# The Raw Prototype Design Document

V5.01

Michael Taylor (mbt@mit.edu)

Department of Electrical Engineering and Computer Science
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

Sept 6, 2004

# TABLE OF CONTENTS

# **1** INTRODUCTION

### 1.0  MANIFEST

In the introduction of this design document, I start by motivating the Raw architecture discipline, from a computer architect's viewpoint.

I then discuss the goals of the Raw prototype processor, a research implementation of the Raw philosophy. I elaborate on the research questions that the Raw group is trying to answer.

In the body of the design document, I will discuss some of the important design decisions in the development of the Raw prototype, and their effects on the overall development.

Finally, at the end of the document comes the complete user's manual, as well as a number of individual design chapters for various components of the Raw prototype.

### 1.1 MOTIVATION FOR A NEW TYPE OF PROCESSOR

### 1.1.1  The sign of the times

The first microprocessor builders designed in a period of famine. Silicon area on die was so small in the early seventies that the great challenge was just in achieving important features like reasonable data and address widths, virtual memory, and support for external I/O.

A decade later, advances in material science provided designers with enough resources that silicon was neither precious nor disposable. It was a period of moderation. Architects looked to advanced, more space consuming techniques like pipelining, out-of-order issue, and caching to provide performance competitive with minicomputers. Most of these techniques were borrowed from supercomputers, and were carefully added from generation to generation as more resources became available.

The next decade brings with it a regime of excess. We will have billions of transistors at our disposal. The new challenge of modern microprocessor architects is very simple: we need to provide the user with an effec-

tive interface to the underlying raw computational resources.

### 1.1.2  An old problem: SpecInt

In this new era, we could continue on as if we still lived in the moderation phase of microprocessor development. We would incrementally add micro-architectural mechanisms to our superscalar and VLIW processors, one by one, carefully measuring the benefits.

For today's programs, epitomized by the SpecInt95 benchmark suite, this is almost certain to provide us with the best performance. Unfortunately, this approach suffers from exponentially growing complexity (measured by development and testing costs and man-years) that is not being sufficiently mitigated by our sophisticated design tools, or by the incredible expertise that we have developed in building these sorts of processors. Unfortunately, this area of research is at a point where increasing effort and increasing area is yielding diminishing returns [Hennessey99].

Instead, we can attack a more fuzzy, less defined goal. We can use the extra resources to expand the scope of problems that microprocessors are skilled at solving. In effect, we redirect our attention from making processors better at solving problems they are already, frankly, quite good at, towards making them better at application domains which they currently are not so good at.

In the meantime, we can continue to rely on the as-yet juggernaut march of the fabrication industry to give us a steady clock speed improvement that will allow our existing SpecInt applications to run faster than ever.

### 1.1.3  A new problem: Extroverted computing

Computers started out as very oblivious, introverted devices. They sat in air-conditioned rooms, isolated from their users and the environment. Although they communicated with EACH OTHER at high speeds, the bandwidth of their interactions with the real world was amazingly low. The primary input devices, keyboards, provided at most tens of characters per second. The output bandwidth was similarly pathetic.

With the advent of video display and sound synthesis, the output bandwidth to the real world has blossomed to 10s of megabytes per second. Soon, with the

advent of audio and video processing, the input bandwidth will match similar levels.

As a result of this, computers are going to become more and more aware of their environments. Given sufficient processing and I/O resources, they will not only become passive recorders and childlike observers of the environment, they will be active participants. In short, computers will turn from recluse introverts to extroverts.

The dawn of the extroverted computing age is upon us. Microprocessors are just getting to the point where they can handle real-time data streams coming in from and out to the real world. Software radios and cell phones can be programmed in a 1000 lines of C++ [Tennenhouse95]. Video games generate real-time video, currently with the help of hardware graphics back ends. Real-time video and speech understanding, searching, generation, encryption, and compression are on the horizon. What once was done with computers for text and integers will soon be done for analog signals. We will want to compose sound and video, search it, interpret it, and translate it.

Imagine, while in Moscow, you could talk to your wrist watch and tell it to listen to all radio stations for the latest news. It would simultaneously tune into the entire radio spectrum (whatever it happens to be in Russia), translate the speech into English, and index and compress any news on the U.S. At the same time, your contact lens display would overlay English translations of any Russian word visible in your sight, compressing and saving it so that you can later edit a video sequence for your kids to see (maybe you'll encrypt the cab ride through the red light district with DES-2048). All of these operations will require massive bandwidth and processing.

### 1.1.4  New problem, old processors?

We could run our new class of extroverted applications on our conventional processors. Unfortunately, these processors are, well, introverted.

First off, conventional processors often treat I/O processing as a second class citizen to memory processing. The I/O requests travel through a hierarchy of slower and slower memory paths, and end up being bottlenecked at the least common denominator. Most of the pins are dedicated to caches, which ironically, are intended to minimize communication with the outside world. These caches, which perform so well on conven-

tional computations, perform poorly on streaming, extroverted, applications which have infinite data streams that are briefly processed and discarded.

Secondly, these new extroverted applications often have very plentiful fine grained parallelism. The conventional ILP architectures have complicated, non-scalable structures (multi-ported or rotating register files, speculation buffers, deferred exception mechanisms, pools of ALUs) that are designed to wrest small degrees of parallelism out of the most twisty code. The parallelism in these new applications does not require such sophistication. It can be exploited on architectures that are easy to design and are scalable to thousands of active functional units.

Finally, the energy efficiency of architectures needs to be considered to evaluate their suitability for these new application domains. The less power microprocessors need, the more and more environments they can exist in. Power requirements create a qualitative difference along the spectrum of processors. Think of the enormous difference among 1) machines that require large air conditioners, 2) ones that need to be plugged in, 3) ones that run on batteries, and ultimately, 4) ones that runs off their tiny green chlorophyllic plastic case.

### 1.1.5  New problems, new processors.

It is not unlikely that existing processors can be modified to have improved performance on these new applications. In fact, the industry has already made some small baby steps with the advent of the Altivec and MAX-2 technologies [Lee96].

