



**Hardware/Software Tradeoffs
for Increased Performance**

John Hennessy, Norman Jouppi,
Forest Baskett, Thomas Gross,
John Gill, and Steven Przybylski

Technical Report No. 22.8

February 1983

The MIPS project has been supported by the Defense Advanced Research Projects Agency under contract # MDA903-79-C-0680. Thomas Gross is supported by an IBM Graduate Fellowship.

Hardware/Software Tradeoffs for Increased Performance

John **Hennessy**, Norman Jouppi,
Forest **Baskett**, Thomas Gross,
John **Gill**, and Steven **Przybylski**

Technical Report No. 228

February 1983

Computer Systems Laboratory
Departments of Electrical Engineering and Computer Science
Stanford University
Stanford, California 94305

Abstract

Most new computer architectures are concerned with maximizing performance by providing suitable instruction sets for compiled code, and support for systems functions. We argue that the most effective design methodology must make simultaneous tradeoffs across all three **arcas**: hardware, software support, and systems support. Recent trends lean towards extensive hardware support for both the compiler and operating systems software. However, consideration of all possible design **tradeoffs** may often lead to less hardware support. Several examples of this approach are presented, including: omission of condition codes, **word**-addressed machines, and imposing pipeline interlocks in software. The specifics and performance of these approaches are examined with respect to the MIPS processor.

Key Words and Phrases: Computer architecture, instruction set design, VLSI, compiler design.

A version of this report appears in the Proceedings of the ACM Symposium on **Architcctural** Support for Programming Languages and Operating Systems, March 1982.

1 Introduction

Until recently, the design of new computer architectures and the demands of high level language compilation were not well integrated. Instead, architecture design was dominated by concerns of assembly language programs, and maintaining capability with features contained in old architectures. Recently, designers have been more conscious of architectures' major role as host machines for high level language programs. This development is evident in the modern microprocessors (e.g. the Z8000 and the 68000) and larger machines such as the VAX.

These architectures certainly address the concerns of running compiled programs more than past efforts have done. However, they assume a fixed and existent compiler technology. In most cases, compilers for other machines are used as a basis for examining the instruction set design (10,161. While this methodology is certainly better than a nonquantitative approach, it is inherently flawed. We propose a methodology whereby the compiler, and in fact even new compiler technology, are major inputs to the instruction set design process. The notion of the compiler and its code generator playing a vital role in the instruction set design process has been used in the design of a processor to run C [8]. In this paper, we will advocate more radical applications of compiler technology.

Adequate support for constructing systems software (e.g., the operating system), and for programs compiled from a high level language, is a necessity. Although instruction set design issues are largely ones of performance, cost and reliability, inadequate architectural support for the operating system often leads to real limitations in using the processor. Such limitations can determine which memory management techniques (swapping versus paging) are possible or feasible, and might restrict any system containing such a processor. On the other hand, including elaborate support in the processor architecture can often cause an overall decrease in performance because of the increased overhead. In many architectures, systems issues (e.g., orthogonal handling of interrupts and traps) are given low consideration. This trend prevails in many current microprocessor architectures; for example, in several architectures page faults are not accommodated. In other architectures, systems issues are the focus of the design to the detriment of performance, as in the Intel. iAPX-432.

We will argue that architecture/compiler/system **tradeoffs** are inevitable, and that by correctly making choices in all three areas the designer can create a simpler design which will have increased performance and reliability for lower cost.. We demonstrate several applications of this approach and give specific examples and empirical performance data from the MIPS project [7]. We are primarily concerned with procedural languages (e.g., ALGOL, Pascal, C. and Ada) and with von Neuman architectures.

2 Architectural support for compiled languages

Architectures can support compiled languages in many ways. For high performance both in executed code and within the compilers themselves, architectures can best support high level languages by providing simple, fast instruction sets that are easy for a compiler to use.

2.1 Simple and fast instructions

This emphasis on a simpler instruction set, and the attractiveness of this approach, have been argued and demonstrated by the RISC project [14]. From the compiler viewpoint, there are three advantages of a simple instruction **set**:

1. Because the instruction set is simpler, individual instructions are executed faster.
2. The compiler is not forced to attempt the utilization of a very sophisticated instruction which doesn't quite fit any particular high level construct. Besides slowing down the execution of other instructions, using overly sophisticated instructions is sometimes slower than using a customized sequence of simpler instructions [13].
3. Although these architectures may require more sophisticated compiler technology, the potential performance improvements to be obtained from faster machines and better compilers are substantial.

2.2 Load/store architectures

Compilers find load/store architectures to be a natural problem decomposition: first get the operands, then use them. Requiring compilers to do both at the same time when the architecture is not completely orthogonal is more complex. Load/store architectures can also increase the performance of compiled languages.

Load/store architectures can yield performance increases if frequently-used operands are kept in registers. Not only is redundant memory traffic decreased, but addressing calculations are saved as well. Software support for load/store **architectures** does not appear to be a problem either: there are efficient register allocation algorithms which produce good assignments [3, 1]. Although these algorithms require powerful and sophisticated compiler technology, they pay off by yielding very dense and near-optimal register assignments. Simpler register allocation algorithms are possible. However, these algorithms are not as effective because they do not consider all variables as equal candidates to be assigned to a register, and because they must conservatively avoid difficulties that can arise from potential aliasing.

