

Instruction-Level Parallel Processing: History, Overview and Perspective

B. Ramakrishna Rau, Joseph A. Fisher
Computer Systems Laboratory
HPL-92-132
October, 1992

instruction-level
parallelism,
VLIW processors,
superscalar
processors, pipelining,
multiple operation
issue, speculative
execution, scheduling,
register allocation

Instruction-level Parallelism (ILP) is a family of processor and compiler design techniques that speed up execution by causing individual machine operations to execute in parallel. Although ILP has appeared in the highest performance uniprocessors for the past 30 years, the 1980s saw it become a much more significant force in computer design. Several systems were built, and sold commercially, which pushed ILP far beyond where it had been before, both in terms of the amount of ILP offered and in the central role ILP played in the design of the system. By the end of the decade, advanced microprocessor design at all major CPU manufacturers had incorporated ILP, and new techniques for ILP have become a popular topic at academic conferences. This article provides an overview and historical perspective of the field of ILP and its development over the past three decades.

1 Introduction

Instruction-level Parallelism (ILP) is a family of processor and compiler design techniques that speed up execution by causing individual machine operations, such as memory loads and stores, integer additions and floating point multiplications, to execute in parallel. The operations involved are normal RISC-style operations, and the system is handed a single program written with a sequential processor in mind. Thus an important feature of these techniques is that like circuit speed improvements, but unlike traditional multiprocessor parallelism and massive parallel processing, they are largely transparent to users. VLIWs and superscalars are examples of processors that derive their benefit from instruction-level parallelism, and software pipelining and trace scheduling are example software techniques that expose the parallelism that these processors can use.

Although small amounts of ILP has been present in the highest performance uniprocessors of the past 30 years, the 1980s saw it became a much more significant force in computer design. Several systems were built, and sold commercially, which pushed ILP far beyond where it had been before, both in terms of the amount of ILP offered and in the central role ILP played in the design of the system. By the early 1990s, advanced microprocessor design at all major CPU manufacturers incorporated ILP, and new techniques for ILP became a popular topic at academic conferences. With all of this activity, we felt that, in contrast to a report on suggested future techniques, there would be great value in gathering, in an archival reference, reports on experience with real ILP systems and reports on the measured potential of ILP. Thus this special issue of the Journal of Supercomputing.

1.1 ILP Execution

A typical ILP processor has the same type of execution hardware as a normal RISC machine. The differences between a machine with ILP and one without is that there may be more of that hardware, for example several integer adders instead of just one, and that the control will allow, and possibly arrange, simultaneous access to whatever execution hardware is present.

Consider the execution hardware of a simplified ILP processor consisting of four functional units and a branch unit connected to a common register file (Table 1). Typically ILP execution hardware allows multiple-cycle operations to be pipelined, so we may assume that a total of four

operations can be initiated each cycle. If in each cycle the longest latency operation is issued, this hardware could have 10 operations "in flight" at once, which would give it a maximum possible speed-up of a factor of 10 over a sequential processor with similar execution hardware. As the papers in this issue show, this execution hardware resembles that of several VLIW processors that have been built and used commercially, though it is more limited in its amount of ILP. Several superscalar processors now being built also offer a similar amount of ILP.

Table 1: Execution hardware for a simplified ILP processor

Functional Unit	Operations Performed	Latency
Integer Unit 1	Integer ALU operations	1
	Integer multiplication	2
	Loads	2
	Stores	1
Integer Unit 2 / Branch Unit	Integer ALU operations	1
	Integer multiplication	2
	Loads	2
	Stores	1
	Test-and-branch	1
Floating-point Unit 1 Floating-point Unit 2	Floating-point operations	3

There is a large amount of parallelism available even in this simple processor. The challenge is to make good use of it--we will see that with the technology available today, an ILP processor is unlikely to achieve nearly as much as a factor of 10 on many classes of programs, though scientific programs, and others, can yield far more than that on a processor which has more functional units. The first question that comes to mind is whether enough ILP exists in programs to make this possible. Then, if this is so, what must the compiler and hardware do to successfully exploit it? In reality, as we shall see in Section 4, the two questions have to be reversed; in the absence of techniques to find and exploit ILP, it remains hidden, and we are left with a pessimistic answer.

Figure 1a shows a very large expression taken from the inner loop of a compute-intensive program. It is presented cycle-by-cycle as it might execute on a processor with similar functional units to those shown in Table 1, but capable of having only one operation in flight at a time. Figure 1b shows the same program fragment as it might be executed on the hardware of Table 1.

```

CYCLE 1 xseed1 = xseed * 1309
CYCLE 2 nop
CYCLE 3 nop
CYCLE 4 yseed1 = yseed * 1308
CYCLE 5 nop
CYCLE 6 nop
CYCLE 7 xseed2 = xseed1 + 13849
CYCLE 8 yseed2 = yseed1 + 13849
CYCLE 9 xseed = xseed2 && 65535
CYCLE 0 yseed = yseed2 && 65535
CYCLE 11 tseed1 = tseed * 1307
CYCLE 12 nop
CYCLE 13 nop
CYCLE 14 vseed1 = vseed * 1306
CYCLE 15 nop
CYCLE 16 nop
CYCLE 17 tseed2 = tseed1 + 13849
CYCLE 18 vseed2 = vseed1 + 13849
CYCLE 19 tseed = tseed2 && 65535
CYCLE 20 vseed = vseed2 && 65535
CYCLE 21 xsq = xseed * xseed
CYCLE 22 nop
CYCLE 23 nop
CYCLE 24 ysq = yseed * yseed
CYCLE 25 nop
CYCLE 26 nop
CYCLE 27 xysumsq = xsq + ysq
CYCLE 28 tsq = tseed * tseed
CYCLE 29 nop
CYCLE 30 nop
CYCLE 31 vsq = vseed * vseed
CYCLE 32 nop
CYCLE 33 nop
CYCLE 34 tvsumsq = tsq + vsq
CYCLE 35 plc = plc + 1
CYCLE 36 tp = tp + 2
CYCLE 37 if xysumsq > radius goto @xy-no-hit

```

(a)

	INT ALU	INT ALU	FLOAT ALU	FLOAT ALU
CYCLE 1	tp=tp+2	plc=plc+1	vseed1=vseed*1306	tseed1=tseed*1307
CYCLE 2			yseed1=yseed*1308	xseed1=xseed*1309
CYCLE 3	nop			
CYCLE 4	vseed2=vseed1+13849	tseed2=tseed1+13849		
CYCLE 5	yseed2=yseed1+13849	xseed2=xseed1+13849		
CYCLE 6	yseed=yseed2&&65535	xseed=xseed2&&65535		
CYCLE 7	vseed=vseed2&&65535	tseed=tseed2&&65535	ysq=yseed*yseed	xsq=xseed*xseed
CYCLE 8			vsq=vseed*vseed	tsq=tseed*tseed
CYCLE 9	nop			
CYCLE 0	xysumsq=xsq+ysq			
CYCLE 11	tvsumsq=tsq+vsq	if xysumsq>radius goto @xy-no-hit		

(b)

Figure 1. (a) An example of the sequential record of execution for a loop. (b) The instruction-level parallel record of execution for the same loop.

Note that several of the cycles in Figure 1a contain no-ops. This is because the sequential processor must await the completion of the 3-cycle latency multiply issued in cycle 1 before issuing the next operation. (These no-ops would not appear in the text of a program, but are shown here as the actual record of what is executed each cycle.) Most instruction-level parallel processors can issue operations during these nop cycles, when previous operations are still in flight, and many can issue more than one operation in a given cycle. In our ILP record of execution (Figure 1b), both effects are evident: in cycle 1, 4 operations are issued; in cycle 2, 2 more operations are issued even though neither multiply in cycle 1 has yet completed execution.

This special issue of The Journal of Supercomputing concerns itself with the technology of systems which try to attain the kind of record of execution in Figure 1b, given a program written with the record of execution in Figure 1a in mind.

1.2 Early History of Instruction-level Parallelism

In small ways, instruction-level parallelism factored into the thinking of machine designers in the 1940s and 1950s. Parallelism that would today be called horizontal microcode appeared in Turing's 1946 design of the Pilot ACE [1], and was carefully described by Wilkes [2]. Indeed, in 1952 Wilkes and Stringer wrote, "In some cases it may be possible for two or more micro-operations to take place at the same time" [3].

The 1960s saw the appearance of transistorized computers. One effect of this revolution was that it became practical to build reliable machines with far more gates than was necessary to build a general-purpose CPU. This led to commercially successful machines which used this available hardware to provide instruction-level parallelism at the machine-language level. In 1963, Control Data Corporation started delivering its CDC 6600 [4, 5], which had 10 functional units--integer add, shift, increment (2), multiply (2), logical, branch, floating-point add and divide. Any one of these could start executing in a given cycle whether or not others were still processing data-independent earlier operations. In this machine the hardware decided, as the program executed, which operation to issue in a given cycle; its model of execution was well along the way toward what we would today call superscalar. Indeed, in many ways it strongly resembled its direct descendant, the scalar portion of the CRAY-1. The CDC 6600 was the scientific supercomputer of its day.

Also during the 1960s, IBM introduced, and in 1967-8 delivered, the 360/91 [6]. This machine, based partly on IBM's instruction-level parallel experimental Stretch processor, offered less instruction-level parallelism than the CDC-6600, having only a single integer adder, a floating-point adder, and a floating-point multiply/divide. But it was far more ambitious than the CDC 6600 in its attempt to rearrange the instruction stream to keep these functional units busy--a key technology in today's superscalar designs. For various non-technical reasons, the 360/91 was not as commercially successful as it might have been, with only about 20 machines delivered [7]. But its CPU architecture was the start of a long line of successful high performance processors. As with the CDC 6600, this ILP pioneer started a chain of superscalar architectures that has lasted into the 1990s.

In the 1960s, research into "parallel processing" often was concerned with the ILP found in these processors. By the mid-1970s, the term was used more often for multiple processor parallelism, and for regular array and vector parallelism. In part, this was due to some very pessimistic results about the availability of ILP in ordinary programs, which we discuss below.

1.3 Modern Instruction-Level Parallelism

In the late 1970s the beginnings of a new style of ILP, called VLIW (for Very Long Instruction Word), emerged on several different fronts. In many ways, VLIWs were a natural outgrowth of horizontal microcode, the first ILP technology, and they were triggered, in the 1980s, by the same changes in semiconductor technology that had such a profound impact upon the entire computer industry.

For sequential processors, as the speed gap between writeable and read-only memory narrowed, the advantages of a small, dedicated, read-only control store began to disappear. One natural effect of this was to diminish the advantage of microcode; it no longer made as much sense to define a complex language as a compiler target, and then interpret this in very fast read-only microcode. Instead, the vertical microcode interface was presented as a clean, simple compiler target. This concept was called RISC [8-10]. In the 1980s the general movement of microprocessor products was towards the RISC concept and instruction-level parallel techniques fell out of favor. In the mini-supercomputer price-bracket though, one innovative superscalar product, the ZS-1, which could issue up to two instructions each cycle, was built and marketed by Astronautics [11].

The same changes in memory technology were having a somewhat different effect upon horizontally microcoded processors. During the 1970s, a large market had grown in specialized signal processing computers. Not aimed at general purpose use, these CPUs hard-wired FFTs and other important algorithms directly into the horizontal control store, gaining tremendous advantages from the instruction-level parallelism available there. When fast, writable memory became available, some of these manufacturers, most notably Floating Point Systems [12], replaced the read-only control store with writable memory, giving users access to instruction-level parallelism in far greater amounts than the early superscalar processors had. These machines were extremely fast, the fastest processors by far in their price ranges, for important classes of scientific applications. However, despite attempts on the part of several manufacturers to market their products for more general, everyday use, they were almost always restricted to a narrow class of applications. This was caused by the lack of good system software, which in turn was caused by the idiosyncratic architecture of processors built for a single application, and by the lack at that time of good code generation algorithms for ILP machines with that much parallelism.

As with RISC, the crucial step was to present a simple, clean interface to the compiler. However, in this case the clean interface was horizontal, not vertical, so as to afford greater ILP [13, 14]. This style of architecture was dubbed VLIW [14]. Code generation techniques, some of which had been developed for generating horizontal microcode, were extended to these general-purpose VLIW machines, so that the compiler could specify the parallelism directly [15, 16].

In the 1980s, VLIW CPUs were offered commercially in the form of capable, general-purpose machines. Three computer startups--Culler, Multiflow and Cydrome--built VLIWs with varying degrees of parallelism [17, 18]. As a group, these companies were able to demonstrate that it was possible to build practical machines which achieved large amounts of ILP on scientific and engineering codes. Although, for various reasons, none was a lasting business success, several major computer manufacturers acquired access to the technologies developed at these startups, and there are several active VLIW design efforts underway. Furthermore, many of the compiler techniques developed with VLIWs in mind, and reported upon in this issue, have been used to compile for superscalar machines as well.

ILP in the 1990s

Just as had happened 30 years ago when the transistor became available, CPU designers in the 1990s now have offered to them more silicon space on a single chip than a RISC processor requires. Virtually all designers have begun to add some degree of superscalar capability, and some are investigating VLIWs as well. It is a safe bet that by 1995 virtually all new CPUs will embody some degree of ILP.

Partly as a result of this commercial resurgence of interest in ILP, research into that area has become a dominant feature of architecture and systems conferences of the 1990s. Unfortunately, those researchers who found themselves designing state-of-the-art products at computer startups did not have the time to document the progress that was made, and the large amount that was learned. Virtually everything that was done by these groups was relevant to what designers wrestle with today. This issue represents an attempt to archive some of those advances.

2 ILP Architectures

The end result of instruction-level parallel execution is that multiple operations are simultaneously in execution, either as a result of having been issued simultaneously or because the time to execute an operation is greater than the interval between the issuance of successive operations. How exactly are the necessary decisions made as to when an operation should be executed and whether an operation should be speculatively executed? The alternatives can be broken down depending on the extent to which these decisions are made by the compiler rather than by the hardware and on the manner in which information regarding parallelism is communicated by the compiler to the hardware via the program.

A computer architecture is a contract between the class of programs that are written for the architecture and the set of processor implementations of that architecture. Usually this contract is concerned with the instruction format and the interpretation of the bits that constitute an instruction, but in the case of ILP architectures it extends to information embedded in the program pertaining to the available parallelism between the instructions or operations in the program. With this in mind, ILP architectures can be classified as follows.

- o **Sequential architectures:** architectures for which the program is not expected to convey any explicit information regarding parallelism. Superscalar processors ([4, 19, 20, 11, 21-28]) are representative of ILP processor implementations for sequential architectures.
- o **Dependence architectures:** architectures for which the program explicitly indicates the dependences that exist between operations. Dataflow processors ([29-31]) are representative of this class.
- o **Independence architectures:** architectures for which the program provides information as to which operations are independent of one another. Very Long Instruction Word (VLIW) processors ([12, 17, 18]) are examples of the class of independence architectures.

In the context of this taxonomy, vector processors [32-34] are best thought of as processors for a sequential, CISC (complex instruction set computer) architecture. The complex instructions are the vector instructions which do possess a stylized form of instruction-level parallelism internal to each vector instruction. Attempting to execute multiple instructions in parallel, whether scalar or vector, incurs all of the same problems that are faced by a superscalar processor. Because of their stylized approach to parallelism, vector processors are less general in their ability to exploit all forms of instruction-level parallelism. Nevertheless, vector processors have enjoyed great commercial success over the past decade. Not being true ILP processors, vector processors are outside the scope of this special issue. (Vector processors have received large amounts of attention, elsewhere, over the past decade and, have been treated extensively in many books and articles, for instance, the survey by Dongarra [35] and the book by Schneck [36].) Also, certain hybrid architectures [37-39], which also combine some degree of multi-threading with ILP, fall outside of this taxonomy for uniprocessors.