The Raw project is creating an extroverted architecture from scratch. We take as our target these data-intensive extroverted applications. Our architecture is extremely simple. Its goal is to expose as much of the copious silicon and pin resources to these applications. The Raw architecture provides a raw, scalable, parallel interface which allows the application to make direct use of every square millimeter of silicon and every I/O pin. The I/O mechanism allows data to be streamed directly in and out of the chip at extraordinary rates.

The Raw architecture discipline also has advantages for energy efficiency. However, they will not be discussed in this design document.

## 1.2 THIS DOCUMENT AND HOW IT RELATES TO RAW

This design document details the decisions and ideas that have shaped the development of a prototype of the new type of processor that our group has developed. This process has been the result of the efforts of many talented people. When I started at MIT three years ago, the Raw project was just beginning. As a result, I have the luxury of having a perspective on the progression of ideas through the group. Initially, I participated in much of the data gathering that refined our initial ideas. As time passed on, I became more and more involved in the development of the architecture. I managed the two simulators, hand-coded a number of applications, worked on some compiler parallelization algorithms, and eventually joined the hardware project. I cannot claim to have originated all of the ideas in this design document; however I can reasonably say that my interpretation of the sequence of events and decisions which lead us to this design point probably is uniquely mine. Also uniquely mine probably is my particular view of what Raw should look like.

Anant Agarwal and Saman Amarasinghe are my fearless leaders. Not enough credit goes out to Jonathan Babb, and Matthew Frank, whose brainstorming planted the first seeds of the Raw project, and who have continued to be a valuable resource. Jason Kim is my partner in crime in heading up the Raw hardware effort. Jason Miller researched I/O interfacing issues, and is designing the Raw handheld board. Mark Stephenson, Andras Moritz, and Ben Greenwald are developing the hardware/software memory system. Ben, our operating systems and tools guru also ported the GNU binutils to Raw. Albert Ma, Mark Stephenson, and Michael Zhang crafted the floating point unit. Sam Larsen wrote the static switch verilog. Rajeev Barua and Walter Lee created our sophisticated compiler technology. Elliot Waingold wrote the original simulator. John Redford and Chris Kappler lent their extensive industry experience to the hardware effort.

### 1.2.1  Design document thesis statement

**The Raw Prototype Design is an effective design for a research implementation of a Raw architecture workstation.**

### 1.2.2  The goals of the prototype

In the implementation of a research prototype, it is important early on to be excruciatingly clear about one's goals. Over the course of the design, many implementa-

tion decisions will be made which will call into question these goals. Unfortunately, the "right" solution from a purely technical standpoint may not be the correct one for the research project. For example, the Raw prototype has a 32-bit architecture. In the commercial world, such a paltry address space is a guaranteed trainwreck in the era of gigabit DRAMs. However, in a research prototype, having a smaller word size gives us nearly twice as much area to further our research goals. The tough part is making sure that the implementation decisions do not invalidate the research's relevance to the real world.

**Ultimately, the prototype must serve to facilitate the exploration and validation of the underlying research hypotheses.**

The Raw project, underneath it all, is trying to answer two key research questions:

### 1.2.3  The Billion Transistor Question

**What should the billion transistor processor of the year 2007 look like?**

The Raw design philosophy argues for an array of replicated tiles, connected by a low latency, high throughput, pipelined network.

This design has three key implementation benefits, relative to existing superscalar and VLIW processors:

First, the wires are short. Wire length has become a growing concern in the VLSI community, now that it takes several cycles for a signal to cross the chip. This is not only because the transistors are shrinking, and die sizes are getting bigger, but because the wires are not scaling with the successive die shrinks, due to capacitive and resistive effects. The luxurious abstraction that the delay through a combinational circuit is merely the sum of its functional components no longer holds. As a result, the chip designer must now worry about both congestion AND timing when placing and routing a circuit. Raw's short wires make for an easy design.

Second, Raw is physically scalable. This means that all of the underlying hardware structures are scalable. All components in the chip are of constant size, and do not grow as the architecture is adapted to utilize larger and larger transistor budgets. Future generations of a Raw architecture merely use more tiles without negatively impacting the cycle time. Although Raw offers scalable computing resources, this does not mean that we will necessarily have scalable performance. That is dependent on the particular application.

Finally, Raw has low design and verification complexity. Processor teams have become exponentially larger over time. Raw offers constant complexity, which does not grow with transistor budget. Unlike today's superscalars and VLIWs, Raw does not require a redesign in order to accommodate configurations with more or fewer processing resources. A Raw designer need only design the smaller region of a single tile, and replicate it across the entire die. The benefit is that the designer can concentrate all of one's resources on tweaking and testing a single tile, resulting in clock speeds higher than that of monolithic processors.

### 1.2.4 The "all-software hardware" question

**What are the trade-offs of replacing conventional hardware structures with compilation and software technology?**

Motivated by advances in circuit compilation technology, the Raw group has been actively exploring the idea of replacing hardware sophistication with compiler smarts. However, it is not enough merely to reproduce the functionality of the hardware. If that were the case, we would just prove that our computing fabric was Turing-general, and move on to the next research project. Instead our goal is more complex. For each alternative solution that we examine, we need to compare its area-efficiency, performance, and complexity to that of the equivalent hardware structure. Worse yet, these numbers need to be tempered by the application set which we are targeting.

In some cases, like in leveraging parallelism, removing the hardware structures allows us to better manage the underlying resources, and results in a performance win. In other cases, as with a floating point unit, the underlying hardware accelerates a basic function which would take many cycles in software. If the target application domain makes heavy use of floating point, it may not be possible to attain similar performance per unit area regardless of the degree of compiler smarts. On the other hand, if the application domain does not use floating point frequently, then the software approach allows the application to apply that silicon area to some other purpose.

### 1.3 SUMMARY

In this section, I have motivated the design of a new family of architectures, the Raw architectures. These architectures will provide an effective interface for the amazing transistor and pin budgets that will come in the next decade. The Raw architectures anticipate the arrival of a new era of extroverted computers. These extroverted computers will spend most of their time interacting with the local environment, and thus are optimized for processing and generating infinite, real-time data streams.

I continued by stating the thesis statement, that the Raw prototype design is an effective design for a research implementation of a Raw architecture workstation. I finished by explaining the central research questions of the Raw project.