Heavy use of **registers** will also improve code density. Code space can increase greatly when each **operation** has multiple memory addresses, even of the form displacement(base). Load and store instructions in MIPS

are at most 32 bits in length, and are of five types: long immediate, absolute, displacement(base), (base,index), and base shifted by n , $0 \leq n \leq 5$. These addressing modes require at most one ALU operation on data in general registers and immediate data from the instruction. The last three forms are less than 32 bits and may be packed with a possibly unrelated ALU or shifter operation. All MIPS instructions, although they may consist of up to two instruction pieces, are 32 bits in length and execute in one data memory cycle time. Because the two instruction pieces are disjoint they may be used for parts of two independent computations.

Orthogonal immediate fields can additionally increase the code density as they reduce the number of those loads which are executed to load a register with a constant. In the MIPS instruction format every operation can optionally contain a four-bit constant in the range 0-15 in place of a register field. Additionally, a move immediate instruction will load an S-bit constant into any register. Table 1 contains the distribution of constants (in magnitudes) found in a collection of Pascal programs including compilers and VLSI design aid software.

<u>Absolute Value</u>	<u>Percentage</u>
0	24.8%
1	19.0%
2	4.1%
3 - 15	20.8%
16 - 255	26.8%
> 255	4.5%

Table 1: Constant distribution in programs

The large majority of the constants in the range 16 - 255 represent character constants. Most of the large constants (>255) represent values that are directly related to the function of the program and do not relate to constants in other programs. Thus, a 4-bit constant should cover approximately 70% of the cases; the special 8-bit constant will catch all but 5%. To obtain small negative constants two approaches are possible: provide for a sign bit in the constants, or provide reverse operators that allow the constants to be treated as negative. MIPS uses the latter approach because it allows more constants to be expressed and eliminates the need for sign extension in the constant insertion hardware.

2.3 Condition codes and control flow

Many architectures have included condition codes as the primary method of implementing conditional control flow. Condition codes provide an implicit communication between two otherwise disjoint instructions. Condition codes can make life difficult for compiler writers, especially in nonorthogonal architectures. Even in reasonably orthogonal architectures (e.g., the M68000), working with condition codes is not simple. Difficulties occur largely because condition codes are side effects of instruction execution.

Condition codes are typically used for conditional control flow breaks, evaluation of boolean expressions, overflow detection and multiprecision arithmetic. Table 2 shows a typical set of features associated with condition codes and various architectures which possess these features. After discussing the disadvantages of condition codes from a hardware viewpoint, we will examine their use in each one of these instances and propose alternatives.

	Has condition code Set on moves	Set on operations	No Condition code
Conditional set Branch Access	VAX	M68000 360	MIPS POP- 10

Table 2: Condition code operations

Condition codes are difficult to implement primarily because they are an irregular structure. As opposed to registers, in which the references of updates of objects are explicit, condition codes are updated as a side-effect. In architectures in which some but not all instructions set the condition code, additional complexity is introduced in decoding and condition code control. Because of these irregularities, implementing condition codes is particularly painful on a heavily pipelined machine. With architectures which set the condition code heavily (such as the VAX), the hardware can assume that all instructions set the condition code; thus the hardware would force branch instructions to wait until the pipe has cleared. This may result in some performance loss when the instructions immediately before the branch do not affect the condition code. The implementation becomes substantially more difficult when the set of instructions that affect the condition code is large but not close to the set of all instructions (as in the 360/370). In this case, a significant amount of hardware is needed to determine whether condition codes are affected by a particular instruction. Because fewer instructions set the condition code, a designer dealing with a branch instruction is tempted to find the last instruction that affected the condition code. However, because the set of instructions that affect the condition code may be nontrivial to determine, finding the **last** instruction can be extremely difficult.

2.3.1 Conditional control flow breaks

Typically, condition codes are used to implement conditional control flow breaks. A compare instruction (or other arithmetic instruction) is used to set the condition code: a conditional branch instruction uses the condition code to **determine** whether to take a branch. The most obvious disadvantage of this approach is that two instructions are required to effect the usual compare and branch. This is not a major disadvantage since a single compare and branch instruction would take longer to execute and more instruction bytes to encode. From a hardware implementation viewpoint, it is also **useful** to know the branch condition explicitly.

There are two primary *hypothetical* advantages for condition codes:

1. They save instructions by allowing branches to use the results of computations that are already done.
2. Condition codes model three way (<, >, =) branches.

Table 3 contains empirical data which show that the number of instructions saved by condition codes is so small as to be essentially useless. Three way branches are a Fortran anomaly introduced into the language because of the architecture of the 704. In any event, we believe that the vast majority of the three way branches in Fortran contain only two branch destinations.

To implement conditional branches on a condition code machine, first the condition is set, and then a conditional branch is used. Most architectures set condition codes by ALU operations only; the VAX sets the condition code on all move operations.

MIPS and a few other architectures, such as the Cray-I and the PDP-10, implement conditional control flow using compare and branch instructions. In MIPS **all** instructions, including the compare and branch instructions, **take** the same amount of **execution** time. Thus, the comparison is to some extent **free**. In cases where explicit comparisons (the dominant case : see Table 3) are-needed on condition code machines, MIPS' approach actually **saves** instructions.

MIPS supports conditional control flow breaks using a compare and branch instruction with one of 16 possible comparisons. The 16 comparisons include both signed and unsigned arithmetic (e.g., $X < Y$) as well as logical (e.g., X and $Y = 0$) comparisons. Most conditional control flow breaks can be implemented with one comparison and the **displacement(base)** branch address available in this instruction. This, combined with the fact that a compare and branch instruction executes in one data memory cycle (as do all other instructions), makes this instruction a very fast means of implementing conditional control flow breaks.