If ILP is to be achieved, between the compiler and the runtime hardware, the following functions must be performed:

- o the dependences between operations must be determined,
- o the operations, that are independent of any operation that has not as yet completed, must be determined, and
- o these independent operations must be scheduled to execute at some particular time, on some specific functional unit, and must be assigned a register into which the result may be deposited.

Figure 2 shows the breakdown of these three tasks, between the compiler and runtime hardware, for the three classes of architecture.

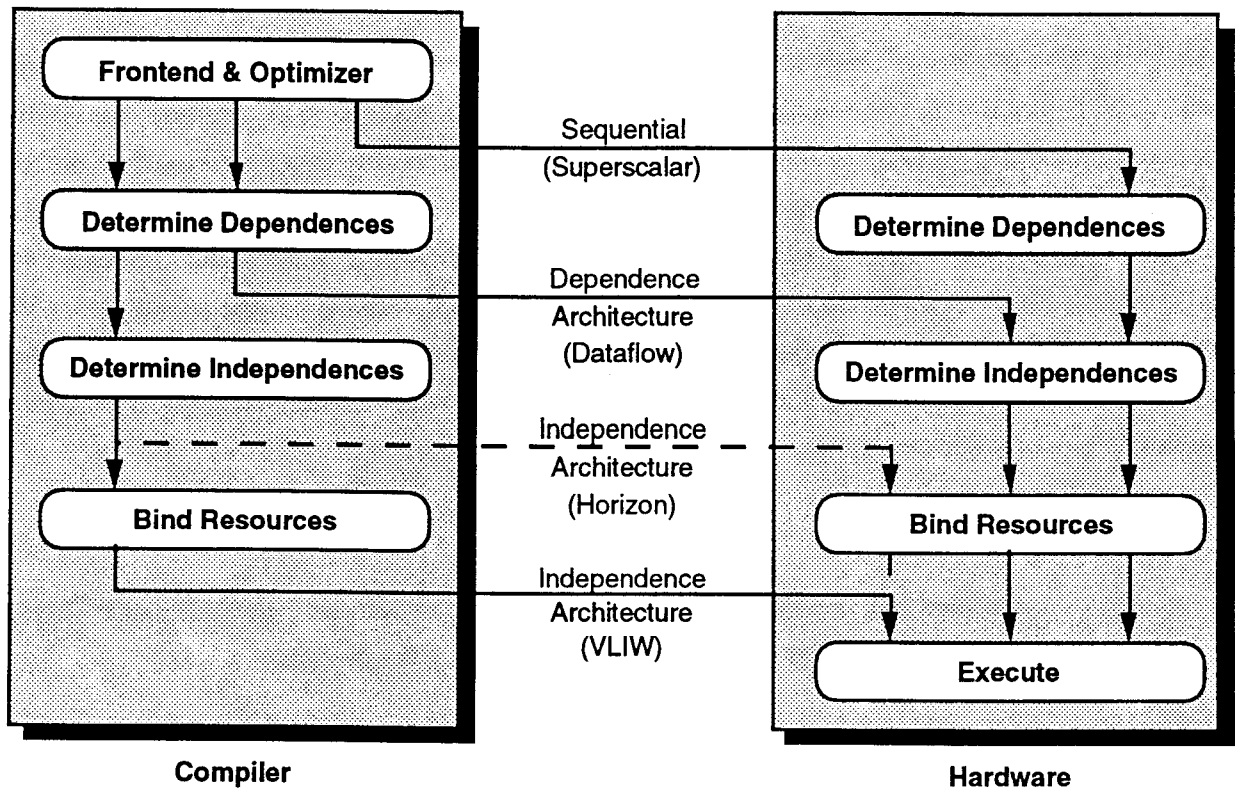


Figure 2. Division of responsibilities between the compiler and the hardware for the three classes of architecture.

2.1 Sequential architectures and superscalar processors

The program for a sequential architecture contains no explicit information regarding the dependences that exist between instructions. Consequently, the compiler need neither identify parallelism nor make scheduling decisions since there is no explicit way to communicate this information to the hardware. (It is true, nevertheless, that there is value in the compiler performing these functions and ordering the instructions so as to facilitate the hardware's task of extracting parallelism.) In any event, if instruction-level parallelism is to be employed, the dependences that exist between instructions must be determined by the hardware. It is only necessary to determine dependences with sequentially preceding operations that are in flight, i.e., those that have been issued but have not yet completed. The operation is now independent of all other operations and may begin execution. At this point, the hardware must make the scheduling decision of when and where this operation is to execute.

A superscalar processor¹ strives to issue an instruction every cycle, so as to execute many instructions in parallel, even though the hardware is handed a sequential program. The problem is that a sequential program is constructed with the assumption only that it will execute correctly when each instruction waits for the previous one to finish, and that is the only order that the architecture guarantees to be correct. The first task, then, for a superscalar processor is to understand, for each instruction, which other instructions it actually is dependent upon. With every instruction that a superscalar processor issues, it must check whether the instruction's operands (registers or memory locations that the instruction uses or modifies) interfere with the operands of any other instruction in flight, i.e., one that is either:

- o already in execution, or
- o has been issued but is waiting for the completion of interfering instructions that would have been executed earlier in a sequential execution of the program.

If either of these conditions is true, the instruction in question must be delayed until the instructions on which it is dependent have completed execution. For each waiting operation, these dependences must be monitored to determine the point in time at which neither condition is

¹ The first machines of this type that were built in the 1960s were referred to as look-ahead processors. Subsequently, machines that performed out-of-order execution, while issuing multiple operations per cycle, came to be termed superscalar processors. Since look-ahead processors are only quantitatively different from superscalar processors, we shall drop the distinction and refer to them, too, as superscalar processors.

true. When this happens, the instruction is independent of all other uncompleted instructions and can be allowed to begin executing at any time thereafter. In the meantime, the processor may begin execution of subsequent instructions which prove to be independent of all sequentially preceding instructions in flight. Once an instruction is independent of all other ones in flight, the hardware must also decide exactly when and on which available functional unit to execute the instruction. The Control Data CDC 6600 employed a mechanism, called the scoreboard, to perform these functions [4]. The IBM System/360 Model 91, built in the early 1960s, utilized an even more sophisticated method known as Tomasulo's Algorithm ([40]) to carry out these functions.

The further goal of a superscalar processor is to issue *multiple* instructions every cycle. The most problematic aspect of so doing is that of determining the dependences between the operations that one wishes to issue simultaneously. Since the semantics of the program, and in particular the essential dependences, are specified by the sequential ordering of the operations, the operations must be processed in this order to determine the essential dependences. This constitutes an unacceptable performance bottleneck in a machine that is attempting parallel execution. On the other hand, eliminating this bottleneck can be very expensive as is always the case when attempting to execute an inherently sequential task in parallel. An excellent reference on superscalar processor design and its complexity is the book by Johnson [41].

A number of superscalar processors have been built during the past decade including the Astronautics' ZS-1 decoupled access mini-supercomputer [11, 42], Apollo's DN10000 personal supercomputer [21, 22] and, most recently, a number of microprocessors [23-28].

Note that an ILP processor need not issue multiple operations per cycle in order to achieve a certain level of performance. For instance, instead of a processor capable of issuing five instructions per cycle, the same performance could be achieved by pipelining the functional units and instruction issue hardware five times as deeply, speeding up the clock rate by a factor of five but issuing only one instruction per cycle. This strategy, which has been termed **superpipelining** ([43]), goes full circle back to the single-issue, superscalar processing of the 1960s. Superpipelining may result in some parts of the processor (such as the instruction unit and communications busses) being less expensive and better utilized and other parts (such as the execution hardware) being more costly and less well used.

2.2 Dependence architectures and dataflow processors

In the case of dependence architectures, the compiler or the programmer identifies the parallelism in the program and communicates it to the hardware by specifying, in the executable program, the dependences between operations. The hardware must still determine, at run-time, when each operation is independent of all other operations and then perform the scheduling. However, the inherently sequential task, of scanning the sequential program in its original order to determine the dependences, has been eliminated.

The objective of a dataflow processor is to execute an instruction at the earliest possible time subject only to the availability of the input operands and a functional unit upon which to execute the instruction [29, 30]. To do so, it counts on the program to provide information about the dependences between instructions. Typically, this is accomplished by including in each instruction a list of successor instructions. (An instruction is a successor of another instruction if it uses as one of its input operands the result of that other instruction.) Each time an instruction completes, it creates a copy of its result for each of its successor instructions. As soon as all of the input operands of an instruction are available, the hardware fetches the instruction, which specifies the operation to be performed and the list of successor instructions. The instruction is then executed as soon as a functional unit of the requisite type is available. This property, whereby the availability of the data triggers the fetching and execution of an instruction, is what gives rise to the name of this type of processor. Because of this property, it is redundant for the instruction to specify its input operands. Rather, the input operands specify the instruction! If there is always at least one instruction ready to execute on every functional unit, the dataflow processor achieves peak performance.

Computation within a basic block typically does not provide adequate levels of parallelism. Superscalar and VLIW processors use control parallelism and speculative execution to keep the hardware fully utilized. (This is discussed in greater detail in Sections 3 and 4.) Dataflow processors have traditionally counted on using control parallelism alone to fully utilize the functional units. A dataflow processor is more successful than the others at looking far down the execution path to find abundant control parallelism. When successful, this is a better strategy than speculative execution since every instruction executed is a useful one and the processor does not have to deal with error conditions raised by speculative operations.

As far as the authors are aware, there have been no commercial products built based on the dataflow architecture, except in a limited sense [44]. There have, however, been a number of research prototypes built, for instance, the ones built at the University of Manchester [31] and at MIT [45].

2.3 Independence architectures and VLIW processors

In order to execute operations in parallel, the system must determine that the operations are independent of one another. Superscalar processors and dataflow processors represent two ways of deriving this information at run-time. In the case of the dataflow processor, the explicitly provided dependence information is used to determine when an instruction may be executed so that it is independent of all other concurrently executing instructions. The superscalar processor must do the same but, since programs for it lack any explicit information, it must also first determine the dependences between instructions. In contrast, for an independence architecture, the compiler identifies the parallelism in the program and communicates it to the hardware by specifying which operations are independent of one another. This information is of direct value to the hardware, since it knows with no further checking which operations it can execute in the same cycle. Unfortunately, for any given operation, the number of operations of which it is independent is far greater than the number of operations on which it is dependent. So, it is impractical to specify all independences. Instead, for each operation, independences with only a subset of all independent operations (those operations that the compiler thinks are the best candidates to execute concurrently) are specified.

By listing operations that could be executed simultaneously, code for an independence architecture may be very close to the record of execution produced by an implementation of that architecture. If the architecture additionally requires that programs specify where (on which functional unit) and when (in which cycle) the operations are executed, then the hardware makes no run-time decisions at all and the code is virtually identical to the desired record of execution. The VLIW processors that have been built to date are of this type and represent the predominant examples of machines with independence architectures. The program for a VLIW processor specifies exactly which functional unit each operation should be executed on and exactly when each operation should be issued so as to be independent of all operations that are being issued at the same time as well as of those that are in execution. A particular processor implementation of a VLIW architecture could choose to disregard the scheduling decisions embedded in the

program, making them at run-time instead. In doing so, the processor would still benefit from the independence information but would have to perform all of the scheduling tasks of a superscalar processor. Furthermore, when attempting to execute concurrently two operations that the program did not specify as being independent of each other, it must determine independence, just as a superscalar processor must.

With a VLIW processor, it is important to distinguish between an instruction and an operation. An operation is a unit of computation, such as an addition, memory load or branch, which would be referred to as an instruction in the context of a sequential architecture. A VLIW instruction is the set of operations that are intended to be issued simultaneously. It is the task of the compiler to decide which operations should go into each instruction. This process is termed scheduling. Conceptually, the compiler schedules a program by emulating at compile-time what a dataflow processor, with the same execution hardware, would do at run-time. All operations that are supposed to begin at the same time are packaged into a single VLIW instruction. The order of the operations within the instruction specifies the functional unit on which each operation is to execute. A VLIW program is a transliteration of a desired record of execution which is feasible in the context of the given execution hardware.

The compiler for a VLIW machine specifies that an operation be executed speculatively merely by performing speculative code motion, that is, scheduling an operation before the branch that determines that it should, in fact, be executed. At run-time, the VLIW processor blindly executes this operation exactly as specified by the program just as it would for a non-speculative operation. Speculative execution is virtually transparent to the VLIW processor and requires little additional hardware. When the compiler decides to schedule an operation for speculative execution, it can arrange to leave behind enough of the state of the computation to assure correct results when the flow of the program requires that the operation be ignored. The hardware required for the support of speculative code motion consists of having some extra registers, of fetching some extra instructions, and of suppressing the generation of spurious error conditions. The VLIW compiler must perform many of the same functions that a superscalar processor performs at run-time to support speculative execution, but it does so at compile-time.

The earliest VLIW processors built were the so-called attached array processors [46-48, 12, 49] of which the best known were the Floating Point Systems products, the AP-120B, the FPS-164 and the FPS-264. The next generation of products were the mini-supercomputers: Multiflow's Trace series of machines [17, 50] and Cydrome's Cydra 5 [51, 18, 52] and the Culler machine

for which, as far as we are aware, there is no published description in the literature. Over the last few years, the VLIW architecture has begun to show up in microprocessors [53-57].

Other types of processors with independence architectures have been built or proposed. A superpipelined machine may issue only one operation per cycle, but if there is no superscalar hardware devoted to preserving the correct execution order of operations, the compiler will have to schedule them with full knowledge of dependencies and latencies. From the compiler's point of view, these machines are virtually the same as VLIWs, though the hardware design of such a processor offers some tradeoffs with respect to VLIWs. Another proposed independence architecture, dubbed *Horizon* ([58]), encodes an integer H into each operation. The architecture guarantees that all of the next H operations in the instruction stream are data-independent of the current operation. All the hardware has to do to release an operation, then, is assure itself that no more than H subsequent operations are allowed to issue before this operation has completed. The hardware does all of its own scheduling, unlike VLIWs and deeply pipelined machines which rely on the compiler, but the hardware is relieved of the task of determining data dependence. The key distinguishing features of these three ILP architectures are summarized in Table 2.

3 Hardware and software techniques for ILP execution

Regardless of which ILP architecture is considered, certain functions must be performed if a sequential program is to be executed in an ILP fashion. The program must be analyzed to determine the dependences; the point in time at which an operation is independent, of all operations that are as yet not complete, must be determined; scheduling and register allocation must be performed; often, operations must be executed speculatively, which in turn requires that branch prediction be performed. All these functions must be performed. The choice is, first, whether they are to be performed by the compiler or by run-time hardware and, second, which specific technique is to be employed. These alternatives are reviewed in the rest of this section.

3.1 Hardware features to support ILP execution

Instruction-level parallelism involves the existence of multiple operations in flight at any one time, i.e., operations that have begun, but not completed, executing. This implies the presence of execution hardware that can simultaneously process multiple operations. This has, historically,

been achieved by two mechanisms: first, providing multiple, parallel functional units and, second, pipelining the functional units. Although both are fairly similar from a compiler's viewpoint--the compiler must find enough independent operations to keep the functional units busy--they have their relative strengths and weaknesses from a hardware viewpoint.

Table 2: A comparison of the instruction-level parallel architecture types discussed in this paper.