# 2 EARLY DESIGN DECISIONS

## 2.0  THE BIRTH OF THE FIRST RAW ARCHITECTURE

### 2.0.1  RawLogic, the first Raw prototype

Raw evolved from FPGA architectures. When I arrived at MIT almost three years ago, Raw was very much in its infancy. Our original idea of the architecture was as a large box of reconfigurable gates, modeled after our million-gate reconfigurable emulation system. Our first major paper, the Raw benchmark suite, showed very positive results on the promise of configurable logic and hardware synthesis compilation. We achieved speedups on a number of benchmarks; numbers that were crazy and exciting [Babb97].

However, the results of the paper actually considerably matured our viewpoint. The term "reconfigurable logic" is really very misleading. It gives one the impression that silicon atoms are actually moving around inside the chip to create your logic structures. But the reality is, an FPGA is an interpreter in much the same way that a processor is. It has underlying programmable hardware, and it runs a software program that is interpreted by the hardware. However, it executes a very small number of very wide instructions. It might even be viewed as an architecture with a instruction set optimized for a particular application; the emulation of digital circuits. Realizing this, it is not surprising that our experiences with programming FPGA devices show that they are neither superior nor inferior to a processor. It is merely a question of which programs run better on which interpreter.

In retrospect, this conclusion is not all that surprising; we already know that FPGAs are better at logic emulation than processors; otherwise they would not exist. Conversely, it is not likely that the extra bit-level flexibility of the FPGA comes for free. And, in fact, it does not. 32-bit datapath operations like additions and multiplies perform much more quickly when optimized by an Intel circuit hacker on a full-custom VLSI process than when they are implemented on a FPGA substrate. And again, it is not much wonder, for the processor's multiplier has been realized directly in silicon, while the multiplier implementation on the FPGA is running under one level of interpretation.

### 2.0.2  Our Conclusions, based on Raw logic

In the end, we identified three major strengths of FPGA logic, relative to a microprocessor:

**FPGAs make a simple, physically scalable parallel fabric.**

For applications which have a lot of parallelism, we can easily exploit it by adding more and more fabric.

**FPGAs allow for extremely fast communication and synchronization between parallel entities.**

In the realm of shared memory multiprocessors, it takes tens to hundreds of cycles for parallel entities to communication and synchronize [Agarwal95]. When a silicon compiler compiles parallel verilog source to an FPGA substrate, the different modules can communicate on a cycle-by-cycle basis. The catch is that this communication often must be statically scheduled.

**FPGAs are very effective at bit and byte-wide data manipulation.**

Since FPGA logic functions operate on small bit quantities, and are designed for circuit emulation, they are very powerful bit-level processors.

We also identified three major strengths of processors relative to FPGAs:

**Processors are highly optimized for datapath oriented computations.**

Processors have been heavily pipelined and have custom circuits for datapath operations. This customization means that they process word-sized data much faster than an FPGAs.

**Compilation times are measure in seconds, not hours [Babb97].**

The current hardware compilation tools are very computationally intensive. In part, this is because the hardware compilation field has very different requirements from the software compilation field. A smaller, faster circuit is usually much more important than fast compilation. Additionally, the problem sizes of the FPGA compilers are much bigger -- a net list of NAND gates is much larger than a dataflow graph of a typical program. This is exacerbated by the fact that the synthesis tools decompose identical macro-operations like 32-bits adds into separately optimized netlists of bit-wise operations.

**Processors are very effective for just getting through the millions of lines of code that AREN'T the inner loop.**

The so-called 90-10 rule says that 90 percent of the time is spent in 10 percent of the program code. Processor caches are very effective at shuffling infrequently used data and code in and out of the processor when it is not needed. As a result, the non-critical program portions can be stored out to a cheaper portion of the memory hierarchy, and can be pulled in at a very rapid rate when needed. FPGAs, on the other hand, have a very small number (one to four) of extremely large, descriptive instructions stored in their instruction memories. These instructions describe operations on the bit level, so a 32-bit add on an FPGA takes many more instruction bits than the equivalent 32-bit processor instruction. It often takes an FPGA thousands or millions of cycles to load a new instruction in. A processor, on the other hand, can store a large number of narrow instruction in its instruction memory, and can load in new instructions in a small number of cycles. Ironically, the fastest way for an FPGA to execute reams of non-loop-intensive code is to build a processor in the FPGA substrate. However, with the extra layer of interpretation, the FPGA's performance will not be comparable to a processor built in the same VLSI process.

### 2.0.3 Our New Concept of a Raw Processor

Based on our conclusions, we arrived at a new model of the architecture, which is described in the September 1997 IEEE Computer "Billion Transistor" issue [Waingold97].

We started with the FPGA design, and added coarse grained functional units, to support datapath operations. We added word-wide data memories to keep frequently used data nearby. We left in some FPGA-like logic to support fine grained applications. We added pipelined sequencers around the functional units to support the reams of non-performance critical code, and to simplify compilation. We linked the sequenced functional units with a statically scheduled pipelined interconnect, to mimic the fast, custom interconnect of ASICs and FPGAs. Finally, we added a dynamic network to support dynamic events.

The end result: a mesh of replicated tiles, each containing a static switch, a dynamic switch, and a small pipelined processor. The tiles are all connected together through two types of high performance, pipelined networks: one static and one dynamic.

Now, two years later, we are on the cusp of building the first prototype of this new architecture.

8

# 3 WHAT WE'RE BUILDING

## 3.0 THE FIRST RAW ARCHITECTURE

In this section, I present a description of the architecture of the Raw prototype, as it currently stands, from an assembly language viewpoint. This will give the reader a more definite feel for exactly how all of the pieces fit together. In the subsequent chapters, I will discuss the progress of design decisions which made the architecture the way it is.

### 3.0.1 A mesh of identical tiles

A Raw processor is a chip containing a 2-D mesh of identical tiles. The tiles are connected to its nearest neighbors by the dynamic and static networks. To program the Raw processor, one programs each of the individual tiles. See the figure entitled "A Mesh of Identical Tiles."