Compares without condition codes	2469
Compares saved using condition codes set by operators only	25 = 1.1%
Compares saved using condi tion codes set by operators and moves	733
Moves used only to set condition code	7 0 6
Total compares saved by condition codes	27
Savi ngs for condition codes set by operators and moves	2.1%

Table 3: Use of condition codes

2.3.2 Evaluating boolean expressions

Handling boolean expressions with numeric comparisons can be difficult in many architectures. For example, consider:

`Found := (Rec = Key) OR (I = 13);`

where each variable except Found is an integer. On a condition code machine where the condition is accessible only by conditional branches (e.g., the VAX), typical code sequences are given in Figure 1. Clearly, early-out evaluation is much better. Early-out evaluation is frequently usable either because the language explicitly permits it or because the absence of side effects makes it possible. The high percentage of branches in this code is perhaps the most disturbing point, since the cost of branches on modern pipelined architectures is far more than the cost of a typical compute-type instruction.

<u>Full Evaluation</u>	<u>Early-out Evaluation</u>
<pre> str 0,r1 comp Rec,Key bne L str 1,r1 vL: comp I,13 bne D str 1,r1 D: str r1,Found </pre>	<pre> str 1,Found comp Rec,Key beq D comp I,13 beq D str 0,Found D: </pre>
<pre> 8 static instructions 2 branches Average of 7 instructions executed Always executes 2 branches </pre>	<pre> 6 static instructions 2 branches Average of 4.25 instructions executed Executes one branch on average </pre>

Figure 1: Evaluating boolean expressions with condition codes

An improvement over this can be gained by including instructions that conditionally set values based on the condition codes (as on the M68000). With such instructions, the improved code sequences in Figure 2 can be used. Although the average dynamic instruction count is slightly higher for this instruction sequence, it would execute faster on almost all machines since it has no branches.

```

comp Rec,Key
seq Found      ; R1 is set to bit that represents equal
comp I,13
seq r1
or  r1,Found

5 static/dynamic instructions
No branches

```

Figure 2: Boolean expression evaluation using conditional set

In an architecture without condition codes, some other method is needed to set the values. One approach is to duplicate the code used for a condition code machine without a conditional set instruction. Instead, MIPS provides a powerful *Set Conditionally* instruction with the same 16 comparisons found in conditional

branches. This instruction performs a comparison, and sets a register to zero or one based on the result. Using this instruction, the code sequence for our example is shown in Figure 3.

```
seq Rec,Key,r1
seq I,13,r2
of r1,r2,Found

3 static and dynamic instructions
No branches
```

Figure 3: Boolean expression evaluation using set conditionally

In addition to boolean expressions that are assigned to variables, boolean expressions appear in conditional tests. In this case, early-out evaluation will result in similar numbers of instruction counts. The conditional set approach will be superior only in the case of complex expressions (more than one boolean operators).

Average operators/boolean expression	1.66
Boolean expressions ending in jumps	80.9%
Boolean expressions ending in stores	19.1%

Table 4: Boolean expressions

Table 4 shows the distribution of boolean expression types for our Pascal data set Table 5 shows the number of operations needed per boolean operator to evaluate boolean expressions using different architectural support.

Compare/Register/Branch instructions per boolean operator	Static Dynamic	
	Set Conditionally instruction	2/1/0
CC and set conditionally based on CC	2/3/0	2/3/0
Only CC and branch, full evaluation	2/2/2	2/2/2
Only CC and branch, early-out	2/0/2	2/0/1.5

Table 5: Operations needed to evaluate a boolean expression

Without a conditional set operation, evaluation of a boolean expression to be stored will require an extra assignment. When evaluating a boolean expression for a conditional branch, the branch instruction will be part of the normal evaluation in the case of a condition branch-branch evaluation, but will be required in addition to the evaluation when conditional set evaluation is used. Using the data from previous tables, Table 6 shows the effectiveness for conditional set assuming that register operations take time 1, compares take time 2, and branches take time 4.

Type	Operations Evaluation	Full	Early-out
Store	Set conditionally/no CC	9.3	9.3
Store	CC/conditional set	14.9	14.9
Store	CC with only branch	27.9	20.5
Jump	Set conditionally/no CC	13.3	13.3
Jump	CC/conditional set	18.9	18.9
Jump	CC with only branch	26.9	19.5
Total	Set conditionally/no CC	12.5	12.5
Total	CC/conditional set	18.0	18.0
Total	CC with only branch	26.9	19.7
Improvement	Conditional set/CC	33.0%	8.6%
Improvement	Set conditionally	53.5%	36.5%

Table 6: Cost of evaluating boolean expressions

2.3.3 Overflow and multiple precision arithmetic

In a machine with condition codes, overflow bits may need to be explicitly tested, as is required on the M68000. This results in significant performance degradation if each result is tested via conditional traps, or a loss in reliability if no or few tests are made. Other machines trap automatically via hardware mechanisms. MIPS traps if overflow detection is enabled and stores the trap type *in a **surprise register***. This is discussed in greater detail in Section 3, Systems Support.

Carry bits are mainly used for multiprecision arithmetic. This is most important for 8- and 16-bit machines, and to a lesser degree for larger machines without floating point hardware. MIPS is in the second category. For intensive floating point applications, the use of a numeric coprocessor such as the Intel 8087 is envisioned. For more common occasional use, multiprecision arithmetic can be synthesized with 31-bit words. This does not cause problems unless precisions of $2w-1, 2w, 3w-2, \dots$ bits (where the wordsize is w) are required. Multiprecision numbers are usually smaller than an integral multiple of the wordsize, such as Fractions in doubleword floating point numbers.

