Guildford

Contents

1	General Purpose Multiprocessing	1
1.1	Why Fine Grain Synchronization	2
1.1.1	Tasks Should be Cheap	3
1.1.2	A Virtual Memory Analogy	4
1.1.3	An Informal Task Model of Parallel Computation	6
1.1.4	Parallelism and Synchronization	8
1.1.5	The Costs of Task State Transitions	9
1.1.6	Dataflow Machines Directly Execute a “Reduced” Task Graph . . .	11
1.2	Roadmap	12
2	The Tagged-Token Dataflow Architecture	14
2.1	A TTDA Primer	14
2.1.1	A Simple Example	15
2.1.2	Tags Distinguish Tokens from Different Activations	18
2.2	Tag Efficiency	19
2.2.1	Recycling Tags	19
2.2.2	Iteration Optimization	20
2.2.3	Context Registers	21
2.3	The TTDA Abstract Pipeline	21
2.3.1	Storage and Names in the TTDA	22
2.3.2	The Waiting-Matching Problem	23
3	The Explicit Token Store	25
3.1	Storage, Tokens and Instructions	29
3.2	Code Blocks	31
3.3	Transformation of Tokens	33
3.4	Activity Generation	34
3.4.1	Dyadic Operators	34
3.4.2	Monadic Operators	35
3.4.3	Dyadic Operators with a Constant Operand	35
3.5	Token Generation	36
3.5.1	The Arithmetic Rule	37
3.5.2	The Send Rule	38
3.5.3	The Extract Rule	38
3.5.4	Combining Rules	39
3.6	The Heap	39
3.7	Parallelism in ETS	41
3.8	Summary	42

4	Compiling for an ETS Dataflow Processor	43
4.1	Basic Instruction Set Equivalences	44
4.1.1	Machine Data Types	44
4.1.2	TTDA Instructions	45
4.1.3	Rewriting Three-Input Instructions	46
4.1.4	Operand Matching Rules	46
4.1.5	Rewriting Instructions with Three or More Destinations	48
4.1.6	Instruction Opcode Classes	49
4.2	Loops	55
4.3	I-structures: Descriptors and Multiple Readers	58
4.3.1	Descriptors	59
4.3.2	Multiple Deferred Reads	61
4.4	Closures	63
4.5	Resource Managers	64
4.6	Operators with More than Two Inputs	66
4.7	Slot Allocation Revisited	67
4.8	Summary	70
5	Compiling Imperative Languages for an ETS	71
5.1	Threads	71
5.2	Translating Quads into Thread Sequences	74
5.3	Multiple Threads Within an Activation	77
5.4	Summary	77
6	Monsoon: An ETS Multiprocessor	79
6.1	Mapping Computation onto Multiple Processing Elements	79
6.1.1	Mapping Activation Frames	81
6.1.2	Mapping Data Structures	82
6.1.3	Mapping Code	83
6.1.4	Mapping Unprocessed Tokens	84
6.2	Processing Element Microarchitecture	85
6.2.1	Tokens	86
6.2.2	Requests	91
6.2.3	Instructions	91
6.2.4	Temporary Registers	92
6.2.5	Exceptions	94
6.2.6	Detailed Pipeline Operation	96
6.3	I-Structure Memory Elements	98
6.4	The Network	98
6.4.1	Bandwidth Requirements	99
6.4.2	The Write-Acknowledgment Problem	99
6.5	Evaluation	101