### 3.0.2 The tile

Each tile has a tile processor, a static switch processor, and a dynamic router. In the rest of this document, the tile processor is usually referred to as "the main pro-

cessor," "the processor," or "the tile processor." "The Raw processor" refers to the entire chip -- the networks and the tiles.
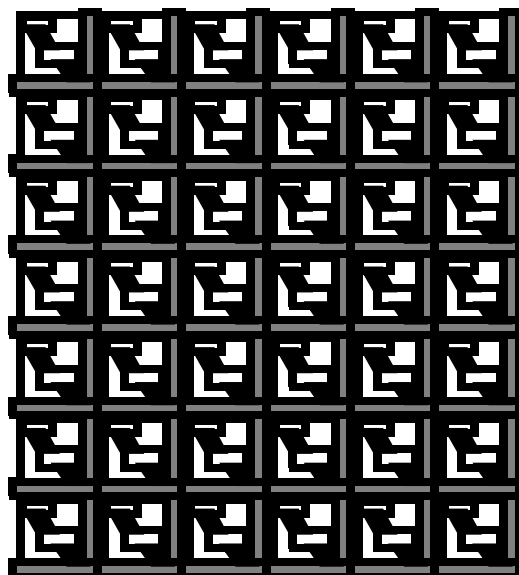
The tile processor uses a 32-bit MIPS instruction set, with some slight modifications. The instruction set is described in more detail in the "Raw User's Manual," which has been appended to the end of this design document.

The switch processor (often referred to as "the switch") uses a MIPS-like instruction set that has been stripped down to contain just moves, branches, jumps and branches. Each instruction also has a ROUTE component, which specifies the transfer of values on the static network between that switch and its neighboring switches.

The dynamic router runs independently, and is under user control only indirectly.

### 3.0.3 The tile processor

The tile processor has a 32 Kilobyte SRAM data memory, and a 32 Kilobyte SRAM instruction memory. The instruction memory is uncached. The data memory can be used in two modes: one cached, and the other uncached. It is the compiler's responsibility to virtualize the memories in software, if caching is disable and 32k is not enough to fit the entire dataset.



**A Mesh of Identical Tiles**



**Logical View of A Raw Tile**

The tile processor communicates with the switch through three ports which have special register names, $csto, $csti2, and $csti. When a data value is written to $csto, it is actually sent to a small FIFO located in the switch. When a data value is read from $csti or $csti2, it is actually read from one of two FIFOs inside the switch. The value is removed from the FIFO when the read occurs. There are two input ports and one output port to match the bandwidth of the ALU.

If a read on $csti or $csti2 is specified, and there is no data available from that port, the processor will block. If a write to $csto occurs, and the buffer space has been filled, the processor will also block.

Here is some sample assembly language:

```
# XOR register 2 with 15,
# and put result in register 31

xori $31,$2,15

# get two values from switch,
# add to register 3, and put
# result in register 9

addu $9,$csti2,$csti

# an ! indicates that the result
# of the operation should also
# be written to $csto

and! $0,$3,$2

# load from address at $csti+25
# put value in register 9 AND
# send it through $csto port
# to static switch

ld! $9,25($csti)

# jump through value specified
# by $csti2

jr $csti2
```

The dynamic network ports operate very similarly. $cgni and $cgno are the general dynamic network ports, while $cmno and $cmni are the memory dynamic network ports and have restricted usage. When writes to $cgno occur, instead of showing up at the static switch, the messages are routed through the chip to their destination tile. This tile is specified by the first word, the header word, that is written into the output port. The

header word also contains mlen, the message length, in words. This and successive words trickle out into the dynamic network, streaming to the other tile. After mlen words are written into the network, the message is complete and the next word that is written into the output port will be interpreted as the destination for a new dynamic message. In all cases, it is preferable that the user send the message words back-to-back, more or less, to minimize the occupancy of the message in the network.

```
# specify a send to tile #15
lui   $3,$0,15

# the message length is 2
ihdr $cgno,$3,0x0200

# put in a couple of datawords,
# one from register 9 and the other
# from the csti network port

or $cgno,$0,$9

# not highly recommend
# for dangling reasons

ld $cgno,$0,$csti

# the message will be
# automatically launched
# into the network

# if we were tile 15, we could
# receive our message with:

# read first word
or $2,$cgni,$0

# read second word,
or $3,$cgni,$0

# the header word is discarded
# by the routing hardware, so
# the recipient does not see it
# there are only two words in
# this message
```

### 3.0.4 The switch processor

The switch processor has a local 8096-instruction instruction memory, but no data memory. This memory is also not cached, and must be virtualized in software by the switch's nearby tile processor.

10

The switch processor executes a very basic instruction set, which consists of only moves, branches, jumps, and nops. It has a small, four element register file. The destinations of all of the instructions must be registers. However, the sources can be network ports. The network port names for the switch processor are $csto, $csti, $csti2, $cNi, $cEi, $cSi, $cWi, $cNi2, $cEi2, $cSi2, $cWi2, $cNo, $cEo, $cSo, $cWo, $cNo2, $cEo2, $cSo2, and $cWo2. These correspond to the main processor's output queue, the main processor's input queues, the input queues coming from the switch's four neighbors, and the output queues going out to the switch's four neighbors.

Each switch processor instruction also has a ROUTE component, which is executed in parallel with the instruction component. If any of the ports specified in the instruction are full (for outputs) or empty (for inputs), the switch processor will stall.

```
# branch instruction
beqz $9, target

# branch if processor
# sends us a zero

beqz $csto, target

# branch if the value coming
# from the west neighbor is a zero

beqz $cWi, target

# store away value from
# east neighbor switch

move $3, $cEi

# same as above, but also route
# the value coming from the north
# port2 to the south port 2

move $3, $cEi   route $cNi2->$cSo2

# all at the same time:
# send value from north neighbor
# to both the south and processor
# input ports.
# send value from processor to west
# neighbor.
# send value from west neighbor to
# east neighbor
```

```
nop  route $cNi->$cSo, $cNi->$csti,
          $csto->$cWo,$cWi->$cEo

#      jump to location specified
# by west neighbor and route that
# location to our east neighbor

jr $cWi     route $cWi->$cEo
```

### 3.0.5 Putting it all together

For each switch-processor, processor-switch, or switch-switch link, the value arrives at the end of the cycle. The code below shows the switch and tile code required for a tile-to-tile send.

```
TILE 0:

    or $csto,$0,$5

SWITCH 0:

    nop  route $csto->$cEo2

SWITCH 1:

    nop  route $cWi2->$csti2

TILE 1:

    and $5, $5, $csti2
```

This code sequence takes five cycles to execute. In the first cycle, tile 0 executes the OR instruction, and the value arrives at switch 0. On the second cycle, switch 0 transmits the value to switch 1. On the third cycle, switch 1 transfers the value to the processor. On the fourth cycle, the value enters the decode stage of the processor. On the fifth cycle, the AND instruction is executed.