3 Architectural Support for Systems

It is essential that the interfaces between the processor and its supporting external hardware be carefully considered at the architectural level if the system is to attain its full performance potential. In this section, these systems aspects of the MIPS processor are examined.

A goal of the MIPS processor was that it perform well in a wide range of environments, from a fully configured personal work station, to a much smaller system more typical of microprocessor applications. This goal implied the following requirements:

1. The support of virtual memory with demand paging
2. Efficient context switches
3. Effective exception handling
4. Privilege enforcement

These requirements must be balanced against the limited silicon resources, and the need for flexibility in the design of the rest of the system.

3.1 Memory mapping

There has been some controversy over whether a microprocessor should implement an on-chip memory management scheme at the expense of die size. Some processor designers have chosen to implement powerful on-chip memory management [15], while other have opted for off-chip structures [11]. In the MIPS architecture we attempt to achieve a good compromise by combining an optional page-level mapping unit off-chip with a simple yet elegant address space segmentation mechanism on-chip. The on-chip unit divides the virtual address space into a variable number of variably sized segments. The sum of the sizes of all segments cannot exceed the virtual address space of 16 million words. This restricts the address spaces of individual processes, and allows an off-chip page map to simultaneously contain entries for many processes without a corresponding increase in the tag field size. The on-chip segmentation is done by masking out the top n bits of every address and inserting an n -bit process identification number.

A process virtual address space thus can range from **65K** words to the full 16M words. It is split into two halves: one residing at the top of the program's virtual 32-bit address space, and the other at the bottom. Any attempt to reference a word between the two valid regions is treated as a page fault. The operating system then has the option of either re-allocating the process identifiers in order to increase the size of the process address space of the offending process; or terminating the offending process.

A rather novel aspect of the MIPS memory architecture is the existence of the free **memory** cycles provision. Since memory cycles are allocated to instructions, just as ALU or register access resources, an instruction that did not include a load or store piece would waste some of the memory bandwidth. Dynamic simulations indicated that the wasted bandwidth came close to 40% of the available bandwidth. To make use of the **otherwise** unused memory slots, a status pin on the processor indicates the **presence** of an upcoming

free memory cycle. Thus, these cycles can be used for DMA, I/O or cache write-backs. In addition, the architecture includes the proper mechanisms for normal block DMA and extended memory references.

3.2 Context switches

Context switches happen very frequently in typical multiprocessing environments. On every exception, interrupt or monitor call, control is taken from the executing process, and transferred (possibly only temporarily) to the kernel of the operating system. Many of these asynchronous events ultimately cause control to be passed to a third process. For conventional machines, the most time consuming aspect of context switching is the saving of registers and other miscellaneous state. Schemes to mitigate this overhead include multiple register sets and “move multiple register” instructions. Neither approach was a practical solution for MIPS because of the limited silicon resources and the constraint that all instructions execute in exactly five pipe stages. However, the dual instruction/data memory interface implies that a sequence of save register instructions could completely utilize the memory bandwidth for storing register contents, and thus would run as fast or faster than a microcoded **move-multiple** instruction.

The context switch time can be greatly increased if there is a large amount of additional state involved with the processor and memory map status. In MIPS, all the miscellaneous state of the processor is encapsulated into a single **surprise register** -- the MIPS equivalent of a processor status word. The surprise register includes the current and previous privilege levels, and enable bits for interrupts, overflow traps and memory mapping. Finally, there are two fields that specify the exact nature of the last exception. The addition of the on-chip **segmentation** means that most context switches do not require changes to the memory map.

To limit the size of the control portion of the processor, the privilege scheme in MIPS is restricted to a simple two level scheme. The only instructions that require supervisor privilege are those that read and write the surprise register and the on-chip segmentation registers. The current privilege level and mapping **state** are available to the rest of the system as part of the virtual address. To complete the protection scheme, the exterior mapping unit and any peripherals on the virtual address bus must be protected from user level processes.

3.3 Exceptional conditions: page faults, interrupts, and traps

The requirement that the processor support demand paging implies that it has the ability to restart instructions that have faulted. This is generally simpler in load/store architectures because instructions with memory references do not have secondary effects. However, in MIPS, the strict load/store framework is complicated by the concept of instruction pieces and the processor pipelining. An instruction can consist of a load or store piece and an ALU piece; the combined instruction can behave much like an auto increment or

decrement addressing mode. This potential difficulty is overcome by requiring an instruction that calls for a memory reference to not allow register writes to take place until after the reference has been committed. This restriction presents only a minor inconvenience in the instruction set design. If a data memory reference faults, the register write of an ALU operation during SX of the same instruction can be prevented. Page faults in the instruction stream do not encounter this difficulty. The absence of condition codes in the machine means that the actual operation has no permanent consequence, provided the store of the result to its destination register has been inhibited.

The proper handling of exceptions in a heavily pipelined machine is a major source of complexity and irregularity [9]. By an *exception* we mean all synchronous and asynchronous events that disrupt the **normal** flow of control. These include interrupts, software traps, both internal and external faults, and unrecoverable errors such as reset. The primary difficulty stems from the requirement that all instructions logically before a faulting instruction must be completed before the exception is handled. Otherwise, the instruction stream would not be easily **restartable**. Once an instruction **begins** a write to a register, it is impossible to stop it in a restartable state; thus an instruction once begun must be completed.

An additional level of complexity is introduced when another exception occurs during the completion of instructions that were partly executed when the original exception occurred. Assume an arithmetic overflow exists in an instruction that is sequentially before an instruction which incurs a mapping error. Due to the pipelining, the instruction fetch and decode of the second instruction occur before the ALU operation of the first. However, the overflow event must be handled before the mapping error because it originates in a first instruction.