7	The Monsoon Macroarchitecture	105
7.1	Overview	105
7.2	Data	106
7.2.1	The Data Field	108
7.2.2	The Type Field	109
7.3	Programming Examples	110
7.3.1	A Simple Expression	110
7.3.2	Forking and Joining Threads	111
7.3.3	Combining Instructions	112
7.3.4	Split Phase	114
7.3.5	Procedure Call	115
7.3.6	Conditional Branch	119
7.3.7	Exceptions	119
7.4	Instruction Set Summary	120
8	Conclusion	121
A	Monsoon Instruction Decoding	124
A.1	First Level Decode	126
A.1.1	Type Map	126
A.1.2	Presence Map	127
A.1.3	Frame Store Operation	127
A.2	Second Level Decode	128
A.2.1	Function Unit Control	129
A.2.2	Floating point, arithmetic and logic unit	130
A.2.3	Pointer increment unit	130
A.2.4	Type Propagation Unit	131
A.2.5	Machine control unit	132
A.2.6	Next Address Control	134
A.2.7	Form Token and Token Queues	135
A.2.8	Form Token Control	135
A.3	Exceptions and Condition Codes	139
A.4	Statistics	141
A.5	Changes from the Monsoon Prototype	141
B	Monsoon Macroinstruction Set	143
B.1	Instruction Set Summary	143
B.1.1	Arithmetic	143
B.1.2	Branch	147
B.1.3	Supervisor Call	149
B.1.4	Split-Phase	150
B.2	Detailed Macroinstruction Encoding	151
B.2.1	Macroinstruction Fields	151
B.2.2	Detailed Encodings for Selected Instructions	153

List of Figures

1.1	The Parallelism–Overhead Tradeoff (from Sarkar)	4
1.2	Speedup vs. Task Size When Task Overhead is Very Low	5
1.3	Simultaneous Tasks Referencing a Shared Object	7
1.4	A Typical Task State Transition Diagram	10
1.5	Task State Transition of A Dataflow Instruction	12
2.1	A Dataflow Graph for Computing Vector Inner Product (From [14])	16
2.2	Dataflow Graph for “ $s + A[j] * B[j]$ ”	16
2.3	A Firing Sequence for “ $s + A[j] * B[j]$ ” (From [14])	17
2.4	The TTDA Token Processing Pipeline	22
3.1	Explicit Matching Operation	26
3.2	Compilation of a Simple Expression for an ETS	27
3.3	An Example ETS Pipeline	28
3.4	ETS Instruction Components	30
3.5	Schematic Relationship Between Storage, Tokens and Instructions	31
3.6	An ETS Code Block	32
3.7	Relationship Between Activity Generation and Token Generation Phases .	33
3.8	The Matching Function for a Dyadic Operator	35
3.9	The <i>sticky</i> Matching Function for a Dyadic Operator with a Constant Input	37
3.10	The Arithmetic, Extract and Send Token Forming Rules	37
3.11	A Single Instruction Code Block That Performs Reads and Writes on Any Location	40
3.12	A Two Phase Read of a Location	41
4.1	TTDA Instruction Components	46
4.2	Rewriting Three-Input TTDA Instructions into Two-Input Instructions . .	47
4.3	The Four Basic Input Operand Forms	48
4.4	A Fanout Tree for an Instruction with Five Destinations	49
4.5	Executing a Loop of n Iterations as a Tail Recursion	56
4.6	Executing a k -Bounded Loop on an ETS	56
4.7	Block Diagram of an ETS Loop Schema	57
4.8	A Set of Contexts Forming a $k = 3$ Bounded Loop	58
4.9	Example of an ETS I-Structure Descriptor Convention	59
4.10	Rewriting <code>form-address</code> to Account for a Non-Zero Lower Bound	60
4.11	Translation of <code>upper-bound</code> and <code>lower-bound</code>	61
4.12	Augmenting an <code>i-fetch</code> to Support Multiple Deferred Readers	62
4.13	A Deferred Read List Comprising Two Readers	62