	Sequential Architecture	Dependence Architecture	Independence Architecture
Additional information required in the program	None	Complete specification of dependences between operations	Minimally, a partial list of independences. Typically, a complete specification of when and where each operation is to be executed
Typical kind of ILP processor	Superscalar	Dataflow	VLIW
Analysis of dependences between operations	Performed by hardware	Performed by the compiler	Performed by the compiler
Analysis of independent operations	Performed by hardware	Performed by hardware	Performed by the compiler
Final operation scheduling	Performed by hardware	Performed by hardware	Typically, performed by the compiler
Role of compiler	Rearranges the code to make the analysis and scheduling hardware more successful	Replaces some analysis hardware	Replaces virtually all the analysis and scheduling hardware

In principle, pipelining is the more cost-effective way of building ILP execution hardware. For the relatively low cost of adding pipeline latches within each functional unit, the amount of ILP can be doubled, tripled, or more. The limiting factor in increasing the performance by this means are the data and clock skews and the latch setup and hold times. These issues were studied during the 1960s and 1970s, and the upper limits on the extent of pipelining were determined [59-63]. However, the upper limit on pipelining is not necessarily best from the viewpoint of achieved performance. Pipelining adds delay to the execution time of individual operations (even though multiple of them can be in flight on the same functional unit). Especially on computations that have small amounts of parallelism, beyond a certain point the increase in the latency counterbalances the benefits of the increase in ILP, yielding lower performance [64]. Parallelism

achieved by adding more functional units does not suffer from this drawback, but has its own set of disadvantages. First, the amount of functional unit hardware goes up in linear proportion to the parallelism. Worse, the cost of the interconnection network and the register files goes up proportional to the square of the number of functional units since, ideally, each functional unit's output bus must communicate with every functional unit's input buses through the register file. Also, as the number of loads on each bus increases, so must the cycle time or the extent of pipelining, both of which degrade performance on computation with little parallelism.

The related techniques of pipelining and overlapped execution were employed as early as in the late 1950s in computers such as IBM's STRETCH computer [65, 66] and UNIVAC's LARC [67]. Traditionally, overlapped execution refers to the parallelism that results from multiple active instructions, each in a different one of the phases of instruction fetch, decode, operand fetch, and execute whereas pipelining is used in the context of functional units such as multipliers and floating-point adders [68, 69]. (A potential source of confusion is that, in the context of RISC processors, overlapped execution and pipelining, especially when the integer ALU is pipelined, have been referred to as pipelining and superpipelining, respectively [43].)

The organization of the register files becomes a major issue when there are multiple functional units operating concurrently. For ease of scheduling, it is desirable that every operation (except loads and stores) be register-register and that the register file be the hub for communication between all the functional units. However, with each functional unit performing two reads and one write per cycle from or to the register file, the implementation of the register file becomes problematic. The chip real-estate of a multi-ported register file is proportional to the product of the number of read ports and the number of write ports. The loading of multiple read ports on each register cell slows down the access time. For these reasons, highly parallel ILP hardware is structured as multiple clusters of functional units, with all the functional units within a single cluster sharing the same multi-ported register files [14, 70, 17, 50]. Communication between clusters is slower and occurs with lower bandwidth. This places a burden upon the compiler to partition the computation intelligently across the clusters; an inept partitioning can result in worse performance than if just a single cluster were used, leaving the rest of them idle.

The presence of multiple, pipelined function units places increased demands upon the instruction issue unit. In a fully sequential processor, each instruction is issued after the previous one has completed. Of course, this totally defeats the benefits of parallel execution hardware. However, if the instruction unit attempts to issue an instruction every cycle, care must be taken to not do so if

an instruction, upon which this one is dependent, is still not complete. The scoreboard in the CDC 6600 [4] was capable of issuing an instruction every cycle until an output dependence was discovered. In the process, instructions following one which was waiting on a flow dependence could begin execution. This was the first implementation of an out-of-order execution scheme. Stalling instruction issue is unnecessary on encountering an output dependence if register renaming is performed. The Tomasulo algorithm [40], which was implemented in the IBM System/360 Model 91 [19], is the classical scheme for register renaming and has served as the model for subsequent variations [71-73, 25, 74]. A different, programmatically-controlled register renaming scheme is obtained by providing rotating register files, i.e., base-displacement indexing into the register file using an instruction-provided displacement off a dedicated base register [12, 51, 75, 18]. Although applicable only for renaming registers across multiple iterations of a loop, rotating registers have the advantage of being considerably less expensive in their implementation than are other renaming schemes.

The first consideration given to the possibility of issuing multiple instructions per cycle from a sequential program was by Tjaden and Flynn [76]. This line of investigation into the logic needed to perform multiple-issue was continued by various researchers [72, 73, 77-81]. This idea, of multiple instruction issue of sequential programs, was probably first referred to as superscalar execution by Agerwala and Cocke [82]. A careful assessment of the complexity of the control logic involved in superscalar processors is provided by Johnson [41]. An interesting variation on multiple-issue, which made use of architecturally-visible queues to simplify the out-of-order execution logic, was the decoupled access/execute architecture proposed by Smith [83] and subsequently developed as a commercial product [11, 42].

A completely different approach to achieving multiple instruction issue which grew out of horizontal microprogramming was represented by attached processor products such as the Floating Point Systems AP-120B [48], the Polycyclic project at ESL [16, 84, 13], the Stanford University MIPS project [85] and the ELI project at Yale [14, 70]. The concept is to have the compiler decide which operations should be issued in parallel and to group them in a single, long instruction. This style of architecture, which was dubbed a Very Long Instruction Word (VLIW) architecture [14], has the advantage that the instruction issue logic is trivial in comparison to that for a superscalar machine but suffers the disadvantage that the set of operations that are to be issued simultaneously is fixed once and for all at compile-time. One of the implications of issuing multiple operations per instruction is that one needs the ability to issue (and process)

multiple branches per second. Various types of multiway branches, each corresponding to a different detailed model of execution or compilation have been suggested [86-88, 17].

The first obstacle that one encounters when attempting ILP computation is the, generally, small size of basic blocks. In light of the pipeline latencies and the inter-operation dependences, little instruction-level parallelism is to be found. It is important that operations from multiple basic blocks be executed concurrently if a parallel machine is to be fully utilized. Since the branch condition, which determines which block is to be executed next, is often resolved only at the end of a basic block, it is necessary to resort to speculative execution, i.e., continuing execution along one or more paths before it is known which way the branch will go. Dynamic schemes for speculative execution [72, 73, 89, 90] must provide ways to

- o terminate unnecessary speculative computation once the branch has been resolved,
- o undo the effects of the speculatively executed operations which should not have been executed,
- o ensure that no exceptions are reported until it is known that the excepting operation should, in fact, have been executed, and
- o preserve enough execution state at each speculative branch point to enable execution to resume down the correct path if the speculative execution happened to proceed down the wrong one.

All this can be expensive in hardware. The alternative is to perform speculative code motion at compile-time, i.e., move operations from subsequent blocks up past branch operations into preceding blocks. These operations will end up being executed before the branch that they were supposed to follow; hence, they are executed speculatively. Such code motion is fundamental to global scheduling schemes such as trace scheduling [91, 15, 92]. The hardware support needed is much less demanding. First, a mechanism to ensure that exceptions caused by speculatively scheduled operations are reported if and only if the flow of control is such that they would have been executed in the non-speculative version of the code [93] and, second, additional architecturally visible registers to hold the speculative execution state. A limited form of speculative code motion is provided by the "boosting" scheme [94, 95].

Since all speculative computation is wasted if the wrong path is followed, it is important that accurate branch prediction be used to guide speculative execution. Various dynamic schemes of

varying levels of sophistication and practicality have been suggested, that gather execution statistics of one form or another while the program is running [96-99]. The alternative is to use profiling runs to gather the appropriate statistics and to embed the prediction, at compile-time, into the program. Trace scheduling and superblock scheduling [100, 101] use this approach to reorder the control flow graph to reflect the expected branch behavior. Hwu, et al., claim better performance than with dynamic branch prediction [100]. Fisher and Freudenberger have examined the extent to which branch statistics gathered using one set of data are applicable to subsequent runs with different data [102]. Although static prediction can be useful for guiding both static and dynamic speculation, it is not apparent how dynamic prediction can assist static speculative code motion.

Predicated execution is an architectural feature that permits the execution of individual operations to be determined by an additional, boolean input. It has been used to selectively squash operations that have been moved up from successor blocks into the delay slots of a branch operation [103, 88]. In its more general form [51, 18, 52], it is used to eliminate branches in their entirety over an acyclic region of a control flow graph [104-106] which has been IF-converted [107].

3.2 ILP Compilation

3.2.1 Scheduling

Scheduling algorithms can be classified based on two broad criteria. The first one is the nature of the control flow graph that can be scheduled by the algorithm. The control flow graph can be described by the following two properties:

- o whether it consists of a single basic block or multiple basic blocks, and
- o whether it is an acyclic or cyclic control flow graph.

Algorithms that can only schedule single acyclic basic blocks are known as **local scheduling** algorithms. Algorithms that jointly schedule multiple basic blocks (even if these are multiple iterations of a single static basic block) are termed **global scheduling** algorithms. Acyclic global scheduling algorithms deal either with control flow graphs that contain no cycles or, more typically, cyclic graphs for which a self-imposed scheduling barrier exists at each back-edge in

the control flow graph. As a consequence of these scheduling barriers, back-edges present no opportunity to the scheduler and are, therefore, irrelevant to it. Acyclic schedulers can yield better performance on cyclic graphs by unrolling the loop, a transformation which though easier to visualize for cyclic graphs with a single back-edge, can be generalized to arbitrary cyclic graphs. The benefit of this transformation is that the acyclic scheduler now has multiple iterations worth of computation to work with and overlap. The penalty of the scheduling barrier is amortized over more computation. Cyclic global scheduling algorithms attempt to directly optimize the schedule across back-edges as well. Each class of scheduling algorithms is more general than the previous one and, as we shall see, attempts to build on the intuition and heuristics of the simpler, less general algorithm. As might be expected, the more general algorithms experience greater difficulty in achieving near-optimality or of even articulating intuitively appealing heuristics.

The second classifying criterion is the type of machine for which scheduling is being performed, which in turn is described by the following assumed properties of the machine:

- o finite vs. unbounded resources,
- o unit latency vs. multiple cycle latency execution, and
- o simple resource usage patterns for every operation (i.e., each operation uses just one resource for a single cycle, typically during the first cycle of the operation's execution) vs. more complex resource usage patterns for some or all of the operations.

Needless to say, real machines have finite resources, generally have at least a few operations that have latencies greater than one cycle, and often have at least a few operations with complex usage patterns. We believe that the value of a scheduling algorithm is proportional to the degree of realism of the assumed machine model.

Finally, the scheduling algorithm can also be categorized by the nature of the process involved in generating the schedule. At one extreme are one-pass algorithms that schedule each operation once and for all. At the other extreme are algorithms that perform an exhaustive, branch-and-bound style of search for the schedule. In between, is a spectrum of possibilities such as iterative but non-exhaustive search algorithms or incremental algorithms that make a succession of elementary perturbations to an existing legal schedule to nudge it toward the final solution. This aspect of the scheduling algorithm is immensely important in practice. The further one diverges

from a one-pass algorithm, the slower the scheduler gets until, eventually, it is unacceptable in a real-world setting.

Local scheduling

Scheduling, as a part of the code generation process, was first studied extensively in the context of microprogramming. Local scheduling is concerned with generating as short a schedule as possible for the operations within a single basic block; in effect a scheduling barrier is assumed to exist between adjacent basic blocks in the control flow graph. Although it was typically referred to as local code compaction², the similarity to the job of scheduling tasks on processors [108-117] was soon understood and a number of notions and algorithms from scheduling theory were borrowed by the microprogramming community. Attempts at automating this task have been made since at least the late 1960s [118-125, 91, 126, 127, 15]. Since scheduling is known to be NP-complete [115], the initial focus was on defining adequate heuristics [118, 128, 129, 116, 130, 91]. The consensus was that list scheduling using the highest-level-first priority scheme [113, 91] is relatively inexpensive computationally (a one-pass algorithm) and near-optimal most of the time. Furthermore, this algorithm has no difficulty in dealing with non-unit execution latencies.

The other dimension in which local scheduling matured was in the degree of realism of the machine model. From an initial model in which each operation used a single resource for a single cycle (the simple resource usage model) and had unit latency, algorithms for local scheduling were gradually generalized to cope with complex resource usage and arbitrary latencies [128, 131, 120, 132, 121, 129, 130] culminating in the fully general resource usage “microtemplate” model proposed by Tokoro, et al. [133], and which was known in the hardware pipeline design field as a reservation table [134]. In one form or another, this is now the commonly used machine model in serious instruction schedulers. This machine model is quite compatible with the highest-level-first list scheduling algorithm and does not compromise the near-optimality of this algorithm [15].

² We shall consistently refer to this code generation activity as scheduling.

Global acyclic scheduling

A number of studies have established that basic blocks are quite short--typically about 5-20 instructions on the average. So, whereas local scheduling can generate a near-optimal schedule, data dependences and execution latencies conspire to make the optimal schedule, itself, rather disappointing in terms of its speedup over the original sequential code. Further improvements require overlapping the execution of successive basic blocks, which is achieved by global scheduling.

Early strategies for global scheduling attempted to automate and emulate the ad hoc techniques that hand coders practiced of first performing local scheduling of each basic block and then attempting to move operations from one block into an empty slot in a neighboring block [135, 133]. The shortcoming of such an approach is that, during local compaction, too many arbitrary decisions have already been made which failed to take into account the needs of and opportunities in the neighboring blocks. Many of these decisions might need to be undone before the global schedule can be improved.

In one very important way, the mindset inherited from microprogramming was an obstacle to progress in global scheduling. Traditionally, code compaction was focused on the objective of reducing the size of the microprogram so as to allow it to fit in the microprogram memory. In the case of individual basic blocks, the objectives of local compaction and local scheduling are aligned. This alignment of objectives is absent in the global case. Whereas global code compaction wishes to minimize the sum of the code sizes for the individual basic blocks, global scheduling must attempt to minimize the total execution time of all the basic blocks. In other words, global scheduling must minimize the sum of the code sizes of the individual basic blocks *weighted by the number of times each basic block is executed*. Thus, effective global scheduling might actually increase the size of the program by greatly lengthening an infrequently visited basic block in order to slightly reduce the length of a high frequency basic block. This difference between global compaction and global scheduling, which was captured neither by the early ad hoc techniques nor by the syntactically-driven hierarchical reduction approach proposed by Wood [136], was noted by Fisher [91, 15].

Furthermore, the focus of Fisher's work was on reducing the length of those *sequences* of basic blocks that are frequently executed by the program. These concepts were captured by Fisher in the global scheduling algorithm known as **trace scheduling** [91, 15]. The procedure is centered

around the concept of a trace, which is an acyclic sequence of basic blocks embedded in the control flow graph. i.e., a path through the program that could conceivably be taken for some set of input data. Traces are selected and scheduled in order of their frequency of execution. The next trace to be scheduled is defined by selecting the highest frequency basic block that has not yet been scheduled as the seed of the trace. The trace is extended forward along the highest frequency edge out of the last block of the trace as long as that edge is also the most frequent edge into the successor block, and as long as the successor block is not already part of the trace. Likewise, the trace is extended backwards, as well, from the seed block. The selected trace is then scheduled as if it were a single basic block, i.e., giving no special consideration to branches, except that they are constrained to remain in their original order. Implicit in the resulting schedule is inter-block code motion along the trace in either the upward or downward direction. Matching off-trace code motions must be performed as prescribed by the rules of inter-block code motion specified by Fisher. This activity is termed bookkeeping. Thereafter, the next trace is selected and scheduled. This procedure is repeated until the entire program has been scheduled. The key property of trace scheduling is that, unlike previous approaches to global scheduling, the decisions as to whether to move an operation from one block to another, where to schedule it, and which register to allocate to hold its result (see Section 3.2.2 below) are all made jointly rather than in distinct compiler phases.