Regardless of the cause or timing of an exception, the sequence of events it triggers remains the same. First, an attempt is made to complete any unfinished instructions. Second, the current status of the machine is saved, and subsequently changed to reflect execution by the operating system in physical address space. Last, the program counter is zeroed so that execution begins at the start of the first physical page. The standard dispatch routine that resides at address zero saves the return addresses, the surprise register, and a small number of the general purpose registers. Three return addresses are saved in order to allow returns to **sequences** that include indirect jumps, which have a branch delay of two. When an instruction following an indirect jump incurs an exception, the first three instructions to be executed in order to resume the code sequence are: the **offending** instruction, its successor, and then the target of the branch.

After saving this minimum state, the dispatch routine looks at the saved surprise register to determine what actually happened, and what routine should be invoked to handle the exception. This involves extracting from the top byte of the surprise register the two exception cause fields, and using the fields as an index into a

jump table Each exception handler can save more registers, enable interrupts, and resume memory mapping as it chooses. The initial dispatch routine performs a task that is typically done in microcode on larger machines. Since it must always be resident (even on the power-up reset exception) it must be put in a ROM on the virtual address bus.

In the case of an interrupt or software trap, a secondary dispatch is used to reach the actual handler. The trap code for software traps is 12 bits long, allowing 4096 different monitor calls. The external interrupt interface to the processor has been left extremely simple. There is a single interrupt line onto the chip; when the line is activated with interrupts enabled, a surprise sequence is initiated. After the first dispatch, the global interrupt handler queries any external prioritization logic to determine which device was requesting service. We decided on this mechanism due to pinout limitations and a desire to keep the control logic simple. Also, it was unclear what interrupt encoding scheme would be best suited to the variety of systems that we had in mind.

Among the main disadvantages of this general surprise scheme are the following:

1. It is slightly slower than a directly mapped, microcode *dispatch implementation. There are additional constraints on the surrounding system: the system must supply a small amount of ROM and RAM to which accesses will never page **fault**.
2. The return from interrupt sequence is complex and requires that the mapping and cache units be able to accept alternating references from two different address and privilege spaces.

The slower response speed is not usually significant, because the difference is small with respect to the total context switch time, and the MIPS approach allows a somewhat finer granularity of dispatching for many interrupts. The last two disadvantages entail complications which are unavoidable within the rigid framework of the fixed length instruction and a pipelined implementation.-

The overwhelming advantages of the scheme are that it can be implemented with a reasonable amount of hardware, both on- and off-chip, and is reasonably simple to use and implement. A secondary advantage is the orthogonality and flexibility of this design. Very few requirements are placed on a minimally configured system, yet a complete multiprocessing workstation would not find the architecture lacking in any significant way.

4 Compiler support for fast architectures

Compilers can help create speedier architectures by relaxing their requirements on the instruction set For example, consider an instruction to perform function F that is fairly complicated with respect to its demands on the hardware. The operation F may be a perfectly natural instruction in terms of a variety of source

languages and compilers. However, incorporating F as an instruction in the architecture may be a bad decision for several reasons:

1. When F is implemented, it may be slower than a custom tailored version of F , which is built from simpler instructions.
2. When F is added to the architecture, all the other instructions may suffer some performance degradation (albeit small in most cases). This phenomenon has been called the “n+ 1 instruction”* phenomenon and has been observed in several VLSI processor designs [12, 7].

Of course, whether or not F slows down the overall performance is dependent both on an implementation of the architecture and on the frequency of use of F . Instances of this type of behavior can be seen in the VAX 11/780. For example, the Index and CallS instructions in most instances are slower than obvious simpler code sequences [4,13].

A clear case can be made for eliminating more complex instructions when the speed benefits and usefulness of F are dubious. Advancing the case one step further, one can argue that even “natural” instructions should be examined for their usefulness and performance, when measured against possibly faster, customized sequences of simpler instructions.

4.1 Word-addressed machines

Most newer computer architectures, from micros up to large machines, have supported some sort of byte addressing. The primary reason for this is that character and logical data is often handled as bytes (at least when it is in an array) and byte addressing simplifies code generation and makes the machine faster.

While it is true that byte addressing makes code generation simpler, it is not at all clear that it has a positive influence on overall performance. Memory interfaces that must support byte addressability as well as word (and probably **halfword**) addressability are significantly more complicated. Our estimates are that a byte addressable memory interface would add **from** 15% to 20% additional overhead to the critical path of the MIPS VLSI processor. Because many processors assume that operand fetch times are constant (ignoring issues such as caches), all operand fetches will pay the cost of this overhead. Most of this overhead results **from** the necessity of doing byte **insert** or **extract operations**. (WC assume only a single memory access is needed and do consider the complexity of each extra read needed to implement byte stores). Other overhead comes from the added control complexity needed to implement byte addressing. If all instructions have byte and **halfword** formats to preserve orthogonality, the size of the control structure may be greatly increased. The **overhead** estimate given below ignores the cost associated with this extra hardware (i.e., the control hardware) and the resultant larger chip area, which will cause additional performance degradation.

An alternative is not to include byte addressability. The performance impact of such a choice depends on three factors:

1. The estimated cost of supporting byte addressing in terms of added operand fetching time for all operands.
2. The percentage of occurrences of byte-sized objects.
3. The cost of performing the byte operations on a word-addressed machine.