4.14	The <i>istr</i> Matching Function for an I-Structure Slot which Supports Deferred Readers	63
4.15	A Closure with Two Arguments Applied	64
4.16	Example Operation of the cs-gate Instruction	65
4.17	Enqueueing Manager Requests with q-gate	66
4.18	A gate Instruction With Three Triggers	67
4.19	A gate Instruction which can be Short-Circuited	68
4.20	Composition of Two Instructions to Make Short-Circuited gate	68
4.21	A) Code Block and B) Corresponding Double Dependence Graph	69
5.1	Viewing an ETS Token as a Sequential Thread Descriptor	72
6.1	Top Level View of the Monsoon Multiprocessor	80
6.2	Partitioning Storage Across Processing Elements	81
6.3	Data Structure Interleaving as a Function on n	83
6.4	Eight Stage Non-Blocking Monsoon Pipeline	87
6.5	Example Use of Temporary Registers for the Expression $(x + y) * (x - y)$	94
6.6	A Race Between an i-store and a Deallocate	100
6.7	Monsoon Processor Board	102
7.1	A Computation Descriptor (CD)	107
7.2	A Sample Monsoon Macrocode Fragment	110
7.3	A Monsoon Macrocode Fragment Using a Temporary Register	111
7.4	Example of Thread Fork and Join	112
7.5	An Example of Optimizing Code by Combining Instructions	113
7.6	Macrocode for the expression $z = (x + y) * (x - y)$	113
7.7	Example of Split-Phase Fetch and Store	115
7.8	Non-strict Procedure Call Convention	117
7.9	An Optimized Version of Non-Strict Procedure Call	118
7.10	An Example Using Conditional Branch	119
7.11	An Example Using Boolean Predicates	119
8.1	The Cost-Performance of Uniprocessors and the Multiprocessor Promise	121
A.1	Instruction Decoding Tables and Maps	125
B.1	Integer Arithmetic Instructions	143
B.2	Integer Bit Manipulation Instructions	144
B.3	Floating Point Arithmetic	144
B.4	Predicate Instructions	145
B.5	Conversion Instructions	146
B.6	Miscellaneous Arithmetic Instruction	147
B.7	Integer Conditional Jump Instructions	148
B.8	Floating Point Conditional Jump Instructions	149
B.9	Conditional Switch Instructions	149
B.10	Unconditional Switch Instructions	149
B.11	Software Exception Instructions	150

B.12 I-Structure Global Memory Instructions	150
B.13 Imperative Global Memory Instructions	151
B.14 Instruction Memory Global Memory Instructions	151
B.15 Generic Global Memory Instructions	151

*To my parents,
Imogen and Michael Papadopoulos*

1 General Purpose Multiprocessing

Across the diverse range of multiprocessor architectures, from small collections of supercomputers to thousands of synchronous single-bit processors, all seem to share one undesirable property: they are hard to use. Programming has taken a giant step *backwards*. The application writer must consider detailed and hard-to-measure interactions of the program with the machine; the resulting codes are difficult to debug [37], are of questionable reliability [46], and are far from portable [31]. Even after arduous work in converting an application for parallel execution, the actual improvement in performance is frequently disappointing. Perhaps we do not yet understand how to express parallelism in a way that is machine independent. Perhaps we need more sophisticated compilers and associated debuggers to better exploit a machine's parallelism. Perhaps we are building the **wrong** machines.

It is a pervasive belief that our lack of real success in general purpose multiprocessing is a software problem. Machine architects adapt sequential processors to the parallel setting by providing an interprocessor communication medium and an *ad hoc* set of synchronization mechanisms, like locks or fetch-and-add. Then the compiler, or worse yet the programmer, is expected to partition the application into tasks that can run in parallel using the supplied synchronization primitives to ensure deterministic behavior. While we certainly share the belief that there is a software problem, we are convinced that there are equally serious defects in the underlying machines. There needs to be a fundamental change in processor architecture before we can expect significant progress to be made in the use of multiprocessors. The nature of this change is the deep integration of synchronization into the processor instruction set [12]. The instruction set must constitute a **parallel machine language**, in which parallel activities are coordinated as efficiently as instructions are scheduled.

The instruction set of a dataflow machine [6] forms such a parallel machine language. An instruction is executed only when all of its required operands are available. Thus, low-level synchronization is performed for every instruction and at the same rate as instructions are issued. It is easy for a compiler to use this fine grain synchronization to produce code which is highly parallel but deterministic for any runtime mapping of instructions onto processors. While there is consensus that dataflow does expose the maximum amount of parallelism, there is considerable debate surrounding efficiency of the execution mechanism. This criticism centers on three points: (1) the number of instructions executed, (2) the relative power of a dataflow instruction and (3) the cost and complexity of a data driven processor.