Fisher and his co-workers at Yale went on to implement trace scheduling in the Bulldog compiler as part of the ELI project [14, 70]. This trace scheduling implementation and other aspects of the Bulldog compiler have been extensively documented by Ellis [92]. The motion of code downwards across branches and upwards across merges results in code replication. Although this is generally acceptable as the price to be paid for better global schedules, Fisher recognized the possibility that the greediness of highest-level-first list scheduling could sometimes cause more code motion and, hence, replication than is needed to achieve a particular schedule length [15]. Su and his colleagues have recommended certain heuristics for the list scheduling of traces to address this problem [137-139]. Experiments over a limited set of test cases indicate that these heuristics appear to have the desired effect.

The research performed in the ELI project formed the basis of the production-quality compiler that was built at Multiflow. One of the enhancements to trace scheduling implemented in the Multiflow compiler was the elimination of redundant copies of operations caused by bookkeeping. When an off-trace path, emanating from a branch on the trace, rejoins the trace lower down, an operation that is moved above the rejoin and all the way to a point above the

branch can make the off-trace copy redundant under the appropriate circumstances. The original version of trace scheduling, oblivious to such situations, retains two copies of the operation. Gross and Ward describe an algorithm to avoid such redundancies [140]. Freudenberger and Ruttenberg discuss the integrated scheduling and register allocation in the Multiflow compiler [141]. Lowney, et al., provide a comprehensive description of the Multiflow compiler [142].

Hwu and his colleagues on the IMPACT project have developed a variant of trace scheduling that they term superblock scheduling [143, 144]. In an attempt to facilitate the task of incorporating profile-driven global scheduling into more conventional compilers, they separate the trace selection and code replication from the actual scheduling and bookkeeping. To do this, they limit themselves to only moving operations up above branches, never down, and never up past merges. To make this possible, they outlaw control flow into the interior of a trace by means of tail duplication, i.e., creating a copy of the trace below the entry point and redirecting the incoming control flow path to that copy. Once this is done for each incoming path, the resulting trace consists of a sequence of basic blocks with branches out of the trace but no incoming branches except to the top of the trace. This constitutes a superblock, also known as an extended basic block in the compiler literature. Chang and Hwu have studied different trace selection strategies and have measured their relative effectiveness [145]. A comprehensive discussion of the results and insights from the IMPACT project are provided in this special issue [101].

Although the global scheduling of linear sequences of basic blocks represents a major step forward, it has been criticized for its total focus on the current trace and neglect of the rest of the program. For instance, if there are two equally frequent paths through the program, that have basic blocks in common, it is unclear as part of which trace these blocks should be scheduled. One solution is to replicate the code as is done for superblock scheduling. The other is to generalize trace scheduling to deal with more general control flow graphs. Linn [146] and Hsu [103] have proposed profile-driven algorithms for scheduling trees of basic blocks in which all but the root basic block have a single incoming path. Nicolau attempted to extend global scheduling to arbitrary, acyclic control flow graphs using percolation scheduling [147, 87]. However, since percolation scheduling assumes unbounded resources, it cannot realistically be viewed as a scheduling algorithm. Percolation scheduling was then extended to non-unit execution latencies (but still with unbounded resources) [148].

The development of practical algorithms for the global scheduling of arbitrary, acyclic control flow graphs is an area of active research. Preliminary algorithms, assuming finite resources have

been defined by Ebcioğlu [149, 150] and by Fisher [151]. These are both generalizations of trace scheduling. However, there are very many difficulties in the engineering of a robust and efficient scheduler of this sort. The challenges in this area of research revolve around the finding of pragmatic engineering solutions to these problems.

A rather different approach to global acyclic scheduling has been pursued in the IMPACT project [106]. An arbitrary, acyclic control flow graph, having a single entry can be handled by this technique. The control flow graph is if-converted [107, 152] so as to eliminate all branches internal to the flow graph. The resulting code, which is similar to a superblock in that it can only be entered at the top but has multiple exits, is termed a hyperblock. This is scheduled in much the same manner as a superblock except that two operations with disjoint predicates (i.e., operations which cannot both be encountered on any single path through the original flow graph) may be scheduled to use the same resources at the same time. After scheduling, reverse if-conversion is performed to regenerate the control flow graph. Portions of the schedule in which m predicates are active yield 2^m versions of the code.

Cyclic scheduling

As with acyclic flow graphs, instruction-level parallelism in loops is obtained by overlapping the execution of multiple basic blocks. With loops, however, the multiple basic blocks are the multiple iterations of the same piece of code. The most natural extension of the previous global scheduling ideas to loops is to unroll the body of the loop some number of times, and to then perform trace scheduling, or some other form of global scheduling, over the unrolled loop body. This approach was suggested by Fisher [153]. A drawback of this approach is that no overlap is sustained across the back-edge of the unrolled loop. Fisher, et al., went on to propose a solution to this problem, which is to continue unrolling and scheduling successive iterations until a repeating pattern is detected in the schedule. The repeating pattern can be re-rolled to yield a loop whose body is the repeating schedule. As we shall see, this approach was subsequently pursued by various researchers. In the meantime, loop scheduling moved off in a different direction which, as is true of most VLIW scheduling work, had its roots in hardware design.

Researchers concerned with the design of pipelined functional units, most notably Davidson and his co-workers, had developed the theory of and algorithms for the design of hardware controllers for pipelines to maximize the rate at which functions could be evaluated [134, 154-158]. The issues considered here were quite similar to those faced by individuals programming

the innermost loops of signal processing algorithms [159-164] on the early peripheral array processors [46-48]. In both cases, the objective was to sustain the initiation of successive function evaluations (loop iterations) before prior ones had completed. Since this style of computation is termed pipelining in the hardware context, it was dubbed **software pipelining** in the programming domain [12].

Early work in software pipelining consisted of ad hoc hand-coding techniques [164, 12]. Both the quality of the schedules generated and the attempts at automating the generation of software pipelined schedules were hampered by the architecture of the early array processors. Nevertheless, Floating Point Systems developed, for the FPS-164 array processor, a compiler that could software pipeline a loop consisting of a single basic block [165]. Weiss and Smith [166] note that a limited form of software pipelining was present both in certain hand-coded libraries for the CDC 6600 and also as a capability in the FORTRAN compiler for the CDC 6600.

The general formulation of the software pipelining process for single basic block loops was stated by Rau, et al., [16, 84] drawing upon and generalizing the theory developed by Davidson and his co-workers on the design of hardware pipelines. This work identified the attributes of a VLIW architecture that make it amenable to software pipelining, most importantly, the availability of conflict-free access to register storage between the output of a functional unit producing a result and the functional unit that uses that result. This provides freedom in scheduling each operation and is in contrast to the situation in array processors where, due to limited register file bandwidth, achieving peak performance required that a majority of the operations be scheduled to start at the same instant that their predecessor operations completed so that they could pluck their operands right off the result buses.

Rau, et al., also presented a condition which has to be met by any legal software pipelined schedule--the "modulo constraint"--and derived lower bounds on the rate at which successive iterations of the loop can be started, i.e., the **initiation interval (II)**. (II is also the length of the software pipelined loop, measured in VLIW instructions, when no loop unrolling is employed.) This lower bound on II, the **minimum initiation interval (MII)**, is the maximum of the lower bound due to the resource usage constraints (**ResMII**) and the lower bound due to the cyclic data dependence constraints caused by recurrences (**RecMII**). This lower bound is applicable both to vectorizable loops as well as those with arbitrary recurrences and for operation latencies of arbitrary length. A simple, deterministic software pipelining algorithm based on list scheduling,

the modulo scheduling algorithm, was shown to achieve the MII, thereby yielding an asymptotically optimal schedule. This algorithm was restricted to DO-loops whose body is a single basic block being scheduled on a machine in which each operation has a simple pattern of resource usage, viz., the resource usage of each operation can be abstracted to the use of a single resource for a single cycle (even though the latency of the operation is not restricted to a single cycle). The task of generating an optimal, resource-constrained schedule for loops with arbitrary recurrences is known to be NP-complete [167, 168] and any practical algorithm must utilize heuristics to guide a generally near-optimal process. These heuristics were only broadly outlined in this work.

Three independent sets of activity took this work and extended it in various directions. The first one was the direct continuation at Cydrome, over the period 1984-1988, of the work done by Rau, et al. [104, 105]. In addition to enhancing the modulo scheduling algorithm to handle loops with recurrences and arbitrary acyclic control flow in the loop body, attention was paid to coping with the very complex resource usage patterns that were the result of compromises forced by pragmatic implementation considerations. Complex recurrences and resource usage patterns make it unlikely that a one-pass scheduling algorithm, such as list scheduling, will be able to succeed in finding a near-optimal modulo schedule, even when one exists, and performing an exhaustive search was deemed impractical. Instead, an iterative scheduling algorithm was employed, that could unreschedule and reschedule operations. This iterative algorithm is guided by heuristics based on dynamic slack-based priorities. The initial attempt is to schedule the loop with the Π equal to the MII. If unsuccessful, the Π is incremented until a modulo schedule is achieved.

Loops with arbitrary acyclic control flow in the loop body are dealt with by performing IF-conversion [107] to replace all branching by predicated (guarded) operations. This transformation, which assumes the hardware capability of predicated execution [51, 18], yields a loop with a single basic block which is then amenable to the modulo scheduling algorithm [104]. A disadvantage of predicated modulo scheduling is that the ResMII must be computed as if all the operations in the body of the loop are executed each iteration whereas, in reality, only those along one of the control flow paths are actually executed. As a result, during execution, some fraction of the operations in an instruction are wasted. Likewise, the RecMII is determined by the worst-case dependence chain across all paths through the loop body. Both contribute to a degree of sub-optimality that depends on the structure of the loop.

Assuming the existence of hardware to support both predicated execution and speculative execution [93], Cydrome's modulo scheduling algorithm has been further extended to handle WHILE-loops and loops with conditional exits [169]. The problem that such loops pose is that it is not known until late in one iteration whether the next one should be started. This eliminates much of the overlap between successive iterations. The solution is to start iterations speculatively, in effect, by moving operations from one iteration into a prior one. The hardware support makes it possible to avoid observing exceptions from operations that should not have been executed, without overlooking exceptions from non-speculative operations.

Independently of the Cydrome work, Hsu proposed a modulo scheduling algorithm for single basic block loops with general recurrences that recognizes each strongly-connected class (SCC) of nodes in the cyclic dependence graph as a distinct entity [167]. Once the nodes in all the SCCs have been jointly scheduled at the smallest possible Π using a combinatorial search, the nodes in a given SCC may only be re-scheduled as a unit and at a time that is displaced by a multiple of Π . This re-scheduling is performed to enable the remaining nodes, that are not part of any SCC, to be inserted into the schedule. Hsu also described an Π extension technique which can be used to take a legal modulo schedule for one iteration and trivially convert it into a legal modulo schedule for a larger Π without performing any scheduling. This works with simple resource usage patterns. With complex patterns, a certain amount of re-scheduling would be required, but less than starting from scratch.

Lam's algorithm, too, utilizes the SCC structure but list schedules each SCC separately, ignoring the inter-iteration dependences [168, 170]. Thereafter, an SCC is treated as a single pseudo-operation with a complex resource usage pattern, employing the technique of hierarchical reduction proposed by Wood [136]. After this hierarchical reduction has been performed, the dependence graph of the computation is acyclic and can be scheduled using modulo scheduling. With an initial value equal to the $M\Pi$, the Π is iteratively increased until a legal modulo schedule is obtained. By determining and fixing the schedule of each SCC in isolation, Lam's algorithm can result in SCCs that cannot be scheduled together at the minimum achievable Π .

On the other hand, the application of hierarchical reduction enables Lam's algorithm to cope with loop bodies containing structured control flow graphs without any special hardware support such as predicated execution. Just as with the SCCs, structured constructs such as IF-THEN-ELSE are list scheduled and treated as atomic objects. Each leg of the IF-THEN-ELSE is list scheduled separately and the union of the resource usages represents that of the reduced IF-THEN-ELSE

construct. This permits loops with structured flow of control to be modulo scheduled. After modulo scheduling, the hierarchically reduced IF-THEN-ELSE pseudo-operations must be expanded. Each portion of the schedule in which m IF-THEN-ELSE pseudo-operations are active must be expanded into 2^m control flow paths with the appropriate branching and merging between the paths.

Since Lam takes the union of the resource usages in a conditional construct while predicated modulo scheduling takes the sum of the usages, the former approach should yield the smaller MII. However, since Lam separately list schedules each leg of the conditional creating pseudo-operations with complex resource usage patterns, the II that she actually achieves should deviate from the MII to a greater extent. Warter, et al., have implemented both techniques and have observed that, on the average, Lam's approach results in smaller MII's but larger II's [171]. This effect increases for processors with higher issue rates. Warter, et al., go on to combine the best of both approaches in their enhanced modulo scheduling algorithm. They derive the modulo schedule as if predicated execution were available, except that two operations from the same iteration are allowed to be scheduled on the same resource at the same time if their predicates are mutually exclusive, i.e., they cannot both be true. This is equivalent to taking the union of the resource usages. Furthermore, it is applicable to arbitrary, possibly unstructured, acyclic flow graphs in the loop body. After modulo scheduling, the control flow graph is re-generated much as in Lam's approach. Enhanced modulo scheduling results in MII's that are as small as for hierarchical reduction, but as with predicated modulo scheduling, the achieved II is rarely more than the MII.

Yet another independent stream of activity has been the work of Su and his colleagues [138, 172]. When limited to loops with a single basic block, Su's URPR algorithm is an ad hoc approximation to modulo scheduling and is susceptible to significant sub-optimality when confronted by non-unit latencies and complex resource usage patterns. The essence of the URPR algorithm is to unroll and schedule successive iterations until the first iteration has completed. Next the smallest contiguous set of instructions, that contain at least one instance of each operation in the original loop, is identified. After deleting multiple instances of all operations, this constitutes the software pipelined schedule. This deletion process introduces "holes" in the schedule and the attendant sub-optimality. Also, for non-unit latencies, there is no guarantee that the schedule, as constructed, can loop back on itself without padding the schedule out with no-op cycles. This introduces further degradation.

Subsequently, Su extended URPR to the GURPR* algorithm for software pipelining loops containing control flow [173-175]. GURPR* consists of first performing global scheduling on the body of the loop and then using a URPR-like procedure, as if each iteration was IF-converted, to derive the repeating pattern. Finally, as with enhanced modulo scheduling, a control flow graph is regenerated. The shortcomings of URPR are inherited by GURPR*. Wharter, et al., who have implemented GURPR* within the IMPACT compiler, have found that GURPR* performs significantly worse than hierarchical reduction, predicated modulo scheduling or enhanced modulo scheduling [171].

The idea proposed by Fisher, et al., of incrementally unrolling and scheduling a loop until the pattern repeats [153] was pursued by Nicolau and his co-workers, assuming unbounded resources, initially for single basic block loops [176] and then, under the title of perfect pipelining, for multiple basic block loops [177, 148]. The latter was subsequently extended to yield a more realistic algorithm assuming finite resources [178]. For single basic block loops, the incremental unrolling yields a growing linear trace, the expansion of which is terminated once a repeating pattern is observed. In practice, there are complications since the various SCCs might proceed at different rates, never yielding a repeating pattern. For multiple basic block loops, the unrolling yields a growing tree of schedules, each leaf of which spawns two further leaves when a conditional branch is scheduled. A new leaf is not spawned if the (infinite) tree, of which it is a root, is identical to another (infinite) tree (of which it might be a sub-tree) whose root has already been generated.