We will address the last two issues. Tables 7 and 8 contain data on storage references in terms of loads and stores. Table 7 allocates all objects as words unless they occur in a packed structure. Table 8 allocates all characters and booleans as bytes. Block movements of data are not included: neither as byte nor as word references. The source for the data is a collection of Pascal programs including compilers, optimizers, and VLSI design aid software; the programs are reasonably involved with text handling, and little or no compute intensive (e.g., floating point) tasks are included. The global activation records of the word-based allocation version average 20% larger.

All data references	-	71.2 % loads.	28.7 % stores
8 bit loads		2.6 %	
32 bit loads or larger		68.6 %	
8 bit stores		2.6 %	
32 bit stores or larger		26.2 %	
Character references	-	66.7 % loads.	33.3 % stores
8 bit character loads		14.7 %	
32 bit character loads		52.0 %	
8 bit character stores		21.5 %	
32 bit character stores		11.8 %	

Table 7: Data reference patterns in word-allocated programs

All data references	-	71.2 % loads.	28.7 % stores
8 bit loads.		6.6 %	
32 bit loads or larger		64.6 %	
8 bit stores		5.9 %	
32 bit stores or larger		22.9 %	

Table 8: Data reference patterns in byte-allocated programs

The major observations we can make from this data are the following:

- Objects allocated as full words dominate byte-allocated objects.
- Character reference patterns have a much higher percentage of stores than do non-character reference patterns.
- When unpacked character data is allocated in words, most of the data references are to the unpacked characters.

To evaluate the effectiveness of byte processing using a word-addressed machine, we need to examine **two** questions:

1. How can the compiler help us?
2. What instructions are available for extracting and isolating words from bytes?

The compiler can help by attempting to transform character at a time processing to word at a time processing. Since many of the operations that deal with characters concern copying and comparing strings, the potential benefits are substantial. This is an interesting code optimization problem, the benefits and difficulties of which are not obvious.

To explore the cost of character processing with only word addressing, we will look at the instructions in the MIPS instruction set and the typical code sequences. Given these and estimates of the overhead associated with byte addressing support, we can get a rough performance comparison.

Although MIPS does not have byte addressing, it does have special instructions for byte objects in packed arrays. Packed byte arrays are typically accessed via a byte index from the beginning of a word-aligned array; byte pointers can be regarded as the special case in which the array is located at memory location 0.

MIPS has four instructions to support byte objects: load and store base shifted; and insert and extract byte operations. Load and store base shifted can be used for accessing packed arrays of 2^n bit objects, where $0 \leq n \leq 15$ (i.e., bits through words). These instructions take a packed array pointer consisting of a $32-(5-n)$ bit word address in the high order bits and a packed array index in the low order n bits. This word is loaded/stored by shifting the packed array pointer n bits and reading/writing from that location as in every other load/store.

Bytes are accessed with insert/extract instructions. These insert or extract the byte specified by the low order two bits of a byte pointer. In the case of extract the byte pointer may be anywhere; for insert the byte pointer must be moved to a special register.

Loading a byte can be performed in two steps: first load the word containing the byte, then extract the byte from one of four locations within the word

If a byte pointer is in RO (the high order 30 bits contain a word address), then the following MIPS code sequence is equivalent to a load byte instruction:

```

;word at (r0/4) into r1
ld (r0>>2),r1
;extract byte from r1 into r1
xc ro,r1,r1

```

Thus, one memory access instruction and one ALU instruction are required to fetch a byte.

Storing into a byte array requires either two or three steps:

1. Fetch into a register the word that contains the destination byte (this step is often not needed because the word is in a register).
2. Replace the desired byte within the word register with the source byte.
3. Store the updated word into memory.

With the aid of an insert character instruction, a MIPS code sequence for store byte becomes:

```

;word at (r0/4) into r2
ld (r0>>2),r2
;r1 into byte selector lo
mov r1,lo
;low order byte of r3 into r2
ic lo,r3,r2
;return word in memory
st r2,(r0>>2)

```

This code sequence utilizes one or two memory reference instructions and two ALU instructions.

The cost of addressing operations using byte-addressed MIPS with/without overhead (15%) and using MIPS byte insert/extract instructions is shown in Table 9. We assume that the cost of an instruction is equal to the number of clock cycles needed to execute that instruction (or instruction piece). We also assume that non-array data can be accessed with the displacement field present in load and store instructions. Because word displacements are larger, word addressing has an advantage when displacements are too large for the byte addressed case.

Operation	Cost with byte operations	Cost with overhead	Cost with MIPS operations
load from array	4	4.6	6
store into array	4	4.6	8-12
load byte	6	6.9	8
store byte	6	6.9	10-18
load word	4	4.6	4
store word	4	4.6	4

Table 9: Cost of various byte operations

Table 10 analyzes the cost of word based addressing versus byte based addressing. The cost of addressing is computed from the cost of each type of addressing times its frequency. Because we ignore the extra addressing range of word offsets, use a minimum overhead factor, and ignore the extra read required to implement byte stores, these figures should be regarded as minimum improvements attributable to word based addressing. When all the factors are considered, improvements in the range of 20% for word-allocated programs and 23% for byte-allocated programs should be expected.