Recently there has been significant progress in compiling scientific codes for dataflow machines [7]. Substantial programs written in the high-level language Id [41][40] and compiled for a dataflow machine yield dynamic instruction mixes (*e.g.* percentage of floating point operations) that are nearly equivalent to the same algorithms compiled from FOR-

TRAN and executing on a reduced-instruction set sequential uniprocessor [20] [8]. Moreover, the dataflow program executes essentially the same number of instructions independent of the number of processors, whereas the parallelization of a program for a conventional multiprocessor invariably incurs non-trivial execution overhead (*e.g.* synchronizing through barriers, task creation [8]) and typically yields *far less* parallel activity.

But how are we to compare the cost, in terms of processor complexity, of executing a dataflow instruction versus executing an instruction from a sequential stream? First, the operation performed by a dataflow instruction is similar in power to an operation on a conventional load/store machine, *i.e.* ADD, MULT, LOAD, STORE, BRANCH *etc.*¹ The difference lies in the way instructions are scheduled. In the von Neumann model, the operands for an instruction are assumed to be available when the program counter points to the instruction. In a dataflow machine an instruction is *asynchronously* scheduled only when its operands have been produced by some other parallel activities, so a dataflow machine must have an **operand matching** mechanism for detecting when an instruction has its required operands. Several general purpose dataflow machines have been built (*e.g.* ETL Sigma-1 [26], Manchester Machine [22]) or extensively simulated (*e.g.* M.I.T. TTDA [14]). But it is clear that these machines are far from commercial practicality². A primary reason for this, and a key criticism of dataflow machines, is the complexity of the operand matching mechanism [21].

Our goal is to discover implementation techniques that improve the cost/performance ratio of dataflow processors. Central to the work presented here is a new approach to operand matching. An **Explicit Token Store (ETS)** machine directly executes dataflow graphs while incorporating a new model of storage. The ETS allows the operand matching storage for the execution of a function invocation to be coalesced into an activation frame which is explicitly managed by the compiler. This enables implementation of the operand matching store with conventional (as opposed to content-addressable) memory technology and permits the realization of well-balanced pipelines.

Although this work focuses on the derivation of the ETS within the realm of dataflow architectures, we take the opportunity by way of this introduction to build a stronger case for machines which support fine grain synchronization. We believe that parallel machine architects eventually will have to apply the same concern for the efficient coordination of parallel activities as they presently do for fast sequential execution within a processor. Only then can we expect significant progress in exploiting parallelism in the general purpose setting.

1.1 Why Fine Grain Synchronization

Most multiprocessors are very bad at managing parallelism. Programmers and compiler writers are painfully aware of this fact. The more finely a program is divided into tasks, the greater the opportunity for parallel execution. However, there is a commensurate increase in the frequency of inter-task communication and associated synchronization. So exposing

¹Here we are restricting our attention to the instruction set of the M.I.T. Tagged-Token Dataflow Architecture (TTDA).

²The ETL Sigma-1 is the best engineered of the group. Presently, a 128 processor 640 MIPS engineering prototype is now operational at the MITI Electro-Technical Laboratory in Japan.

more parallelism, by way of identifying more tasks, does not obviously make the program run faster. In fact, we claim that is largely **unprofitable** to expose most of the latent parallelism in programs unless synchronization and task management are as efficient as the scheduling of primitive operations like add and multiply.

For a given machine, there is a fundamental tradeoff between the amount of parallelism that is profitable to expose and the overhead of synchronization. Sarkar[48] articulates this tradeoff as the competing contributions of the **ideal parallel execution time**, the amount of time required to execute the program in the absence of overhead, versus the **overhead factor**, the extra work required to schedule and coordinate the tasks. The ideal parallel execution time is multiplied by the overhead factor to yield the **actual parallel execution time**, the amount of time required to complete the problem for a given task granularity in the presence of scheduling overhead.