This approach addresses a shortcoming of all the previously mentioned approaches to software pipelining multiple basic block loops. In general, both RecMII and ResMII are dependent upon the specific control flow path followed in each iteration. Whereas the previous approaches had to use a single, constant, conservative value for each one of these lower bounds, the unrolling approach is able to take advantage of the branch history of previous iterations in deriving the schedule for the current one. However, there are some drawbacks as well. One handicap that such unrolling schemes have is a lack of control over the greediness of the process of initiating iterations. Starting successive iterations as soon as possible, rather than at a measured rate that is in balance with the completion rate, cannot reduce the average initiation interval but can increase the time to enter the repeating pattern and the length of the repeating pattern. Both contribute to longer compilation times and larger code size. A second problem with unrolling schemes lies in their implementation; recognizing that one has arrived at a previously visited state, to which one

can wrap back instead of further expanding the search tree, is quite complicated, especially in the context of finite resources, non-unit latencies and complex resource usage patterns.

The cyclic scheduling algorithm developed by the IBM VLIW research project [179-181, 150] might represent a good compromise between the ideal and the practical. Stripped to the essentials, this algorithm defines a cut-set, termed a fence, to the cyclic graph which yields an acyclic graph. This reduces the problem to that of scheduling a general, acyclic graph-- a simpler problem. Once this is done, the fence is moved and the acyclic scheduling is repeated. As this process is repeated, all the cycles in the control flow graph acquire increasingly tight schedules. The acyclic scheduling algorithm used by Ebcioğlu, et al., is a resource-constrained version of percolation scheduling [149, 150].

Software pipelining was also implemented in the compiler for the product line marketed by another mini-supercomputer company, Culler Scientific. Unfortunately, we do not believe that any publication exists, describing their implementation of software pipelining. Quite recently, software pipelining has been implemented in the compilers for HP's PA-RISC line of computers [182].

Scheduling for RISC and superscalar processors

Seemingly conventional scalar processors can sometimes benefit from scheduling techniques. This is due to small amounts of ILP in the form of, for instance, branch delay slots and shallow pipelines. Scheduling for such processors, whether RISC or CISC, has generally been less ambitious and more ad hoc than that for VLIW processors [183-185, 98, 186]. This was a direct consequence of the lack of parallelism in those machines and the corresponding lack of opportunity for the scheduler to make a big difference. Furthermore, the limited number of registers in those architectures made the use of aggressive scheduling rather unattractive. As a result, scheduling was viewed as rather peripheral to the compilation process, in contrast to the central position it occupied for VLIW processors and, to a lesser extent, for more highly pipelined processors [187, 188, 166]. Now, with superscalar processors growing in popularity, the importance of scheduling, as a core part of the compiler, is better appreciated and a good deal of activity has begun in this area [189-193], unfortunately, sometimes unaware of the large body of literature that already exists.

3.2.2 Register Allocation

In conventional, sequential processors, instruction scheduling is not an issue. The program's execution time is barely affected by the order of the instruction, only by the number of instructions. Accordingly, the emphasis of the code generator is on generating the minimum number of instructions and using as few registers as possible [194-199]. However, in the context of pipelined or multiple-issue processors, where instruction scheduling is important, the issue of the phase-ordering between it and register allocation has been a topic of much debate. There are advocates both for performing register allocation before scheduling [185, 200, 192] as well as for performing it after scheduling [183, 201-203]. Each phase-ordering has its advantages and neither one is completely satisfactory.

The most important argument in favor of performing register allocation first is that whereas a better schedule may be desirable, code that requires more registers than are available is just unacceptable. Clearly, achieving a successful register allocation must supersede the objective of constructing a better schedule. The drawback of performing scheduling first, oblivious of the register allocation, is that shorter schedules tend to yield greater register pressure. In the event that a viable allocation cannot be found, spill code must be inserted. At this point, in the case of a statically scheduled processor, the schedule just constructed may no longer be correct. Even if it is, it may be far from the best one possible, for either a VLIW or superscalar machine, since the schedule was built without the spill code in mind. In machines whose load latency is far greater than that of the other operations, the time penalty of the spill code may far exceed the benefits of the better schedule obtained by performing scheduling first.

Historically, the merit of performing register allocation first was that processors had little instruction-level parallelism and few registers. So, whereas there was much to be lost by a poor register allocation, there was little to be gained by good scheduling. It was customary, therefore, to first perform register allocation first, for instance using graph coloring [204-206] followed by a post-pass scheduling step that considered individual basic blocks [185, 200].

From the viewpoint of instruction-level parallel machines, the major problem with performing register allocation first is that it introduces anti-dependences and output dependences that can constrain parallelism and the ability to construct a good schedule. To some extent this is inevitable; the theoretically optimal combination of schedule and allocation might contain additional arcs due to the allocation. The real concern is that, when allocation is done first, an

excessive number of ill-advised and unnecessary arcs might be introduced due to the insensitivity of the register allocator to the scheduling task. On pipelined machines, whose cache access time is as short as or shorter than the functional unit latencies, the benefits of a schedule unconstrained by register allocation may out-weigh the penalties of the resulting spill code.

Scheduling prior to register allocation, known as pre-pass scheduling, was used in the PL.8 compiler [183]. In evolving this compiler to become the compiler for the superscalar IBM RISC System/6000, the sub-optimality of inserting spill code after the creation of the schedule became clear and a second, post-pass scheduling step was added after the register allocation [202]. During the post-pass, the scheduler honors all the dependences caused by the register allocation which, in turn, was aware of the preferred instruction schedule provided by the pre-pass scheduler. The IMPACT project at the University of Illinois has demonstrated the effectiveness of this strategy for multiple-issue processors [203]. Instead of employing the graph coloring paradigm, Hendren, et al., make use of the richer information present in interval graphs, which are a direct temporal representation of the span of the lifetimes [207]. This assumes that the schedule or, at least, the instruction order has already been determined and that a post-pass scheduling step will follow.

Irrespective of which one goes first, a shortcoming of all strategies discussed so far is that the first phase makes its decisions with no consideration of their impact on the subsequent phase. Goodman and Hsu have addressed this problem by developing two algorithms--one, a scheduler that attempts to keep the register pressure below a limit provided to it, and the second, a register allocation algorithm that is sensitive to its effect on the critical path length of the DAG and, thus, to the effect on the eventual schedule [201].

For any piece of code on a given processor, there is some optimal schedule for which register allocation is possible. Scheduling twice, once before and then after register allocation, is an approximation to achieving this ideal. Simultaneous scheduling and register allocation is another strategy for attempting to find a near-optimal schedule and register allocation. Simultaneous scheduling and register allocation is currently understood only in the context of acyclic code, specifically, a single basic block or a linear trace of basic blocks. The essence of the idea is that each time an operation is scheduled, an available register is allocated to hold the result. Also, if this operation constitutes the last use of the contents of one of the source registers, that register is made available once again for subsequent allocation. When no register is available to receive the result of the operation being scheduled, a register must be spilled. The register holding the datum

whose use is furthest away in the future is spilled. This approach was used in the FPS-164 compiler at the level of individual basic blocks [165] as well as across entire traces [92, 141, 142]. An important concept developed by the ELI project at Yale and by Multiflow was that of performing hierarchical, profile-driven, global integrated scheduling and register allocation. Traces are picked in decreasing order of frequency and integrated scheduling and allocation is performed on each. The scheduling and allocation decisions made for traces that have been processed form constraints on the corresponding decisions for the remaining code. This is a far more systematic approach than other ad hoc, priority-based schemes with the same objective. A syntax-based hierarchical approach to global register allocation has been suggested by Callahan and Koblenz [208].

If a loop is unrolled some number of times and then treated as a linear trace of basic blocks [153], simultaneous trace scheduling and register allocation can be accomplished, but with some loss of performance due to the emptying of pipelines across the back-edge. In the case of modulo scheduling, which avoids this performance penalty, no approach has yet been advanced for simultaneous register allocation. Since doing register allocation in advance is unacceptably constraining on the schedule, it must be performed following modulo scheduling. A unique situation encountered with modulo scheduled loops is that the lifetimes are often much longer than the initiation interval. Normally, this would result in a value being over-written before its last use has occurred. One solution is to unroll the kernel of modulo scheduled loop a sufficient number of times to ensure that no lifetime is longer than the length of the replicated kernel [168, 170]. This is known as modulo variable expansion. In addition to techniques such as graph coloring, the heuristics proposed by Hendren, et al., [207] and by Rau, et al., [209] may be applied after modulo variable expansion. The other solution for register allocation is to assume the dynamic register renaming provided by the rotating register capability of the Cydra 5. The entity that the register allocator works with are vector lifetimes, i.e., the entire sequence of (scalar) lifetimes defined by a particular operation over the execution of the loop [104, 105, 209]. Lower bounds on the number of registers needed for a modulo scheduled loop have been developed by Mangione-Smith, et al. [210]. The strategy for recovering from a situation, in which no allocation can be found for the software pipelined loop, is not well understood. Some options have been outlined [209], but their detailed implementation, effectiveness and relative merits have as yet to be investigated.

3.2.3 Other ILP Compiler Topics

Although scheduling and register allocation are at the heart of ILP compilation, a number of other analyses, optimizations and transformations are crucial to the generation of high quality code. Currently, schedulers treat a procedure call as a barrier to code motion. Thus, inlining of intrinsics and user procedures is very important in the high frequency portions of the program [105, 142, 211].

Certain loop-oriented analyses and optimizations are specific to modulo scheduling. IF-conversion and the appropriate placement of predicate-setting operations is needed to modulo schedule loops with control flow [107, 104, 152, 105]. The elimination of subscripted loads and stores that are redundant across multiple iterations of a loop can have a significant effect upon both the ResMII and the RecMII [212, 213, 105]. This is important for trace scheduling unrolled loops as well [142]. Recurrence back-substitution and other transformations that reduce the RecMII have a major effect on the performance of all software pipelined loops [105]. Most of these transformations and analyses are facilitated by the dynamic single assignment representation for inner loops [213, 105].

On machines with multiple, identical clusters, such as the Multiflow TRACE machines, it is necessary to decide which part of the computation will go on each cluster. This is a non-trivial task; whereas increased parallelism argues in favor of spreading the computation over the clusters, this also introduces inter-cluster move operations into the computation, whose latency can degrade performance if the partitioning of the computation across clusters is not done carefully. An algorithm for performing this partitioning was developed by Ellis [92] and was incorporated into the Multiflow compiler [142].

An issue of central importance to all ILP compilation is the disambiguation of memory references., i.e., deciding whether two memory references definitely are to the same memory location or definitely are not. Known as dependence analysis, this has become a very well developed topic in the area of vector computing over the past twenty years [214]. For vector computers, the compiler is attempting to prove that two references in different iterations are *not* to the same location. No benefit is derived if it is determined that they *are* to the same location since such loops cannot be vectorized. Consequently, the nature of the analysis, especially in the context of loops containing conditional branching, has been approximate. This is a shortcoming from the point of view of ILP processors which can benefit both if the two references are or are

not to the same location. A more precise analysis than dependence analysis, involving data flow analysis, is required. Also, with ILP processors, memory disambiguation is important outside of loops as well as within them. Memory disambiguation within traces was studied in the ELI project [215, 92] and was implemented in the Multiflow compiler [142]. Memory disambiguation, in the context of innermost loops was implemented in the Cydra 5 compiler [213, 105] and was studied by Callahan and Koblenz [212].

4 Available ILP

4.1 Limit Studies and Their Shortcomings

Many experimenters have attempted to measure the maximum parallelism available in programs. The goal of such limit studies is to:

Throw away all considerations of hardware and compiler practicality and measure the greatest possible amount of ILP inherent in a program.

Limit studies are simple enough to describe: take an execution trace of the program, and build a data precedence graph on the operations, eliminating false anti-dependences caused by the write-after-read usage of a register or other piece of hardware storage. The length in cycles of the serial execution of the trace gives the serial execution time on hardware with the given latencies. The length in cycles of the critical path through the data dependence graph gives the shortest possible execution time. The quotient of these two is the available speedup. (In practice, an execution trace is not always gathered. Instead, the executed stream is processed as the code runs, greatly reducing the computation and/or storage required.)

These are indeed maximum parallelism measures in some sense, but they have a critical shortcoming that causes them to miss accomplishing their stated goal; they do not consider transformations that a compiler might make to enhance ILP. Although we mostly mean transformations of a yet-unknown nature that researchers may develop in the future, even current state-of-the-art transformations are rarely reflected in limit studies. Thus we have had, in recent years, the anomalies of researchers stating an "upper limit" on available parallelism in programs that is lower than what has already been accomplished with those same programs, or of new results that show the maximum available parallelism to be significantly higher than it was a few years ago, before a new set of code transformations was considered.

There is a somewhat fatuous argument that demonstrates just how imprecise limit studies must be; recalling that infinite hardware is available, we can replace computations in the code with table look-ups. In each case, we will replace a longer--perhaps very long--computation with one that takes a single step. While this is obviously impractical for most computations with operands that span the (finite, but large) range of integers or floating point numbers representable on a system, it is only impractical in the very sense in which practicality is to be discarded in limit studies. And even on practicality grounds, one cannot dismiss this argument completely; in a sense it really does capture what is wrong with these experiments. There are many instances of transformations, some done by hand, others automatically, that reduce to this concept. Arithmetic and transcendental functions are often sped up significantly by the carefully selected use of table lookups at critical parts of the computation. Modern compilers can often replace a nested set of IF-THEN tests with a single lookup in which hardware does an indirect jump through a lookup table. Limit studies have no way of capturing these transformations, the effect of which could be a large improvement in available ILP.

Even in current practice, the effect of ignoring sophisticated compiling is extreme. Transformations such as tree height reduction, loop conditioning, loop exchange, etc. can have a huge effect on the parallelism available in code. A greater unknown is the research future of data structure selection to improve ILP. A simple example can show this effect. The following code finds the maximum element of a linked list of data.

```
    this-ptr = head-ptr;
    max-so-far = most-neg-number;
    while this-ptr {
        if this-ptr.data > max-so-far
            then max-so-far = this-ptr.data;
        this-ptr = this-ptr.next }
```

From simple observation, the list of elements chained from head-ptr cannot be circular. If the compiler had judged it worthwhile, it could have stored these elements in an array and done the comparisons pairwise, in parallel, without having to chase the pointers linearly. This example is not as far-fetched as it might seem. Vectorization took 20 years to go from the ability to recognize the simplest loop to the sophisticated vectorizers we have today. There has been virtually no work done on compiler transformations to enhance ILP.

Limit studies, then, are in some sense finding the maximum parallelism available, but in other ways are finding the minimum. In these senses, they find the maximum parallelism:

- o Disambiguation can be done perfectly, well beyond what is practical.
- o There are infinitely many functional units available.
- o There are infinitely many registers available.
- o Rejoins can be completely unwound.

In other senses, they represent a minimum, or an existence proof that at least a certain amount of parallelism exists, since potentially important processes have been left out:

- o Compiler transformations to enhance ILP have not been done.
- o Intermediate code generation techniques which boost ILP have not been done.

Perhaps it is more accurate to say that a limit study shows that the maximum parallelism available, in the absence of practicality considerations, is *at least* the amount measured.

Early Experiments

The very first ILP limit studies demonstrated the effect we wrote of above: the experimenters' view of the techniques by which one could find parallelism was limited to the current state of the art, and they missed a technique that is now known to provide most of the available ILP: the motion of operations between basic blocks of code. Experiments done by Tjaden and Flynn [76] and by Foster and Riseman [216] (and, anecdotally, elsewhere) found that there was only a small amount (about of factor of 2-3) of improvement due to ILP available in real programs. This was dubbed the **Flynn bottleneck**. By all accounts, these pessimistic and, in a sense, erroneous experiments had a tremendous damping effect on the progress of ILP research. The experiments were only erroneous in the sense of missing improvements; certainly they did correctly what they said they did.