Operations	Word-allocated cost	Byte-allocated cost
byte loads on MIPS	.156	.476
byte stores on MIPS	.208-.312	.486-.75
word loads on MIPS	2.744	2.584
word stores on MIPS	1.048	.916
Total loads and stores on MIPS	4.156-4.26	4.162-4.426
byte loads on byte-addressed MIPS	.12	.396
byte stores on byte-addressed MIPS	.12	.347
word loads on byte-addressed MIPS	3.202	2.972
word stores on byte-addressed MIPS	1.205	1.053
Total loads and stores on byte-addressed MIPS	4.647	4.768
Byte addressing performance penalty	9% - 11.8%	7.7 - 14.6%

Table 10: Cost of byte and word addressed based architectures

4.2 Applying better compiler technology

Another approach to faster, cheaper architectures is to require more software support in the compiler. This is an attempt to shift the burden of the cost from hardware to software. The shifting of the complexity from hardware to software has several major advantages:

- The complexity is paid for only once during compilation. When a user runs his program on a complex architecture, he pays the cost of the architectural overhead each time he runs his program.
- It allows the concentration of energies on the software, rather than on constructing a complex hardware engine, which is hard to design, debug, and **efficiently** utilize. Software is not necessarily easier to construct, but VLSI-based implementations make hardware simplicity important.

Naturally, the purpose of placing more emphasis on the software aspects is to improve the overall performance. WC also hope to improve the cost effectiveness by using compiler technology that is more powerful but not much more complex.

4.2.1 Software-imposed pipeline interlocks

The interlock hardware in a pipeline processor normally provides these **functions**:

- If an operand is fetched from memory, the **interlock** mechanism delays subsequent instructions that **attempt to reference** the data until the operand is available.

- The pipeline is cleared after the execution of a flow-control instruction. If a branch is taken, the interlock mechanism guarantees that the next instruction executed is the instruction at the destination of the branch.
- Arithmetic operations may require different amounts of time to execute, e.g., multiply and divide. The interlock hardware will prevent the next instruction from executing if there are source-destination or destination-source dependencies.

We will discuss a software-based implementation of the first two functions, and present an algorithm for the first function and some empirical data on its performance.

In a processor with interlock hardware that provides the functions listed above, pipelining can be seen as an optimization implemented by hardware. Basically, the hardware will execute the program faster, subject to the interlocks which prevent illegal optimizations. This approach allows the compiler (or other user of the machine-level instruction set) to make simple assumptions about the execution of individual machine instructions.

An alternative approach is to move these optimizations from hardware to software. In that case there is no hardware interlock mechanism. Instead, the functions described above have to be provided by software, either by rearranging the code sequence or by inserting no-ops. This approach has the potential of producing code which executes faster, at the expense of the additional effort required to reorganize the instruction stream. No-ops will only explicitly delay the execution as compared to the invisible delays imposed by the hardware in an architecture with interlocks.

Such a reorganization scheme makes use of knowledge about the interdependencies of the individual instructions. The techniques developed for code optimization can be adapted to handle the requirements of the reorganization algorithms. The MIPS architecture employs the approach outlined here: there are no hardware interlocks. The current scheme provides the reorganization as a post-processing of the code generator's output. This reorganizer performs several major functions:

1. It takes the pipeline constraints into account and reorganizes the code to avoid interlocks when possible, and otherwise inserts no-ops.
2. It packs instruction pieces into one 32-bit word.
3. It assembles instructions.

Thus, the reorganizer also works on programmer-written assembly language code and reorganizes, packs, and assembles it.

Since the code reorganization process is part of every compilation, we must concentrate on solutions which

have acceptable run-time performance and still produce good results in most cases. Finding an optimal code sequence is very expensive, as the problem is NP-hard [6]. All code reorganization is done on a basic block basis. The algorithm is discussed in detail in [6].

The basic steps in the algorithm are the following:

1. Read in a basic block and create a machine-level dag that represents the dependencies between individual instruction pieces.
2. Given the set of instructions generated so far, determine sets of instructions that can be generated next.
3. Eliminate any sets that cannot be started immediately.
4. If there are no sets left, emit a no-op and return to step 2. Otherwise, choose from among the sets remaining.

The choice of the next instruction to be scheduled (step 4 above) is made heuristically from the set of legal instructions. Typically, an instruction that fits in a hole in a **nonfull** instruction is preferred; this provides the instruction packing.

When determining which sets of instructions may be processed, the reorganizer must examine both the interlocks from earlier instructions and register usage in parallel **dags**. The use of registers in parallel **dags** partially determines what set of code reorderings are possible. The algorithm must also avoid reordering loads and stores that might be aliased.

All branches in MIPS are delayed branches with a single instruction delay. If instruction ***i*** is a branch to ***L*** and the branch is taken, then the sequence of instructions executed is ***i, i + 1, L***. There are three major schemes for dealing with delayed branches of delay ***n***:

1. Move ***n*** instructions from before the branch till **after** the branch.
2. If the branch is a **backwards loop** branch, then duplicate the first ***n*** instructions in the loop and branch to the ***n + 1*** instruction.
3. If the branch is conditional, move the next ***n*** sequential instructions so they immediately follow the branch.

Of course, if the branch is conditional the outcome of the test must not depend on any of the moved instructions. Often the front end of the compiler is able to handle delayed branches better than the reorganizer; in this case it emits a pseudo-op which tells the reorganizer that this sequence is not to be touched. The branch delay optimization algorithm and its performance are discussed in [5].

Figure 4 shows how a code fragment is affected by reorganization. In this case, it is assumed that r2 is “dead” outside of the section shown, therefore it can be modified even if the branch in line 2 is taken. Note also that the store instruction is not moved, as it affects memory.

Legal Code with No-ops	Reorganized Code
ld 2(ap), r0 ble r0, #1, L11 No-op No-op sub #1, r0, r2 st r2, 2(sp)	ld 2(ap), r0 ble r0, #1, L11 No-op sub #1, r0, r2 st r2, 2(sp)
. ld 3(sp), r5 add r5, r0 add #1, r4 bra L3	bra L3 ld 3(sp), r5 add #1, r4
L3: . . .	L3: . . .