Since two of those cycles were spent performing useful computation, the send-to-use latency is three cycles.

More information on programming the Raw architecture can be found in the User's Manual at the end of this design document. More information on how our compiler parallelizes sequential applications for the Raw architecture can be found in [Lee98] and [Barua99].

### 3.1 RAW MATERIALS

Before we decided what we were going to build for the prototype, we needed to find out what resources we had available to us. Our first implementation decision, at the highest level, was to build the prototype as a standard-cell CMOS ASIC (application specific integrated circuit) rather than as full-custom VLSI chip.

In part, I believe that this decision reflects the fact that the group's strengths and interests center more on systems architecture than on circuit and micro-architectural design. If our research shows that our software systems can achieve speedups on our micro-architecturally unsophisticated ASIC prototype, it is a sure thing that the micro-architects and circuit designers will be able to carry the design and speedups even further.

### 3.1.1  The ASIC choice

When I originally began the project, I was not entirely clear on the difference between an ASIC and full-custom VLSI process. And indeed, there is a good reason for that; the term ASIC (application specific integrated circuit) is vacuous.

As perhaps is typical for someone with a liberal arts background, I think the best method of explaining the difference is by describing the experience of developing each type of chip.

In a full-custom design, the responsibility of every aspect of the chip lies on designer's shoulders. The designer starts with a blank slate of silicon, and specifies as an end result, the composition of every unit volume of the chip. The designer may make use of a pre-made collection of cells, but they also are likely to design their own. They must test these cells extensively to make sure that they obey all of the design rules of the process they are using.

These rules involve how close the oxide, poly, and metal layers can be to each other. When the design is finally completed, the designer holds their breath and hopes that the chip that comes back works.

In a standard-cell ASIC process, the designer (usually called the customer) has a library of components that have been designed by the ASIC factory. This library often includes RAMs, ROMs, NAND type primitives, PLLs, IO buffers, and sometimes datapath operators. The designer is not typically allowed to use any other components without a special dispensation. The designer is restricted from straying too far from edge triggered design, and there are upper bounds on the quantity of components that are used (like PLLs). The end product is a netlist of those components, and a floorplan of the larger modules. These are run through a variety of scripts supplied by the manufacturer which insert test structures, provide timing numbers and test for a large number of rule violations. At this point, the design is given to the ASIC manufacturer, who converts this netlist (mostly automatically) into the same form that the full-custom designer had to create.

If everything checks out, the ASIC people and the customer shake hands, and the chip returns a couple of months later. Because the designer has followed all of the rules, and the design has been checked for the violation of those rules, the ASIC manufacturer GAURANTEES that the chip will perform exactly as specified by the netlist.

In order to give this guarantee however, their libraries tend to be designed very conservatively, and cannot achieve the same performance as the full custom versions.

The key difference between an ASIC and full custom VLSI project is that the designer gives up degrees of flexibility and performance in order to attain the guarantee that their design will come back "first time right". Additionally, since much of the design is created automatically, it takes less time to create the chip.

### 3.1.2  IBM: Our ASIC foundry

Given the fact that we had decided to do an ASIC, we looked for an industry foundry. This is actually a relatively difficult feat. The majority of ASIC developers are not MIT researchers building processor prototypes. Many are integrating an embedded system onto one chip in order to minimize cost. Closer to our group in terms of performance requirements are the graphics chips designers and the network switch chip designers. They at least are quite concerned with pushing performance envelope. However, their volumes are measured in the hundreds of thousands, while the Raw group probably will be able to get by on just the initial 30 prototype chips that the ASIC manufacturer gives us. Since the ASIC foundry makes its money off of the volume of the chips produced, we do not make for great profits. Instead, we have to rely on the generosity of the vendor and on other, less tangible incentives to entice a partnership.

We were fortunate enough to be able to use IBM's extraordinary SA-27E ASIC process. It is IBM's latest

ASIC process. It is considered to be a "value" process, which means that some of the parameters have been tweaked for density rather than speed. The "premium," higher speed version of SA-27E is called SA-27.

Please note that all of the information that I present about the process is available off of IBM's website (www.chips.ibm.com) and from their databooks. No proprietary information is revealed in this design document.

**Table 1: SA-27E Process**

| Param | Value |
|---|---|
| $L_{eff}$ | .11 micron |
| $L_{drawn}$ | .15 micron |
| Core Voltage | 1.8 Volts |
| Metallization | 6 layers, copper |
| Gates | Up to 24 Million 2-input NANDs, based on die size |
| Embedded Dram | SRAM MACRO<br> 1 MBit = 8mm$^2$<br><br>DRAM MACRO<br> first 1 MBit = 3.4 mm$^2$<br> addt'l MBits = 1.16 mm$^2$<br> 50 MHz random access |
| I/O | C4 Flip Chip Area I/O up to 1657 pins on CCGA (1124 signal I/Os)<br><br>Signal technologies: SSTL, HSTL, GTL, LVTTL AGP, PCI... |

The 24 million gates number assumes perfect wireability, which although we do have many layers of metal in the process, is unlikely. Classically, I have heard of wireability being quoted at around %35 - %60 for older non-IBM processes.

This means that between %65 and %40 of those gates are not realizable when it comes to wiring up the

design. Fortunately, the wireability of RAM macros is at %100, and the Raw processor is mostly SRAM!

We were very pleasantly surprised by the IBM process, especially with the available gates, and the abundance of I/O. Also, later, we found that we were very impressed with the thoroughness of IBM's LSSD test methodology.