Interestingly, one of the research teams doing these experiments saw that under the hypothesis of free and infinite hardware, one wouldn't necessarily have to stop finding ILP at basic block boundaries. In a companion paper to the one mentioned above, Riseman and Foster [217] put forward a hardware-intensive solution to the problem of doing operations speculatively: they

measured what would happen if one used duplicate hardware at conditional jumps, and disregarded the one that went in the wrong direction. They found a far larger amount of parallelism: indeed, they found more than an order of magnitude more than they could when branches were a barrier. Some of the programs they measured could achieve arbitrarily large amounts of parallelism, depending only on data set size. But in an otherwise insightful and visionary piece of work, the researchers lost sight of the fact that they were doing a limit study, and in their tone and abstract emphasized how impractical it would be to implement the hardware scheme they had suggested. (They found that to get a factor of 10 ILP speedup, one had to be prepared to cope with 16 unresolved branches at the worst point of a typical program. Their scheme would require, then, 2^{16} sets of hardware to do so. Today, as described in most of the papers in this issue, we try to get much of the benefit of the same parallelism without the hardware cost by doing code motions that move operations between blocks and having the code generator make sure that the correct computation is ultimately done once the branches settle.)

Contemporary Experiments

We know of no other ILP limit studies published between then and the 1980s. In 1981, Nicolau and Fisher [218, 219] used some of the apparatus being developed for the Yale Bulldog Compiler to repeat the experiment done by Riseman and Foster, and found virtually the same results.

In the late 1980s, architects began to look at superscalar microprocessors and again started a series of limit studies. Interestingly, the most notorious of these [220] again neglected the possibility of code motions between blocks. Unsurprisingly, the Flynn bottleneck appeared again, and only the factor of 2-3 parallelism found earlier was found. Two years later, Wall [221] did the most thorough limit study to date, and accounted for speculative execution, memory disambiguation, and other factors. He built an elaborate model, and published available ILP speedup under a great many scenarios, yielding a wealth of valuable data, but no simple answers. The various scenarios allow one to try to bracket what really might be practical in the near future, but are subject to quite a bit of interpretation. In examining the various scenarios presented, we find that settings that a sophisticated compiler might approach during the coming decade could yield speedups ranging from 7-60 on the sample programs, which are taken from the SPEC suite and other standard benchmarks. (It's worth noting that Wall himself is much more pessimistic. In the same results he sees an average ceiling of about 5, and the near impossibility of attaining even that much.) Lam and Wilson [222] did an experiment to measure

the effects of different methods of eliminating control flow barriers to parallelism. When their model agreed with Wall's, their results were similar. Butler and Patt [223] considered models with a large variety of numbers of functional units, and found that with good branch prediction schemes and speculative execution, a wide range of speedup was available.

4.2 Experiments That Measure Attained Parallelism

In contrast to the limit studies, some people have built real or simulated ILP systems, and have measured their speedup against real or simulated non-parallel systems. When simulated systems have been involved, they have been relatively realistic systems, or systems that the researchers have argued would abstract the essence of realistic systems in such a way that the system realities should not lower the attained parallelism. Thus these experiments represent something closer to true lower bounds on available parallelism.

Ellis [92] used the Bulldog Compiler to generate code for a hypothetical machine. His model was unrealistic in several aspects, most notably the memory system, but realistic implementations should have little difficulty exploiting the parallelism he found. Ellis measured the speedups obtained on 12 small scientific programs for both a "realistic" machine (corresponding to one under design at Yale), and an "ideal" machine, with limitless hardware and single-cycle functional units. He found speedups ranging from no speedup to 7.6 times speedup for the real model, and a range of 2.7 to 48.3 for the ideal model.

In this issue there are three papers that add to our understanding of the performance of ILP systems. The paper by Hwu, et al. [101], considers the effect of a realistic compiler which uses superblock scheduling. Lowney, et al. [142] and Schuette and Shen [224] compare the performance of the Multiflow Trace 14/300 with current microprocessors from MIPS and IBM, respectively.

Fewer studies have been done to measure the attained performance of software pipelining. Wharter et al. [171] consider a set of 30 *doall* loops with branches found in the Perfect and SPEC benchmark sets. Relative to a single-issue machine without modulo scheduling, they find a 6 times speedup on a hypothetical 4-issue machine, 10 times speedup on a hypothetical 8-issue machine. Lee, et al. [225], combined superblock scheduling and software pipelining for a machine capable of issuing up to 7 operations per cycle. On a mix of loop-intensive (e.g.,

LINPACK) and "scalar" (e.g., Spice) codes, they found an average of 1-4 operations issued per cycle, with 2-7 operations in flight.

5 An introduction to this Special Issue

In this Special Double Issue of The Journal of Supercomputing we have attempted to capture the most significant work, that took place during the 1980s, in the area of instruction-level parallel processing. The intent is to document both the theory and the practice of ILP computing. Consequently, our emphasis is on projects that resulted in implementations of serious scope, since it is this reduction to practice that exposes the true merit and the real problems of ideas that sound good on paper.

During the 1980s, the bulk of the advances in ILP occurred in the form of VLIW processing and this special issue reflects it with papers on Multiflow's TRACE family and on Cydrome's Cydra 5. The paper by Lowney, et al., [142] provides an overview of the TRACE hardware and an in-depth discussion of the compiler. The paper by Schuette and Shen [224] reports on an evaluation performed by the authors of the TRACE 14/300 and a comparison of it to the superscalar IBM RS/6000. The Cydra 5 effort is documented by two papers: one by Beck, Yen and Anderson [52] on the Cydra 5 architecture and hardware implementation, and the other by Dehnert and Towle [105] on the Cydra 5 compiler. (While reading the descriptions of these large and bulky mini-supercomputers, it is worthwhile to bear in mind that they could easily fit on a single chip in the near future!) The only important superscalar product of the 1980s was Astronautics' ZS-1 mini-supercomputer. Although we wanted to include a paper on it in this special issue, that did not come to pass. The final paper in this special issue reports on IMPACT [101], the most thorough implementation of an ILP compiler that has occurred in academia.

6 References

1. B. E. Carpenter and R. W. Doran (Editor). A. M. Turing's ACE Report of 1946 and Other Papers. (MIT Press, Cambridge, Massachusetts, 1986).
2. M. V. Wilkes. The best way to design an automatic calculating machine. Proc. Manchester University Computer Inaugural Conference (Manchester, England, July 1951), 16-18.
3. M. V. Wilkes and J. B. Stringer. Microprogramming and the design of the control circuits in an electronic digital computer. Proc. The Cambridge Philosophical Society, Part 2 (April 1953), 230-238.
4. J. E. Thornton. Parallel operation in the Control Data 6600. Proc. AFIPS Fall Joint Computer Conference (1964), 33-40.
5. J. E. Thornton. Design of a Computer - The Control Data 6600. (Scott, Foresman and Co., Glenview, Illinois, 1970).
6. (Special issue on the System/360 Model 91). IBM Journal of Research and Development 11, 1 (January 1967).
7. C. G. Bell and A. Newell. Computer Structures: Readings and Examples. 1st edn. (McGraw-Hill, New York, 1971).
8. D. A. Patterson and C. H. Sequin. RISC I: A reduced instruction set VLSI computer. Proc. 8th Annual Symposium on Computer Architecture (Minneapolis, Minnesota, May 1981), 443-450.
9. G. Radin. The 801 minicomputer. Proc. Symposium on Architectural Support for Programming Languages and Operating Systems (Palo Alto, California, March 1982), 39-47.
10. J. Hennessy, N. Jouppi, F. Baskett, T. Gross and J. Gill. Hardware/software tradeoffs for increased performance. Proc. Symposium on Architectural Support for Programming Languages and Operating Systems (Palo Alto, California, March 1982), 2-11.
11. J. E. Smith, G. E. Dermer, B. D. Vanderwarn, S. D. Klinger, C. M. Roszewski, D. L. Fowler, K. R. Scidmore and J. P. Laudon. The ZS-1 central processor. Proc. Second International Conference on Architectural Support for Programming Languages and Operating Systems (Palo Alto, California, October 1987), 199-204.
12. A. E. Charlesworth. An approach to scientific array processing: the architectural design of the AP-120B/FPS-164 Family. Computer 14, 9 (1981), 18-27.
13. B. R. Rau, C. D. Glaeser and E. M. Greenawalt. Architectural support for the efficient generation of code for horizontal architectures. Proc. Symposium on Architectural Support for Programming Languages and Operating Systems (Palo Alto, March 1982), 96-99.
14. J. A. Fisher. Very long instruction word architectures and the ELI-512. Proc. Tenth Annual International Symposium on Computer Architecture (Stockholm, Sweden, June 1983), 140-150.
15. J. A. Fisher. Trace scheduling: a technique for global microcode compaction. IEEE Transactions on Computers C-30, 7 (July 1981), 478-490.

16. B. R. Rau and C. D. Glaeser. Some scheduling techniques and an easily schedulable horizontal architecture for high performance scientific computing. Proc. Fourteenth Annual Workshop on Microprogramming (October 1981), 183-198.
17. R. P. Colwell, R. P. Nix, J. J. O'Donnell, D. B. Papworth and P. K. Rodman. A VLIW architecture for a trace scheduling compiler. IEEE Transactions on Computers C-37, 8 (August 1988), 967-979.
18. B. R. Rau, D. W. L. Yen, W. Yen and R. A. Towle. The Cydra 5 departmental supercomputer: design philosophies, decisions and trade-offs. Computer 22, 1 (January 1989).
19. D. W. Anderson, F. J. Sparacio and R. M. Tomasulo. The System/360 Model 91: machine philosophy and instruction handling. IBM Journal of Research and Development 11, 1 (January 1967), 8-24.
20. R. M. Keller. Look-ahead processors. Computing Surveys 7, 4 (December 1975), 177-196.
21. The Series 10000 Personal Supercomputer: Inside a New Architecture. Publication No. 002402-007 2-88. Apollo Computer, Inc., Chelmsford, Massachusetts, 1988.
22. R. Bahr, S. Ciavaglia, B. Flahive, M. Kline, P. Mageau and D. Nickel. The DN10000TX: a new high-performance PRISM processor. Proc. COMPCON '91 (1991).
23. 80960CA User's Manual. Publication No. 270710-001. Intel Corporation, Santa Clara, California, 1989.
24. (Special issue on IBM RISC System/6000 processor). IBM Journal of Research and Development 34, 1 (January 1990).
25. V. Popescu, M. Schultz, J. Spracklen, G. Gibson, B. Lightner and D. Isaman. The Metaflow architecture. IEEE Micro 11, 3 (June 1991), 10.
26. G. Blanck and S. Krueger. The SuperSPARC™ microprocessor. Proc. COMPCON '92 (1992), 136-141.
27. K. Diefendorff and M. Allen. Organization of the Motorola 88110 superscalar RISC microprocessor. IEEE Micro 12, 2 (April 1992), 40-63.
28. E. DeLano, W. Walker, J. Yetter and M. Forsyth. A high speed superscalar PA-RISC processor. Proc. COMPCON '92 (February 1992), 116-121.
29. Arvind and V. Kathail. A multiple processor dataflow machine that supports generalised procedures. Proc. Eighth Annual Symposium on Computer Architecture (May 1981), 291-302.
30. Arvind and K. Gostelow. The U-interpreter. Computer 15, 2 (February 1982).
31. J. Gurd, C. C. Kirkham and I. Watson. The Manchester prototype dataflow computer. Communications of the ACM 28, 1 (January 1985), 34-52.
32. W. J. Watson. The TI ASC -- a highly modular and flexible super computer architecture. Proc. AFIPS Fall Joint Computer Conference (1972), 221-228.
33. R. G. Hintz and D. P. Tate. Control Data STAR-100 processor design. Proc. COMPCON '72 (September 1972), 1-4.
34. R. M. Russell. The CRAY-1 computer system. Communications of the ACM 21 (1978), 63-72.

35. J. J. Dongarra. A survey of high performance computers. Proc. COMPCON '86 (March 1986), 8-11.
36. P. B. Schneck. Supercomputer Architecture. (Kluwer Academic Publishers, Norwell, Massachusetts, 1987).
37. M. Danelutto and M. Vanneschi. VLIW in-the-large: a model for fine grain parallelism exploitation of distributed memory multiprocessors. Proc. 23rd Annual Workshop on Microprogramming and Microarchitecture (November 1990), 7-16.
38. A. Wolfe and J. P. Shen. A variable instruction stream extension to the VLIW architecture. Proc. Fourth International Conference on Architectural Support for Programming Languages and Operating Systems (Santa Clara, April 1991), 2-14.
39. M. Franklin and G. S. Sohi. The expandable split window paradigm for exploiting fine-grain parallelism. Proc. 19th Annual International Symposium on Computer Architecture (Gold Coast, Australia, May 1992), 58-67.
40. R. M. Tomasulo. An efficient algorithm for exploiting multiple arithmetic units. IBM Journal of Research and Development 11, 1 (January 1967), 25-33.
41. M. Johnson. Superscalar Microprocessor Design. (Prentice-Hall, Englewood Cliffs, New Jersey, 1991).
42. J. E. Smith. Dynamic instruction scheduling and the Astronautics ZS-1. Computer 22, 1 (January 1989), 21-35.
43. N. P. Jouppi. The nonuniform distribution of instruction-level and machine parallelism and its effect on performance. IEEE Transactions on Computers C-38, 12 (December 1989), 1645-1658.
44. U. Schmidt and K. Caesar. Datawave: a single-chip multiprocessor for video applications. IEEE Micro 11, 3 (June 1991), 22.
45. G. M. Papadopoulos and D. E. Culler. Monsoon: an explicit token store architecture. Proc. Seventeenth International Symposium on Computer Architecture (Seattle, Washington, May 1990), 82-91.
46. J. F. Ruggiero and D. A. Coryell. An auxiliary processing system for array calculations. IBM Systems Journal 8, 2 (1969), 118-135.
47. IBM 3838 Array Processor Functional Characteristics. Publication No. 6A24-3639-0, File No. S370-08. IBM Corporation, Endicott, New York, 1976.
48. FPS AP-120B Processor Handbook. Floating Point Systems, Inc., Beaverton, Oregon, 1979.
49. i860 64-Bit Microprocessor Programmer's Reference Manual. Publication No. 240329-001. Intel Corporation, Santa Clara, California, 1989.
50. R. P. Colwell, W. E. Hall, C. S. Joshi, D. B. Papworth, P. K. Rodman and J. E. Tornes. Architecture and implementation of a VLIW supercomputer. Proc. Supercomputing '90 (1990), 910-919.
51. B. R. Rau. Cydra 5 Directed Dataflow architecture. Proc. COMPCON '88 (San Francisco, March 1988), 106-113.
52. G. R. Beck, D. W. L. Yen and T. L. Anderson. The Cydra 5 mini-supercomputer: architecture and implementation. The Journal of Supercomputing (this issue) (1992).