Figure 4: Reorganization, packing, and branch delay

To show the effectiveness of these optimizations, we ran versions of a program that does reorganization, packing, and branch delay elimination of three input programs. The input programs consist of an implementation of computing Fibonacci numbers and two implementations of the Puzzle benchmark [2]. All the programs were written in C and compiled to instruction pieces by a version-of the Portable C Compiler. The data in Table 11 show the improvements in static instruction counts.

<u>Optimization</u>	<u>Fibonacci</u>	<u>Puzzle 0</u>	<u>Puzzle 1</u>
None (no-ops inserted)	63	843	1219
Reorganization	63	834	1113
Packing	55	776	992
Branch delay	50	634	791
Total Improvement	20.6%	24.8%	35.1%

Table 11: Cumulative improvements with postpass optimization

5 Conclusions

We contend that the most effective performance can be obtained by a design methodology that makes tradeoffs across the boundaries between hardware and software. This approach is just the opposite of some architectures that advocate extensive hardware support. Several experimental projects, including MIPS and RISC, are pursuing the goal of a simplified hardware implementation.

We examine the issue of instruction set design to support the execution of compiled code. As opposed to “language-oriented” developments such as stack machines, we advocate two major alterations: the use of load/store architectures, and the absence of condition codes. Both of these design alternatives have two major advantages: simpler hardware implementation, and potentially higher individual instruction efficiency.

Second, we discuss the issue of supporting the construction of systems by providing the necessary primitive functions in hardware. We note that in many cases the minimum functionality is not available. This limits the type of design approach that can be used. We consider the complete support of page faults and interrupts in detail, and discuss the disadvantages of approaches that provide extensive system support services.

The issue of word-based versus byte-based addressing is explored. Based on a set of empirical data, we conclude that for many applications architectures will have higher performance with word addressing. Word based addressing gains its advantages from two primary points: it has a lower overhead associated with each fetch or store, and word references occur much more frequently than byte references. To make a word based approach feasible, special support for accessing bytes (as in the MIPS instruction set) are needed.

Lastly, we discuss approaches that rely extensively on improved compiler technology. The compiler is taken into account both in terms of the instruction set and in terms of radically simplified hardware designs. The software imposition of interlocks is presented as an example of this approach.

We explore the concept of simultaneously dealing with the hardware implementation, systems requirements, and the compiler technology. Following such an approach may lead to a reduction in hardware, as opposed to additional hardware support. Whether such an approach is effective and efficient depends on a wide variety of factors, including resultant improvements in hardware **performance**, the usefulness of the feature, and the alternative cost without “hardware **support**.” Using several specific examples, we have shown that such tradeoffs can produce significant improvements in overall performance.

References

1. Aho, A.V. and Ullman J.D. ***Principles of Cotnpiler Design***. Addison-Wesley, Menlo Park, 1977.
2. Baskett, F. Puzzle: an informal compute bound benchmark. Widely circulated and run.
3. Chaitin, G.J., Auslander, M.A., Chandra, A.K., Cocke, J., Hopkins, M.E., Markstein, P.W. Register Allocation by Coloring. Research Report 8395, IBM Watson Rescarch Center, 1981.
4. ***DECVAX11 Architecture Handbook*** Digital Equipment Corp., Maynard, MA., 1979.
5. Gross, T.R. and Hennessy, J.L. Optimizing Delayed Branches. Proceedings of Micro-15, IEEE, October, 1982, pp. 114-120.
6. Hennessy, J.L. and Gross, T.R. Code Generation and Reorganization in the Presence of Pipeline Constraints. Proc. Ninth POPL Conference, ACM, January, 1982, pp. 120-127.
7. Hennessy, J.L., Jouppi, N., Baskett, F., and Gill, J. MIPS: A VLSI Processor Architecture. Proc. CMU Conference on VLSI Systems and Computations, October, 1981, pp. 337-346.
8. Johnson, S.C. A 32-Bit Processor Design. Tech. Rept Computing Science #80, Bell Labortories, Murray Hill, April, 1979.
9. Lampson, B.W., McDaniel, G.A. and SM. Omstein. An Instruction Fetch Unit for a High Performance Personal Computer. Tech. Rept CSL-81-1, Xerox PARC, January, 1981.
10. Lunde, A. "Empirical Evaluation of Some Features of Instruction Set Processor Architectures." *CACM* 20, 3 (March 1977), 143-152.
11. ***MC68000 Users Manual***. 2nd edition, Motorola Inc., 1980.
12. Murphy, B.T. and Molinelli, J.J. A 32-Bit Single Chip CMOS Microprocessor. Seminar given at the Integrated Circuits Laboratory, Stanford University, May 22, 1981.
13. Patterson, D.A. and Ditzel, D.R. "The Case for the Reduced Instruction Set Computer." ***Computer Architecture News*** 8, 6 (October 1980), 25 - 33.
14. Patterson, D.A. and Sequin C.H. RISC-I: A Reduced Instruction Set VLSI Computer. Proc. of the Eighth Annual Symposium on Computer Architecture, Minneapolis, Minn., May, 1981, pp. 443 - 457.
15. Pollack, F., Cox, G., Hammerstrom, D., Kahn, K., Lai, K., Rattner, J. Supporting Ada Memory Management in the iAPX-432. Proc. Sym. on Architectural Support for Programming Languages and Operating Systems, ACM, Palo Alto, Ca., March, 1982, pp. 117-131..
16. Shustek, L.J. ***Analysis and Performance of Computer Instruction Sets***. Ph.D. Th., Stanford University, May 1977. Also published as SLAC Report 205.