### 3.1.3 Back of the envelope: A 16 tile Raw chip

To start out conservatively, we started out with a die size which was roughly 16 million gates, and assume 16 Raw tiles. The smaller die size gives us some slack at the high end should we make any late discoveries or have any unpleasant realizations. This gave us roughly 1 million gates to allocate to each tile. Of that, we allocated half the area to memory. This amounts to roughly 32 kWords of SRAM, with 1/2 million gates left to dedicate to logic. Interestingly, the IBM process also allows us to integrate DRAM on the actual die. Using the embedded DRAM instead of the SRAM would have allowed us to pack about four times as much memory in the same space. However, we perceived two principal issues with using DRAM:

First, the 50 MHz random access rate would require that we add a significant amount of micro-architectural complexity to attain good performance. Second, embedded DRAM is a new feature in the IBM ASIC flow, and we did not want to push too many frontiers at once.

We assume a pessimistic utilization of %45 for safeness, which brings us to 225,000 "real" gates. My preferred area metric of choice, the 32-bit Wallace tree multiplier, is 8000 gates. My estimate of a processor (with multiplier) is that it takes about 10 32 bit multipliers worth of area. A pipelined FPU would add about 4 multipliers worth of area.

The rest remains for the switch processor and crossbars. I do not have a good idea of how much area they will take (the actual logic is small, but the congestion due to the wiring is of concern) We pessimistically assign the remaining 14 multipliers worth of area to these components.

Based on this back-of-the-envelope calculation, a 16 tile Raw system looks eminently reasonable.

This number is calculated using a very conservative wireability ratio for a process with so many layers of metal. Additionally, should we require it, we have the

34*8 wires

Switch Bus

Switch MEMORY
(8k x 64)

Partial Crossbar

Switch
Processor

FPU

Status Ram (512x8)

Processor

Tag Ram
(512x40)

Processor
Instr Mem
(8kx32)

Data
Memory
(4kx64)

~4 mm

**A Preliminary Tile Floorplan**

possibility of moving up to a larger die. Note however, that these numbers do not include the area required for I/O buffers and pads, or the clock tree. The addition area due to LSSD (level sensitive scan design) is included.
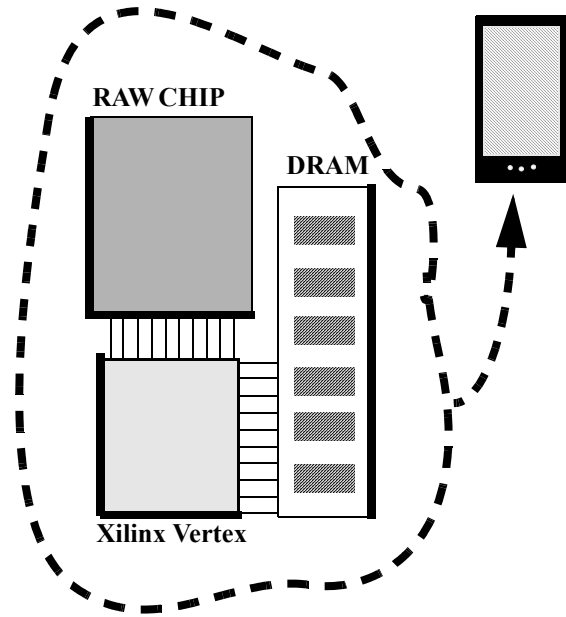
The figure "A Preliminary Tile Floorplan" is a possible floorplan for the Raw tile. It is optimistic because it assumes some flexibility with memory footprints, and the sizes of logic are approximate. It may well be necessary that we reduce the size of the memories to make things fit. Effort has been made to route the large buses over the memories, which is possible in the SA-27E process. This should improve the routability of the processor greatly, because there are few global wires. Because I am not sure of the area required by the crossbar, I have allocated a large area based on the assumption that crossbar area will be proportional to the square of the width of input wires.

In theory, we could push and make it up to 32 tiles. However, I believe that we would be stretching ourselves very thinly -- the RAMs need to be halved (a big problem considering much of our software technology has code expansion effects), and we would have to assume a much better wireability factor, and possibly dump the FPU.

**Table 2: Ballpark clock calculation**

| Structure | Propagation Delay |
|---|---|
| 8192x32 SRAM read | 2.50 ns |
| 2-1 Mux | 0.20 ns |
| Register | 0.25 ns |
| Required slack | 0.50 ns (estimated) |
| | |
| Total | 3.45 ns |

For an estimate on clock speed, we need to be a bit more creative because memory timing numbers are not yet available in the SA-27E databooks. We approximate by using the SA-27 "premium" process databook numbers, which should give us a reasonable upper bound. At the very least, we need to have a path in our processor which goes from i-memory to a 2-1 mux to a register. From the databook, we can see the total in the "Ballpark clock calculation" table.



**A Raw Handheld Device**

The slack is extra margin required by the ASIC manufacturer to account for routing anomalies, PLL jitter, and process variation. The number given is only an estimate, and has no correlation with the number actually required by IBM.

This calculation shows that, short of undergoing micro-architectural heroics, 290 Mhz is a reasonable strawman UPPER BOUND for our WORST-case clock rate.

## 3.2 THE TWO RAW SYSTEMS

Given an estimate of what a Raw chip would look like; we decided to target two systems, a Raw Handheld device, and a Raw Fabric.

### 3.2.1 A Raw Handheld Device

The Raw handheld device would consist of one Raw chip, a Xilinx Vertex, and 128 MB of SDRAM. The FPGA would be used to interface to a variety of peripherals. The Xilinx part acts both as glue logic and as a signal transceiver. Since we are not focusing on the issue of low-power at this time, this handheld device would not actually run off of battery power (well, perhaps a car battery.).

This Raw system serves a number of purposes. First, it is a simple system, which means that it will

15

**A Raw Fabric**

## 3.3 SUMMARY

In this chapter, I described the architecture of the Raw prototype. I elaborated on the ASIC process that we are building our prototype in. Finally, I described the two systems that we are planning to build: a hand-held device, and the multi-chip supercomputer.

make a good test device for a Raw chip. Second, it gets people thinking of the application mode that Raw chips will be used in -- small, portable, extroverted devices rather than large workstations. One of the nice aspects of this device is that we can easily build several of them, and distribute them among our group members. There is something fundamentally more exciting about having a device that we can toss around, rather than a single large prototype sitting inaccessible in the lab. Additionally, it means that people can work on the software required to get the machine running without jockeying for time on a single machine.