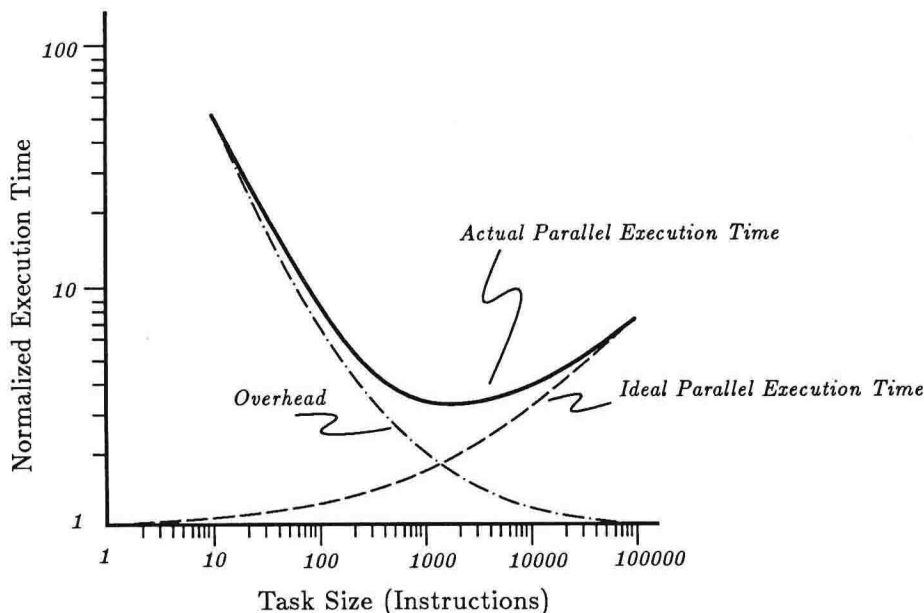
Figure 1.1 illustrates the general characteristic of the parallelism-overhead tradeoff for a typical program running on a machine with ten processors. This plot is suggestive of what was experienced running various programs on contemporary multiprocessors, but it does not express the data from a specific machine or application. The normalized execution time is the ratio n/s where n is the number of processors and s is the actual speedup relative to a single processor. The normalized ideal parallel execution time increases from 1 to n as the task granularity increases from 1, a single instruction, to 100,000 when the entire program executes as a single task. The task overhead factor is given by $(g + o)/g$ where g is the task size in instructions and o is the per-task overhead, in this case 1000 instructions³.

In this example the actual execution time is a minimum for a task size of about 2,000 instructions yielding a normalized execution time of three — that is, three times slower than the best ideal time. And this plot is optimistic. It is not possible, in general, to partition a given program into tasks of equal size, and the task overhead is a complicated expression that depends upon how the program was partitioned into tasks (*e.g.* the communication overhead component is a function of the number and sizes of data structures shared by two tasks). Thus, achieving optimal performance in the presence of overhead is much more difficult than simply finding the intersection of the ideal execution and overhead factor curves in Figure 1.1.

1.1.1 Tasks Should be Cheap

Look again at Figure 1.1. We think there is something fundamentally wrong — why should the task overhead be so extraordinarily large? Is it an inherent aspect of parallel computation or an artifact of the processor architecture used to construct the multiprocessor? We believe the latter. It seems counterproductive to force a programmer or compiler to expend so much effort working around what amounts to a basic deficiency in the processor architecture. If the task overhead were essentially zero, more parallelism would be exposed and exploited and compiling would be far easier. In fact, the entire partitioning and scheduling problem solved by Sarkar for the functional language SISAL [38] would be moot on a machine that had very low task overhead.

³The example employs a per-task overhead of 1,000 instructions which is very *low* as compared to the tasking costs on current commercial multiprocessors and operating systems.