53. J. Labrousse and G. A. Slavenburg. CREATE-LIFE : a design system for high performance VLSI circuits. Proc. International Conference on Circuits and Devices (1988).
54. L. Kohn and N. Margulis. Introducing the Intel i860 64-bit microprocessor. IEEE Micro 9, 4 (August 1989), 15-30.
55. J. Labrousse and G. A. Slavenburg. A 50 MHz microprocessor with a VLIW architecture. Proc. ISSCC '90 (San Francisco, 1990).
56. J. Labrousse and G. A. Slavenburg. CREATE-LIFE : a modular design approach for high performance ASIC's. Proc. COMPCON '90 (San Francisco, 1990).
57. C. Peterson, J. Sutton and P. Wiley. iWarp: a 100-MOPS, LIW microprocessor for multicomputers. IEEE Micro 11, 3 (June 1991), 26.
58. M. R. Thistle and B. J. Smith. A processor architecture for Horizon. Proc. Supercomputing '88 (Orlando, Florida, November 1988), 35-41.
59. L. W. Cotten. Circuit implementation of high-speed pipeline systems. Proc. AFIPS Fall Joint Computing Conference (1965), 489-504.
60. L. W. Cotten. Maximum-rate pipeline systems. Proc. AFIPS Spring Joint Computing Conference (1969), 581-586.
61. T. C. Chen. Parallelism, pipelining, and computer efficiency. Computer Design 10, 1 (January 1971), 69-74.
62. T. G. Hallin and M. J. Flynn. Pipelining of arithmetic functions. IEEE Transactions on Computers C-21, 8 (August 1972), 880-886.
63. B. K. Fawcett. Maximal Clocking Rates for Pipelined Digital Systems. M.S. thesis. University of Illinois, Urbana-Champaign, 1975.
64. S. R. Kunkel and J. E. Smith. Optimal pipelining in supercomputers. Proc. 13th Annual International Symposium on Computer Architecture (Tokyo, Japan, June 1986), 404-411.
65. E. Bloch. The engineering design of the STRETCH computer. Proc. Eastern Joint Computer Conference (1959), 48-59.
66. W. Buchholz (Editor). Planning A Computer System: Project Stretch. (McGraw-Hill, New York, 1962).
67. J. P. Eckert, J. C. Chu, A. B. Tonik and W. F. Schmitt. Design of UNIVAC-LARC System: I. Proc. Eastern Joint Computer Conference (1959), 59-65.
68. T. C. Chen. Overlap and pipeline processing, in Introduction to Computer Architecture. H. S. Stone (Editor). (Science Research Associates, Chicago, Illinois, 1975), 375-431.
69. P. M. Kogge. The Architecture of Pipelined Computers. (McGraw-Hill, New York, 1981).
70. J. A. Fisher, J. R. Ellis, J. C. Ruttenberg and A. Nicolau. Parallel processing: a smart compiler and a dumb machine. Proc. ACM SIGPLAN '84 Symposium on Compiler Construction (Montreal, Canada, June 1984), 37-47.
71. S. Weiss and J. E. Smith. Instruction issue logic for pipelined supercomputers. Proc. 11th Annual International Symposium on Computer Architecture (1984), 110-118.
72. W. W. Hwu and Y. N. Patt. HPSm, a high performance restricted data flow architecture having minimal functionality. Proc. 13th Annual International Symposium on Computer Architecture (Tokyo, Japan, June 1986), 297-306.

73. W. W. Hwu and Y. N. Patt. Checkpoint repair for out-of-order execution machines. IEEE Transactions on Computers C-36, 12 (December 1987), 1496-1514.
74. R. R. Oehler and M. W. Blasgen. IBM RISC System/6000: architecture and performance. IEEE Micro 11, 3 (June 1991), 14.
75. Am29000 Users Manual. Publication No. 10620B. Advanced Micro Devices, Inc., Sunnyvale, California, 1989.
76. G. S. Tjaden and M. J. Flynn. Detection and parallel execution of parallel instructions. IEEE Transactions on Computers C-19, 10 (October 1970), 889-895.
77. G. S. Tjaden and M. J. Flynn. Representation of concurrency with ordering matrices. IEEE Transactions on Computers C-22, 8 (August 1973), 752-761.
78. R. D. Acosta, J. Kjelstrup and H. C. Torng. An instruction issuing approach to enhancing performance in multiple function unit processors. IEEE Transactions on Computers C-35, 9 (September 1986), 815-828.
79. H. Dwyer and H. C. Torng. An out-of-order superscalar processor with speculative execution and fast, precise interrupts. Proc. 25th Annual International Symposium on Microarchitecture (Portland, Oregon, 1992).
80. R. G. Wedig. Detection of Concurrency in Directly Executed Language Instruction Streams. Ph.D. thesis. Stanford University, Stanford, California, 1982.
81. A. K. Uht. An efficient hardware algorithm to extract concurrency from general-purpose code. Proc. Nineteenth Annual Hawaii Conference on System Sciences (January 1986), 41-50.
82. T. Agerwala and J. Cocke. High performance reduced instruction set processors. Technical Report RC12434 (#55845). IBM Thomas J. Watson Research Center, 1987.
83. J. E. Smith. Decoupled access/execute architectures. Proc. Ninth Annual International Symposium on Computer Architecture (April 1982), 112-119.
84. B. R. Rau, C. D. Glaeser and R. L. Picard. Efficient code generation for horizontal architectures: compiler techniques and architectural support. Proc. Ninth Annual International Symposium on Computer Architecture (April 1982), 131-139.
85. J. Hennessy, N. Jouppi, S. Przybylski, C. Rowen, T. Gross, F. Baskett and J. Gill. MIPS: A microprocessor architecture. Proc. 15th Annual Workshop on Microprogramming (Palo Alto, California, October 1982), 17-22.
86. J. A. Fisher. 2^N -way jump microinstruction hardware and an effective instruction binding method. Proc. 13th Annual Workshop on Microprogramming (Colorado Springs, Colorado, November 1980), 64-75.
87. A. Nicolau. Percolation scheduling: a parallel compilation technique. Technical Report TR 85-678. Department of Computer Science, Cornell, 1985.
88. K. Ebcioglu. Some design ideas for a VLIW architecture for sequential-natured software, in Parallel Processing (Proc. IFIP WG 10.3 Working Conference on Parallel Processing, Pisa, Italy), M. Cosnard and e. al. (Editor). (North Holland, 1988), 3-21.
89. G. S. Sohi and S. Vajapayem. Instruction issue logic for high-performance, interruptible pipelined processors. Proc. 14th Annual Symposium on Computer Architecture (Pittsburgh, Pennsylvania, June 1987), 27-36.

90. J. E. Smith and A. R. Pleszkun. Implementing precise interrupts in pipelined processors. IEEE Transactions on Computers C-37, 5 (May 1988), 562-573.
91. J. A. Fisher. The Optimization of Horizontal Microcode Within and Beyond Basic Blocks: An Application of Processor Scheduling with Resources. Ph.D. thesis. New York University, 1979.
92. J. R. Ellis. Bulldog: A Compiler for VLIW Architectures. (The MIT Press, Cambridge, Massachusetts, 1985).
93. S. A. Mahlke, W. Y. Chen, W. W. Hwu, B. R. Rau and M. S. Schlansker. Sentinel scheduling for VLIW and superscalar processors. Proc. Fifth International Conference on Architectural Support for Programming Languages and Operating Systems (Boston, Massachusetts, October 1992), 238-247.
94. M. D. Smith, M. S. Lam and M. A. Horowitz. Boosting beyond static scheduling in a superscalar processor. Proc. Seventeenth International Symposium on Computer Architecture (June 1990).
95. M. D. Smith, M. Horowitz and M. Lam. Efficient superscalar performance through boosting. Proc. Fifth International Conference on Architectural Support for Programming Languages and Operating Systems (Boston, Massachusetts, October 1992), 248-259.
96. J. E. Smith. A study of branch prediction strategies. Proc. Eighth Annual International Symposium on Computer Architecture (May 1981), 135-148.
97. J. K. F. Lee and A. J. Smith. Branch prediction strategies and branch target buffer design. Computer 17, 1 (January 1984), 6-22.
98. S. McFarling and J. Hennessy. Reducing the cost of branches. Proc. Thirteenth International Symposium on Computer Architecture (Tokyo, Japan, June 1986), 396-403.
99. T. Y. Yeh and Y. N. Patt. Alternative implementations of two-level adaptive branch prediction. Proc. Nineteenth International Symposium on Computer Architecture (Gold Coast, Australia, May 1992), 124-134.
100. W. W. Hwu, T. M. Conte and P. P. Chang. Comparing software and hardware schemes for reducing the cost of branches. Proc. 16th Annual International Symposium on Computer Architecture (May 1989), 224-233.
101. W. W. Hwu, S. A. Mahlke, W. Y. Chen, P. P. Chang, N. J. Warter, R. A. Bringmann, R. G. Ouellette, R. E. Hank, T. Kiyohara, G. E. Haab, J. G. Holm and D. M. Lavery. The superblock: an effective technique for VLIW and superscalar compilation. The Journal of Supercomputing (this issue) (1992).
102. J. A. Fisher and S. M. Freudenberger. Predicting conditional jump directions from previous runs of a program. Proc. Fifth International Conference on Architectural Support for Programming Languages and Operating Systems (Boston, Mass., October 1992), 85-95.
103. P. Y. T. Hsu and E. S. Davidson. Highly concurrent scalar processing. Proc. Thirteenth Annual International Symposium on Computer Architecture (1986), 386-395.
104. J. C. Dehnert, P. Y.-T. Hsu and J. P. Bratt. Overlapped loop support in the Cydra 5. Proc. Third International Conference on Architectural Support for Programming Languages and Operating Systems (Boston, Mass., April 1989), 26-38.
105. J. C. Dehnert and R. A. Towle. Compiling for the Cydra 5. The Journal of Supercomputing (this issue) (1992).

106. S. A. Mahlke, D. C. Lin, W. Y. Chen, R. E. Hank and R. A. Bringmann. Effective compiler support for predicated execution using the hyperblock. Proc. 25th Annual International Symposium on Microarchitecture (1992).
107. J. R. Allen, K. K. C. Porterfield and J. Warren. Conversion of control dependence to data dependence. Proc. Tenth Annual ACM Symposium on Principles of Programming Languages (January 1983).
108. T. C. Hu. Parallel sequencing and assembly line problems. Operations Research 9, 6 (1961), 841-848.
109. E. G. Coffman and R. L. Graham. Optimal scheduling for two processor systems. Acta Informatica 1, 3 (1972), 200-213.
110. C. V. Ramamoorthy, K. M. Chandy and M. J. Gonzalez. Optimal scheduling strategies in a multiprocessor system. IEEE Transactions on Computers C-21, 2 (February 1972), 137-146.
111. E. B. Fernandez and B. Bussel. Bounds on the number of processors and time for multiprocessor optimal schedule. IEEE Transactions on Computers C-22, 8 (August 1973), 745-751.
112. K. R. Baker. Introduction to Sequencing and Scheduling. (John Wiley, New York, 1974).
113. T. L. Adam, K. M. Chandy and J. R. Dickson. A comparison of list schedules for parallel processing systems. Communications of the ACM 17, 12 (December 1974), 685-690.
114. W. H. Kohler. A preliminary evaluation of the critical path method for scheduling tasks on multiprocessor systems. IEEE Transactions on Computers C-24, 12 (December 1975), 1235-1238.
115. J. R. Coffman (Editor). Computer and Job-Shop Scheduling Theory. (John Wiley, New York, 1976).
116. M. J. Gonzalez. Deterministic processor scheduling. ACM Computing Surveys 9, 3 (September 1977), 173-204.
117. H. Kasahara and S. Narita. Practical multiprocessor scheduling algorithms for efficient parallel processing. IEEE Transactions on Computers C-33, 11 (November 1984), 1023-1029.
118. C. V. Ramamoorthy and M. J. Gonzalez. A survey of techniques for recognizing parallel processable streams in computer programs. Proc. AFIPS Fall Joint Computing Conference (1969), 1-15.
119. R. L. Kleir and C. V. Ramamoorthy. Optimization strategies for microprograms. IEEE Transactions on Computers C-20, 7 (July 1971), 783-794.
120. M. Tsuchiya and M. J. Gonzalez. An approach to optimization of horizontal microprograms. Proc. Seventh Annual Workshop on Microprogramming (Palo Alto, California, 1974), 85-90.
121. D. J. DeWitt. A control word model for detecting conflicts between microprograms. Proc. 8th Annual Workshop on Microprogramming (Chicago, Illinois, September 1975), 6-12.
122. M. Tsuchiya and M. J. Gonzalez. Toward optimization of horizontal microprograms. IEEE Transactions on Computers C-25, 10 (October 1976), 992-999.

123. T. Agerwala. Microprogram optimization: a survey. IEEE Transactions on Computers C-25, 10 (October 1976), 962-973.
124. M. Tokoro, E. Tamura, K. Takase and K. Tamaru. An approach to microprogram optimization considering resource occupancy and instruction formats. Proc. 10th Annual Workshop on Microprogramming (Niagara Falls, New York, November 1977), 92-108.
125. G. Wood. On the packing of micro-operations into micro-instruction words. Proc. 11th Annual Workshop on Microprogramming (Asilomar, California, November 1978), 51-55.
126. D. Landskov, S. Davidson, B. Shriver and P. W. Mallett. Local microcode compaction techniques. ACM Computing Surveys 12, 3 (September 1980), 261-294.
127. S. Davidson, D. Landskov, B. D. Shriver and P. W. Mallett. Some experiments in local microcode compaction for horizontal machines. IEEE Transactions on Computers C-30, 7 (1981), 460-477.
128. C. V. Ramamoorthy and M. Tsuchiya. A high level language for horizontal microprogramming. IEEE Transactions on Computers C-23 (1974), 791-802.
129. S. Dasgupta and J. Tartar. The identification of maximal parallelism in straight-line microprograms. IEEE Transactions on Computers C-25, 10 (October 1976), 986-991.
130. P. W. Mallett. Methods of Compacting Microprograms. Ph.D. thesis. University of Southwestern Louisiana, Lafayette, 1978.
131. S. S. Yau, A. C. Schowe and M. Tsuchiya. On storage optimization of horizontal microprograms. Proc. Seventh Annual Workshop on Microprogramming (Palo Alto, California, 1974), 98-106.
132. R. L. Kleir. A representation for the analysis of microprogram operation. Proc. 7th Annual Workshop on Microprogramming (September 1974), 107-118.
133. M. Tokoro, E. Tamura and T. Takizuka. Optimization of microprograms. IEEE Transactions on Computers C-30, 7 (July 1981), 491-504.
134. E. S. Davidson. The design and control of pipelined function generators. Proc. 1971 International IEEE Conference on Systems, Networks, and Computers (Oaxtepec, Mexico, January 1971), 19-21.
135. M. Tokoro, T. Takizuka, E. Tamura and I. Yamaura. A technique of global optimization of microprograms. Proc. 11th Annual Workshop on Microprogramming (Asilomar, California, November 1978), 41-50.
136. G. Wood. Global optimization of microprograms through modular control constructs. Proc. 12th Annual Workshop on Microprogramming (Hershey, Pennsylvania, 1979), 1-6.
137. R. Grishman and S. Bogong. A preliminary evaluation of trace scheduling for global microcode compaction. IEEE Transactions on Computers C-32, 12 (December 1983), 1191-1194.
138. B. Su, S. Ding and L. Jin. An improvement of trace scheduling for global microcode compaction. Proc. 17th Annual Workshop on Microprogramming (New Orleans, Louisiana, October 1984), 78-85.
139. B. Su and S. Ding. Some experiments in global microcode compaction. Proc. 18th Annual Workshop on Microprogramming (Asilomar, California, November 1985), 175-180.