### 3.2.2 A Multi-chip Raw Fabric, or Supercomputer

This device would incorporate 16 Raw Chips onto a single board, resulting in 256 MIPS processor equivalents on one board. The static and dynamic networks of these chips will be connected together via high-speed I/O running at the core ASIC speed. In effect, the programmer will see one 256-tile Raw chip.

This would give the logical semblance of the Raw chip that we envisioned for the year 2007, where hundreds of tiles fit on a single die. This system will give us the best simulation of what it means to have such an enormous amount of computing resources available. It will help us answer a number of questions. What sort of applications can we create to utilize these processing resources? How does our mentality and programming paradigm change when a tile is a small percentage of the total processing power available to us? What sort of issues exist in the scalability of such a system? We believe that the per-tile cost of a Raw chip will be so low in the future that every handheld device will actually have hundreds of tiles at their disposal.

16

# 4 STATIC NETWORK DESIGN

## 4.0  STATIC NETWORK

The best place to start in explaining the design decisions of the Raw architecture is with the static network.

The static network is the seed around which the rest of the Raw tile design crystallizes. In order to make efficient fine-grained parallel computation feasible, the entire system had to be designed to facilitate high-bandwidth, low latency communication between the tiles. The static network is optimized to route single-word quantities of data, and has no header words. Each tile knows in advance, for each data word it receives, where it must be sent. This is because the compiler (whether human or machine) generated the appropriate route instructions at compile time.

The static network is a point-to-point 2-D mesh network. Each Raw tile is connected to its nearest neighbors through a series of separate, pipelined channels -- one or more channels in each direction for each neighbor. Every cycle, the tile sequences a small, per-tile crossbar which transfers data between the channels. These channels are pipelined so that no wire requires more than one cycle to traverse. This means that the Raw network can be physically scaled to larger numbers of tiles without reducing the clock rate, because the wire lengths and capacitances do not change with the number of tiles. The alternative, large common buses, will encounter scalability problems as the number of tiles connected to those buses increases. In practice, a hybrid approach (with buses connecting neighbor tiles) could be more effective; however, doing so would add complexity and does not seem crucial to the research results.

The topology of the pipelined network which connects the Raw tiles is a 2-D mesh. This makes for an efficient compilation target because the two dimensional logical topology matches that of the physical topology of the tiles. The delay between tiles is then strictly a linear function of the Manhattan distances of the tiles. This topology also allows us to build a Raw chip by merely replicating a series of identical tiles.

### 4.0.1  Flow Control

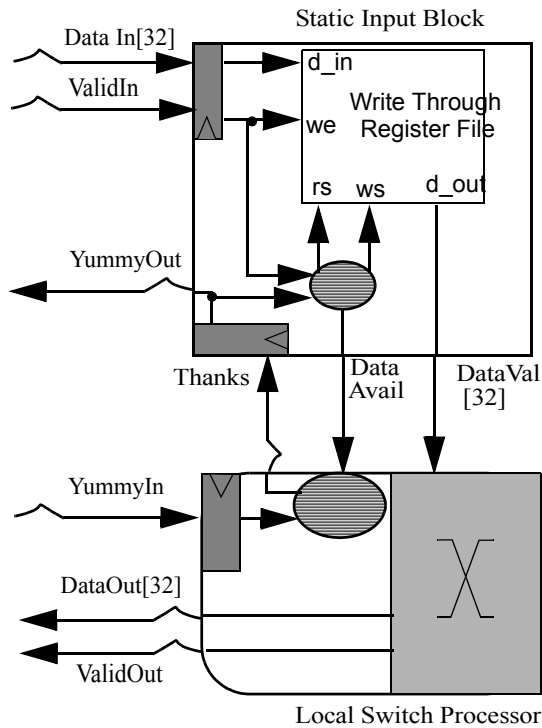Originally, we envisioned that the network would be precisely cycle-counted -- on each cycle, we would know exactly what signal was on which wire. If the compiler were to incorrectly count, then garbage would be read instead, or the value would disappear off of the wire. This mirrors the behaviour of the FPGA prototype that we designed. For computations that have little or no variability in them, this is not a problem. However, cycle-counting general purpose programs that have more variance in their timing behaviour is more difficult. Two classic examples are cache misses and unbalanced if-then-else statements. The compiler could schedule the computation pessimistically, and assume the worst case, padding the best case with special multi-cycle noop instructions. However, this would have abysmal performance. Alternatively, the compiler could insert explicit flow control instructions to handshake between tiles into the program around these dynamic points. This gets especially hairy if we want to support an interrupt model in the Raw processor.

We eventually moved to a flow-control policy that was somewhere between cycle-counting and a fully dynamic network. We call this policy *static ordering* [Waingold97, 2]. Static ordering is a handshake between crossbars which provides flow control in the static network. When the sequencer attempts to route a dataword which has not arrived yet, it will stall until it does arrive. Additionally, the sequencer will stall if a destination port has no space. Delivery of data words in the face of random delays can then be guaranteed. Each tile still knows a priori the destination and order of each data word coming in; however, it does not know exactly which cycle that will be. This contrasts with a dynamic network, where neither timing nor order are known a priori. Interestingly, in order to obtain good performance, the compiler must cycle count when it schedules the instructions across the Raw fabric. However, with static ordering, it can do so without worrying that imperfect knowledge of program behaviour will violate program correctness.

The main benefits of adding flow control to the architecture are the abstraction layer that it provides and the added support for programs with unpredictable timing. Interestingly, the Warp project at CMU started without flow control in their initial prototypes, and then added it in subsequent revisions [Gross98]. In the next section, we will examine the static input block, which is the hardware used to implement the static ordering protocol.

### 4.0.2  The Static Input Block

The static input block (SIB) is a FIFO which has both backwards and forwards flow control. There is a

Static Input Block

**Static Input Block Design**

local SIB at every input port on the switch's crossbar. The switch's crossbar also connects to a remote input buffer that belongs to another tile. The figure "Static Input Block Design" shows the static input block and switch crossbar design. Note that an arrow that begins with a squiggle indicates a signal which will arrive at its destination at the end of the cycle. The basic operation of the SIB is as follows:

1. Just before the clock edge, the `DataIn` and `ValidIn` signals arrive at the input flops, coming from the remote switch that the SIB is connected to. The `Thanks` signal arrives from the local switch, indicating if the SIB should remove the item at the head of the fifo. The `Thanks` signal is used to calculate the `YummyOut` signal, which gives the remote switch an idea of how much space is left in the fifo.