10 Processors, 100000 Instructions, Overhead = 1000 Instructions/Task

Figure 1.1: The Parallelism-Overhead Tradeoff (from Sarkar)

By “low overhead” we mean on the order of one to ten instructions, basically **orders of magnitude** better than contemporary multiprocessors. We think the parallelism-overhead tradeoff should look like Figure 1.2. In this ideal world, task overhead is not a first order issue. Instead, the objective is to expose the maximum amount of parallelism by subdividing the computation as finely as possible.

1.1.2 A Virtual Memory Analogy

In a way, our argument for hardware support of fine grain synchronization is not unlike the case for architectural changes in order to efficiently implement demand-paged virtual memory. Demand-paging is a convenience for the programmer. It is very tedious for a programmer to manage overlays of code or data, and overlays clash with the semantics of modern programming languages. When data structures get large and access patterns are unpredictable, the programmer using overlays essentially emulates a virtual memory system, with the attendant loss of efficiency — both of the machine and the programmer.

While it is possible to look for the compiler-*forte* that can sift through the access patterns of programs and insert page management code, the problem is so intractable that we have come to insist on hardware support. The essential changes to the hardware are simple to describe, but the effect on the machine architecture is pervasive. Address translation and page fault detection must correctly occur on *every* memory reference, and each instruction that touches memory must be restartable. These two requirements

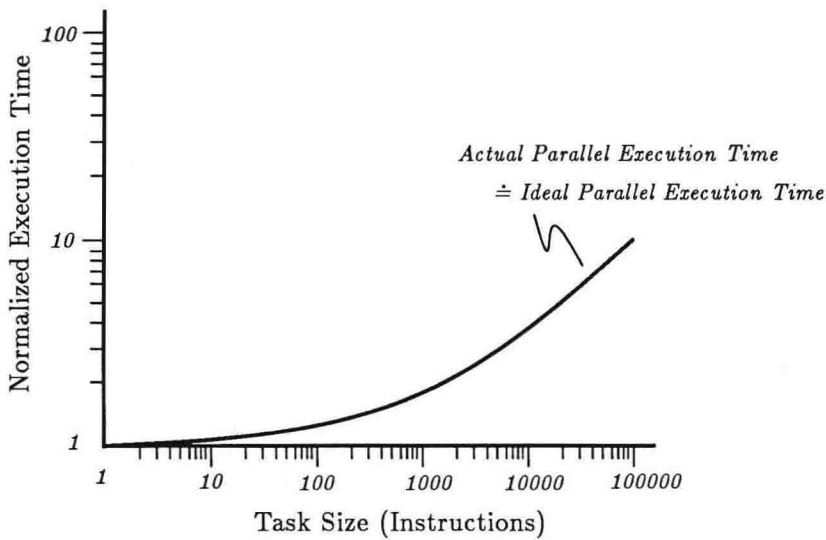


Figure 1.2: Speedup vs. Task Size When Task Overhead is Very Low

affect almost every aspect of the implementation from cache design to the basic instruction interpretation mechanisms.

Why should the management of complex interactions among concurrently executing entities be any *easier*? At least the demand paging analysis need only focus upon the address patterns of a single locus of control, whereas intertask synchronization analysis must discover the *interaction* of addressing patterns across multiple control loci. We think that a parallel machine should provide, at a minimum, real time synchronization checks on every memory load and store.

We also maintain that compiling for a parallel machine is greatly simplified when the task namespace is “virtualized”; that is, the namespace is much larger than the number of physical processors, so tasks are not bound to limited processor resources like multiple register sets.

An analogy for a machine that supports only a small fixed number of contexts (*cf* the Denelcor HEP [50]) would be a virtual memory system where the hardware supports a small number (say, 64) of simultaneous page translations, *but* where issuing an address that is not presently in the translation buffer causes an unrecoverable error (as opposed to a restartable fault). This puts the burden on the compiler to manage the translation buffer, explicitly inserting and deleting entries. It is not clear that this is any easier than performing demand paging on a machine with *no* hardware support. Similarly, a machine that provides a virtualized task namespace relieves the compiler and runtime system of the tough job of managing a small set of active tasks drawn from the large set of *possibly* active tasks.