140. T. Gross and M. Ward. The suppression of compensation code, in Advances in Languages and Compilers for Parallel Computing, A. Nicolau, D. Gelernter, T. Gross and D. Padua (Editor). (Pitman/MIT Press, London, 1990), 260-273.
141. S. M. Freudenberger and J. C. Ruttenberg. Phase ordering of register allocation and instruction scheduling, in Code Generation -- Concepts, Tools, Techniques, Proceedings of the International Workshop on Code Generation, May 1991, R. Giegerich and S. L. Graham (Editor). (Springer-Verlag, London, 1992), 146-172.
142. P. G. Lowney, S. M. Freudenberger, T. J. Karzes, W. D. Lichtenstein, R. P. Nix, J. S. O'Donnell and J. C. Ruttenberg. The Multiflow trace scheduling compiler. The Journal of Supercomputing (this issue) (1992).
143. W. W. Hwu and P. P. Chang. Exploiting parallel microprocessor microarchitectures with a compiler code generator. Proc. 15th Annual International Symposium on Computer Architecture (Honolulu, Hawaii, May 1988), 45-53.
144. P. P. Chang, S. A. Mahlke, W. Y. Chen, N. J. Warter and W. W. Hwu. IMPACT: an architectural framework for multiple-instruction-issue processors. Proc. 18th Annual International Symposium on Computer Architecture (Toronto, Canada, May 1991), 266-275.
145. P. P. Chang and W. W. Hwu. Trace selection for compiling large C application programs to microcode. Proc. 21st Annual Workshop on Microprogramming and Microarchitectures (San Diego, California, November 1988), 21-29.
146. J. L. Linn. Horizontal microcode compaction, in Microprogramming and Firmware Engineering Methods, S. Habib (Editor). (Van Nostrand Reinhold, New York, 1988), 381-431.
147. A. Nicolau. Uniform parallelism exploitation in ordinary programs. Proc. International Conference on Parallel Processing (August 1985), 614-618.
148. A. Nicolau and R. Potasman. Realistic scheduling: compaction for pipelined architectures. Proc. 23th Annual Workshop on Microprogramming and Microarchitecture (Orlando, Florida, November 1990), 69-79.
149. K. Ebcioglu and A. Nicolau. A *global* resource-constrained parallelization technique. Proc. 3rd International Conference on Supercomputing (Crete, Greece, June 1989), 154-163.
150. S.-M. Moon and K. Ebcioglu. An efficient resource-constrained global scheduling technique for superscalar and VLIW processors. Proc. 25th Annual International Symposium on Microarchitecture (Portland, Oregon, December 1992).
151. J. A. Fisher. Trace Scheduling-2, an extension of trace scheduling. Technical Report. Hewlett-Packard Laboratories, 1992.
152. J. C. H. Park and M. S. Schlansker. On predicated execution. Technical Report HPL-91-58. Hewlett Packard Laboratories, 1991.
153. J. A. Fisher, D. Landskov and B. D. Shriver. Microcode compaction: looking backward and looking forward. Proc. 1981 National Computer Conference (1981), 95-102.
154. E. S. Davidson. Scheduling for pipelined processors. Proc. 7th Hawaii Conference on Systems Sciences (1974), 58-60.

155. A. T. Thomas and E. S. Davidson. Scheduling of multiconfigurible pipelines. Proc. 12th Annual Allerton Conference on Circuits and Systems Theory (Allerton, Illinois, 1974), 658-669.
156. E. S. Davidson, L. E. Shar, A. T. Thomas and J. H. Patel. Effective control for pipelined computers. Proc. COMPCON '90 (San Francisco, February 1975), 181-184.
157. J. H. Patel. Improving the Throughput of Pipelines with Delays and Buffers. Ph.D. thesis. University of Illinois, Urbana-Champaign, 1976.
158. J. H. Patel and E. S. Davidson. Improving the throughput of a pipeline by insertion of delays. Proc. 3rd Annual Symposium on Computer Architecture (January 1976), 159-164.
159. P. M. Kogge and H. S. Stone. A parallel algorithm for the efficient solution of a general class of recurrence equations. IEEE Transactions on Computers C-22, 8 (August 1973), 786-793.
160. P. M. Kogge. Maximal rate pipelined solutions to recurrence problems. Proc. First Annual Symposium on Computer Architecture (University of Florida, December 1973), 71-76.
161. P. M. Kogge. Parallel solution of recurrence problems. IBM Journal of Research and Development 18, 2 (March 1974), 138-148.
162. P. M. Kogge. Algorithm development for pipelined processors. Proc. 1977 International Conference on Parallel Processing (August 1977), 217.
163. P. M. Kogge. The microprogramming of pipelined processors. Proc. The 4th Annual Symposium on Computer Architecture (March 1977), 63-69.
164. D. Cohen. A methodology for programming a pipeline array processor. Proc. 11th Annual Microprogramming Workshop (Asilomar, California, November 1978), 82-89.
165. R. F. Touzeau. A FORTRAN compiler for the FPS-164 scientific computer. Proc. ACM SIGPLAN '84 Symposium on Compiler Construction (Montreal, Canada, 1984), 48-57.
166. S. Weiss and J. E. Smith. A study of scalar compilation techniques for pipelined supercomputers. Proc. Second International Conference on Architectural Support for Programming Languages and Operating Systems (Palo Alto, October 1987), 105-109.
167. P. Y. T. Hsu. Highly Concurrent Scalar Processing. Ph.D. thesis. University of Illinois, Urbana-Champaign, 1986.
168. M. S.-L. Lam. A Systolic Array Optimizing Compiler. Ph.D. thesis. Carnegie Mellon University, 1987.
169. P. Tirumalai, M. Lee and M. S. Schlansker. Parallelization of loops with exits on pipelined architectures. Proc. Supercomputing '90 (November 1990), 200-212.
170. M. Lam. Software pipelining: an effective scheduling technique for VLIW machines. Proc. ACM SIGPLAN '88 Conference on Programming Language Design and Implementation (June 1988), 318-327.
171. N. J. Warter, J. W. Bockhaus, G. E. Haab and K. Subramanian. Enhanced modulo scheduling for loops with conditional branches. Proc. The 25th Annual International Symposium on Microarchitecture (Portland, Oregon, December 1992).
172. B. Su, S. Ding and J. Xia. URPR - An extension of URCR for software pipelining. Proc. 19th Annual Workshop on Microprogramming (New York, New York, October 1986), 104-108.

173. B. Su, S. Ding, J. Wang and J. Xia. GURPR--a method for global software pipelining. Proc. 20th Annual Workshop on Microprogramming (Colorado Springs, Colorado, December 1987), 88-96.
174. B. Su and J. Wang. Loop-carried dependence and the General URPR software pipelining approach. Proc. 24th Annual Hawaii International Conference on System Sciences (Hawaii, January 1991).
175. B. Su and J. Wang. GURPR*: a new global software pipelining algorithm. Proc. 24th Annual International Symposium on Microarchitecture (Albuquerque, New Mexico, November 1991), 212-216.
176. A. Aiken and A. Nicolau. Perfect pipelining: a new loop parallelization technique, in Proceedings of the 1988 European Symposium on Programming. (Springer Verlag, New York, New York, 1988)
177. A. Aiken and A. Nicolau. Optimal loop parallelization. Proc. SIGPLAN'88 Conference on Programming Language Design and Implementation (Atlanta, Georgia, 1988), 308-317.
178. A. Aiken and A. Nicolau. A realistic resource-constrained software pipelining algorithm, in Advances in Languages and Compilers for Parallel Processing, A. Nicolau, D. Gelernter, T. Gross and D. Padua (Editor). (Pitman/The MIT Press, London, 1991), 274-290.
179. K. Ebcioglu and T. Nakatani. A new compilation technique for parallelizing loops with unpredictable branches on a VLIW architecture, in Languages and Compilers for Parallel Computing, D. Gelernter, A. Nicolau and D. Padua (Editor). (Pitman/The MIT Press, London, 1989), 213-229.
180. F. Gasperoni. Compilation techniques for VLIW architectures. Technical Report RC 14915. IBM Research Division, T. J. Watson Research Center, 1989.
181. T. Nakatani and K. Ebcioglu. Using a lookahead window in a compaction-based parallelizing compiler. Proc. 23th Annual Workshop on Microprogramming and Microarchitecture (Orlando, Florida, November 1990), 57-68.
182. S. Ramakrishnan. Software pipelining in PA-RISC compilers. Hewlett-Packard Journal (July 1992), 39-45.
183. M. Auslander and M. Hopkins. An overview of the PL.8 compiler. Proc. ACM SIGPLAN Symposium on Compiler Construction (Boston, Massachusetts, June 1982), 22-31.
184. T. R. Gross and J. L. Hennessy. Optimizing delayed branches. Proc. 15th Annual Workshop on Microprogramming (October 1982), 114-120.
185. J. L. Hennessy and T. Gross. Post-pass code optimization of pipelined constraints. ACM Transactions on Programming Languages and Systems 5, 3 (July 1983), 422-448.
186. W.-C. Hsu. Register allocation and code scheduling for load/store architectures. Computer Science Technical Report No. 722. University of Wisconsin, 1987.
187. R. L. Sites. Instruction ordering for the CRAY-1 Computer. Technical Report 78-CS-023. University of California, 1978.
188. J. Rymarczyk. Coding guidelines for pipelined processors. Proc. Symposium on Architectural Support for Programming Languages and Operating Systems (Palo Alto, California, March 1982), 12-19.
189. M. C. Golumbic and V. Rainish. Instruction scheduling beyond basic blocks. IBM Journal of Research and Development 344, 1 (January 1990), 93-97.

190. D. Bernstein, D. Cohen and H. Krawczyk. Code duplication: an assist for global instruction scheduling. Proc. 24th Annual International Symposium on Microarchitecture (Albuquerque, New Mexico, 1991), 103-113.
191. D. Bernstein and M. Rodeh. Global instruction scheduling for superscalar machines. Proc. SIGPLAN '91 Conference on Programming Language Design and Implementation (June 1991), 241-255.
192. S. Jain. Circular scheduling: a new technique to perform software pipelining. Proc. ACM SIGPLAN '91 Conference on Programming Language Design and Implementation (June 1991), 219-228.
193. M. Smotherman, S. Krishnamurthy, P. S. Aravind and D. Hunnicutt. Efficient DAG construction and heuristic calculation for instruction scheduling. Proc. 24th Annual International Workshop on Microarchitecture (Albuquerque, New Mexico, November 1991), 93-102.
194. R. Sethi and J. D. Ullman. The generation of optimal code for arithmetic expressions. Journal of the ACM 17, 4 (October 1970), 715-728.
195. R. Sethi. Complete register allocation problems. SIAM Journal of Computing 4, 3 (1975), 226-248.
196. A. V. Aho and S. C. Johnson. Optimal code generation for expression trees. Journal of the ACM 23, 3 (July 1976), 488-501.
197. J. L. Bruno and R. Sethi. Code generation for a one-register machine. Journal of the ACM 23, 3 (July 1976), 502-510.
198. A. V. Aho, S. C. Johnson and J. D. Ullman. Code generation for machines with multiregister operations. Proc. Fourth ACM Symposium on Principles of Programming Languages (1977), 21-28.
199. A. V. Aho, S. C. Johnson and J. D. Ullman. Code generation for expressions with common subexpressions. Journal of the ACM 24, 1 (January 1977), 146-160.
200. P. B. Gibbons and S. S. Muchnick. Efficient instruction scheduling for a pipelined architecture. Proc. ACM SIGPLAN '86 Symposium on Compiler Construction (July 1986), 11-16.
201. J. R. Goodman and W.-C. Hsu. Code scheduling and register allocation in large basic blocks. Proc. 1988 International Conference on Supercomputing (St. Malo, July 1988).
202. H. S. Warren. Instruction scheduling for the IBM RISC System/6000 processor. IBM Journal of Research and Development 34, 1 (January 1990), 85-92.
203. P. P. Chang, D. M. Lavery and W. W. Hwu. The importance of prepass code scheduling for superscalar and superpipelined processors. Technical Report No. CRHC-91-18. Center for Reliable and High-Performance Computing, University of Illinois, Urbana-Champaign, 1991.
204. G. J. Chaitin. Register allocation and spilling via graph coloring. Proc. ACM SIGPLAN Symposium on Compiler Construction (June 1982).
205. F. Chow and J. Hennessy. Register allocation by priority-based coloring. Proc. ACM SIGPLAN Symposium on Compiler Construction (1984), 222-232.
206. F. C. Chow and J. L. Hennessy. The priority-based coloring approach to register allocation. ACM Transactions on Programming Languages and Systems 12 (October 1990), 501-536.

207. L. J. Hendren, G. R. Gao, E. R. Altman and C. Mukerji. Register allocation using cyclic interval graphs: a new approach to an old problem ACAPS Technical Memo 33. Advanced Computer Architecture and Program Structures Group, McGill University, 1992.
208. D. Callahan and B. Koblenz. Register allocation via hierarchical graph coloring. Proc. SIGPLAN '91 Conference on Programming Language Design and Implementation (Toronto, Ontario, Canada, June 1991), 192-203.
209. B. R. Rau, M. Lee, P. Tirumalai and M. S. Schlansker. Register allocation for software pipelined loops. Proc. SIGPLAN'92 Conference on Programming Language Design and Implementation (San Francisco, June 17-19 1992).
210. W. Mangione-Smith, S. G. Abraham and E. S. Davidson. Register requirements of pipelined processors. Proc. International Conference on Supercomputing (Washington, D.C., July 1992).
211. P. P. Chang and W. W. Hwu. Profile-guided automatic inline expansion for C programs. Software Practice and Experience 22, 5 (May 1992), 349-376.
212. D. Callahan, S. Carr and K. Kennedy. Improving register allocation for subscripted variables. Proc. ACM SIGPLAN '90 Conference on Programming Language Design and Implementation (1990), 53-65.
213. B. R. Rau. Data flow and dependence analysis for instruction level parallelism, in Fourth International Workshop on Languages and Compilers for Parallel Computing, U. Banerjee, D. Gelernter, A. Nicolau and D. Padua (Editor). (Springer-Verlag, 1992), 236-250.
214. H. Zima and B. Chapman. Supercompilers for Parallel and Vector Computers. (Addison-Wesley, Reading, Massachusetts, 1990).
215. A. Nicolau. Parallelism, memory-anti-aliasing and correctness for trace scheduling compilers. Ph.D. thesis. Yale University, 1984.
216. C. C. Foster and E. M. Riseman. Percolation of code to enhance parallel dispatching and execution. IEEE Transactions on Computers C-21, 12 (December 1972), 1411-1415.
217. E. M. Riseman and C. C. Foster. The inhibition of potential parallelism by conditional jumps. IEEE Transactions on Computers C-21, 12 (December 1972), 1405-1411.
218. A. Nicolau and J. A. Fisher. Using an oracle to measure parallelism in single instruction stream programs. Proc. Fourteenth Annual Microprogramming Workshop (October 1981), 171-182.
219. A. Nicolau and J. A. Fisher. Measuring the parallelism available for very long instruction word architectures. IEEE Transactions on Computers C-33, 11 (November 1984), 968-976.
220. N. P. Jouppi and D. Wall. Available instruction level parallelism for superscalar and superpipelined machines. Proc. Third International Conference on Architectural Support for Programming Languages and Operating Systems (April 1989), 272-282.
221. D. W. Wall. Limits of instruction-level parallelism. Proc. Fourth International Conference on Architectural Support for Programming Languages and Operating Systems (April 1991), 176-188.
222. M. S. Lam and R. P. Wilson. Limits of control flow on parallelism. Proc. Nineteenth International Symposium on Computer Architecture (Gold Coast, Australia, May 1992), 46-57.

223. M. Butler, T. Yeh, Y. Patt, M. Alsup, H. Scales and M. Shebanow. Single instruction stream parallelism is greater than two. Proc. Eighteenth Annual International Symposium on Computer Architecture (Toronto, 1991).
224. M. A. Schuette and J. P. Shen. Instruction-level experimental evaluation of the Multiflow TRACE 14/300 VLIW computer. The Journal of Supercomputing (this issue) (1992).
225. M. Lee, P. P. Tirumalai and T.-F. Ngai. Software pipelining and superblock scheduling: compilation techniques for VLIW machines. Proc. 26th Annual Hawaii International Conference on System Sciences (Hawaii, 1993).