In the remainder of this chapter we look more carefully at the overhead incurred in the parallel execution of a program. We build an informal model of parallel computation based upon concurrently executing tasks that can communicate through shared memory. While a

conventional processor executes the sequential portion of these tasks very well, it provides little support for moderating the interaction among tasks. Our objective is to sensitize the reader to two of our axioms: (1) a machine should support lots of simultaneous tasks and (2) synchronizing and scheduling these tasks should be very cheap. Not surprisingly, we believe that dataflow machines possess these properties. This work presents a dataflow processor architecture which, we believe, is a simple and appropriate building block for general purpose multiprocessors.

1.1.3 An Informal Task Model of Parallel Computation

We present a model of parallel execution to build our intuition about the cause of high overhead in parallel computation. We do this to motivate architectural support for fine grain synchronization, *not* to put forward an exact, or even complete, performance model for multiprocessors.

We are strongly biased towards programming parallel machines in high level languages which support a dynamic model of storage (*i.e.* a heap). At present, we do not see how modern languages can be effectively compiled for machines that do not directly supply a uniform address space without forsaking the ability to freely reference shared objects from within the language. In **message passing** machines (*e.g.* the Cosmic Cube [49] and Intel iPSC) the only way to share data is for the programmer to explicitly code commands to move data from one processor to another, so references to shared objects are emulated by the programmer and are not part of the machine language⁴. We focus instead upon **uniform address space** machines (*e.g.* the BBN Butterfly [47], IBM's RP3 [44], Alliant and Cedar [33], Cray-XMP) which provide hardware interpretation of references to addresses which are non-local, meaning there is no requirement for a compile-time distinction between local and non-local storage.

Suppose that we represent an executing parallel program as a **task graph** of inter-dependent sequential tasks. The nodes of the graph represent activations of sequential tasks, the local storage for which will be contained by an **activation frame**. The edges represent inter-task control and data dependences. We further assume that the execution dependences form a tree, although tasks can also communicate through shared objects on a heap for which references are passed as arguments and results. The shape of the task tree and the connectivity of objects on the heap, in general, cannot be determined statically, as we also permit fully recursive application and dynamic allocation of shared objects. For example, consider the following program which initially invokes **f**, which in turn allocates and returns as a result the object **A**, and makes calls to **g** and **h**;

```
def f() =                                def g(A) =                            def h(A) =
{ allocate A;                            { call g1();                            { call h1(A)};
  call g(A);                             call g2(A)};
  call h(A);
  return A };
```

⁴The driving concern is efficiency. It is certainly possible to *emulate* a shared address space on a message passing machine. The inefficiencies of such a scheme should be obvious, notwithstanding the efforts of intensive compile-time analysis.

We have not yet taken a position on how f is to be computed. Namely, what order is to be imposed on the computation of g and h ? There are roughly three possibilities.

1. **Sequential.** The familiar sequential order: f calls g and f suspends; g computes, terminates and then continues f ; f calls h and suspends; h computes, terminates and continues f ; f terminates. These rules are applied recursively to g , h , $g1$, $g2$ and $h1$. So at any time there is a stack of activation frames, of which only the top one is active.
2. **Fringe Parallel.** The calls to g and h are made in parallel: f forks g and h as independent tasks, and f suspends; g and h compute; when *both* g and h have terminated f is continued (*i.e.* g and h are joined with f); f terminates. These rules are applied recursively to g , h , $g1$, $g2$ and $h1$, but notice that only the leaves of the calling tree are active simultaneously.
3. **Fully Parallel.** A parallel call is made to g and h except that f is not suspended. f terminates after g and h terminate. The recursive application of this rule enables all tasks to run concurrently in the tree of active tasks.

We are not concerned with how the programmer indicates which evaluation method to use, nor whether the programming language semantics are functional or imperative. In either the parallel call or fully parallel case, the task graph unfolds as shown in Figure 1.3. We note that if calls proceed in parallel, there is no way *a priori* to allocate and deallocate local storage for the simultaneous tasks from a stack. In this case, tasks must be given separate activation records which, in general, must be managed more nearly like a heap.

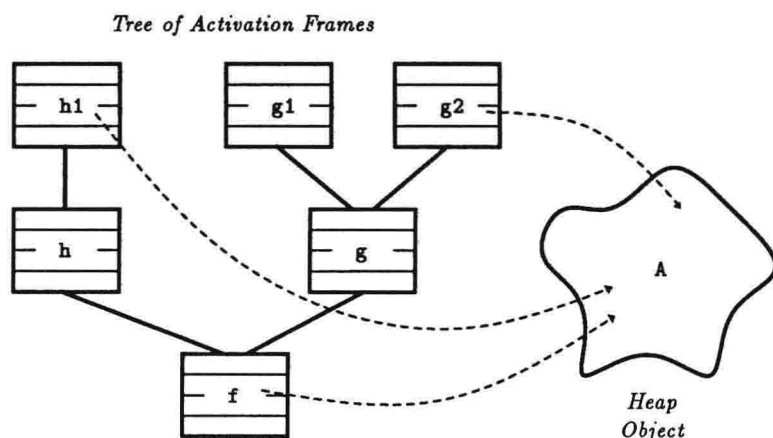


Figure 1.3: Simultaneous Tasks Referencing a Shared Object

The decomposition of the program into parallel tasks is usually accomplished explicitly by the programmer by annotating the source program [31], although functional language (*e.g.* FP [15], SISAL [38], Id [41][40]) programs can be partitioned automatically by the compiler. There is also considerable effort in the semi-automatic partitioning of programs written in sequential languages, notably FORTRAN, into parallel tasks [3][2].

1.1.4 Parallelism and Synchronization

The opportunity for parallel execution arises by noticing that all tasks whose dependences are satisfied can be freely scheduled. A dependence may either be a parameter to a task (*e.g.* the pointer to an object), a control prerequisite (*e.g.* the termination of another task, the acquisition of a resource lock), or dynamic data dependence through the heap (*e.g.* an element of an array).

The determination of the set of executable tasks is fundamentally a problem of **synchronization**, the act of translating the implicit or explicit assertion of a dependence by one task into a decision either to schedule or block another dependent task. The forms of synchronization required to support various dependences found in parallel programs include the following basic operations:

1. **Producer-Consumer.** A task produces a data structure that is read by another task. If the tasks are executed in parallel, synchronization is needed to avoid the **read-before-write** race.
2. **Forks and Joins.** A parallel call forks the thread of control into two tasks which is subsequently joined back together.
3. **Mutual Exclusion.** Concurrent procedures may emit requests that must be processed one at a time, *e.g.* updating an object or the serialization in the use of a resource.

All forms of synchronization require the *naming* of a synchronization event. At least one bit of state must be associated with this name, indicating whether the event has occurred or is **pending**. When the event occurs, the synchronization can complete. For producer-consumer synchronization this means that a waited-for value has become valid and the consumer can read it; for a join this means that both threads have reached the join and that one may proceed; for mutual exclusion this means that the exclusive resource has been freed for use by a pending requester. We say that a task is **blocked** when it is waiting for a pending event.

The more finely a program is divided into independent tasks, the greater the opportunity for exploiting parallelism; but there is typically a similar increase in the frequency of synchronization. Each pending synchronization event requires a unique name. Thus, as the number of potentially concurrent activities increase, so does the number of simultaneously named synchronization events. If the total number of concurrent tasks are identified dynamically (*e.g.* parallel execution of nested loops, recursive execution, or even separate compilation) then maximum number of synchronization events is difficult, or impossible, to predict statically. If the synchronization namespace is small, limited for example by the number of registers in a processor, then exposing parallelism is apt to be more difficult as the synchronization namespace must be carefully managed. This requires either static analysis at compile time (tantamount to global register allocation) or fairly expensive runtime (*e.g.* operating system) management.

Aside from how a synchronization point is named, the most important property of an implementation is how the event is *related* to the completion of synchronization. In an **event driven** system, the event directly causes the blocked task(s) to be scheduled,