2. If `ValidIn` is set, then this is a data word which must be stored in the register file. The protocol ensures that data will not be sent if there is no space in the circular fifo.

3. `DataAvail` is generated based on whether the fifo is empty. The head data word of the queue is propagated out of `DataVal`. These signals travel to the switch.

4. The switch uses `DataAvail` and `DataVal` to perform its route instructions. It also uses the `YummyIn` information to determine if there is space on the remote side of the queue. The `DataOut` and `ValidOut` signals will arrive at a remote input buffer at the end of the cycle.

5. If the switch used the data word from the SIB, it asserts `Thanks`.

The subtlety of the SIB comes from that fact that it is a distributed protocol. The receiving SIB is at least one cycle away from the switch that is sending the value. This means that the sender does not have perfect information about how much space is available on the receiver side. As a result, the sender must be conservative about when to send data, so as not to overflow the fifo. This can result in suboptimal performance for streams of data that are starting out, or are recovering from a blockage in the network. The solution is to add a sufficient number of storage elements to the FIFO.

The worksheets "One Element Fifo" and "Three Element Fifo" help illustrate this principle. They show the state of the system after each cycle. The left boxes are a simplified version of the switch circuit. The right boxes are a simplified version of a SIB connected to a remote switch. The top arrow is the `ValidIn` bit, and the bottom arrow is the "`Yummy`" line. The column of numbers underneath "PB" (perceived buffers) are the switch's conservative estimate of the number of elements in the remote SIB at the beginning of the cycle. The column of numbers underneath "AB" (actual buffers) are the actual number of elements in the fifo at the beginning of the cycle.

The two figures model the "Balanced Producer-Consumer" problem, where the producer is capable of producing data every cycle, and the consumer is capable of consuming it every cycle. This would correspond to a stream of data running across the Raw tiles. Both figures show the cycle-by-cycle progress of the communication between a switch and its SIB.

We will explain the "One Element Fifo" figure so the reader can get an idea of how the worksheets work. In the first cycle, we can see that the switch is asserting its `ValidOut` line, sending a data value to the SIB. On the second cycle, the switch stalls because it knows that the Consumer has an element in its buffer, and may not have space if it sends a value. The `ValidOut` line is thus held low. Although it is not indicated in the diagram, the Consumer consumes the data value from the previous cycle. On the third cycle, the SIB asserts the YummyOut line, indicating that the value had been con-

**One Element Fifo**

**Three Element Fifo**

sumed. However, the Switch does not receive this value until the next cycle. Because of this, the switch stalls for another cycle. On the fourth cycle, the switch finally knows that there is buffer space and sends the next value along. The fifth and sixth cycles are exactly like the second and third.

Thus, in the one element case, the static switch is stalling because it cannot guarantee that the receiver will have space. It unfortunately has to wait until it receives notification that the last word was consumed.

In the three element case, the static network and SIBs are able to achieve optimal throughput. The extra storage allows the sender to send up to three times before it hears back from the input buffer that the first value was consumed. It is not a coincidence that this is also the round trip latency from switch to SIB. In fact, if Raw were moved to a technology where it took multiple cycles to cross the pipelined interconnect between tiles (like for instance, for the Raw multi-chip system), the number of buffers would have to be increased to match

the new round trip latency. By looking at the diagram, you may think that perhaps two buffers is enough, since that is the maximum perceived element size. In actuality, the switch would have to stall on the third cycle because it perceives 2 elements, and is trying to send a third out before it received the first positive "Yummy-Out" signal back.

The other case where it is important that the SIBs perform adequately is in the case where there is head-of-line blocking. In this instance, data is being streamed through a line of tiles, attaining the steady state, and then one of the tiles at the head becomes blocked. We want the SIB protocol to insure that the head tile, when unblocked, is capable of reading data at the maximum rate. In other words, the protocol should insure that no bubbles are formed later down the pipeline of producers and consumers. The "Three Element Fifo, continued" figure forms the basis of an inductive proof of this property.

**Three Element Fifo, continued**
Starts at Steady State, then Head blocks (stalls) for four cycles

I will elaborate on "Three Element Fifo, continued," some more. In the first cycle, the "BLOCK" indicates that no value is read from the input buffer at the head of the line on that cycle. After one more BLOCKs, in cycle three, the switch behind the head of the line STALLs because it correctly believes that its consumer has run out of space. This stall continues for three more cycles, when the switch receives notice that a value has been dequeued from the head of the queue. These stalls ripple down the chain of producers and consumers, all offsetted by two cycles.

It is likely that even more buffering will provide greater resistance to the performance effects of blockages in the network. However, every element we add to the FIFO is an element that will have to be exposed for draining on a context switch. More simulation results could tell us if increased buffering is worthwhile.

**The Unified Approach**

63      MIPS instruction      32      route instruction      0

| op | S | rs | rt | imm | N | E | S | W | P | other |
|----|---|----|----|-----|---|---|---|---|---|-------|

**The Unified Approach**

---

**The Slave Processor Approach**

63      48      32

| op | S | rs | rt | imm |
|----|---|----|----|-----|

MIPS Instruction

63      48      32      26

| op | imm | | | | N | E | S | W | P |
|----|-----|---|---|---|---|---|---|---|---|

Switch Instruction

**The Slave Processor Approach**

### 4.0.3 Static Network Summary

The high order bit is that adding flow control to the network has resulted in a fair amount of additional complexity and architectural state. Additionally, it adds logic to the path from tile to tile, which could have performance implications. With that said, the buffering allows our compiler writers some room to breath, and gives us support for events with unpredictable timing.

### 4.1 THE SWITCH (SLAVE) PROCESSOR

The switch processor is responsible for controlling the tile's static crossbar. It has very little functionality -- in some senses one might call it a "slave parallel move processor," since all it can do is move values between a small register file, its PC, and the static crossbar.

One of the main decisions that we made early on was whether or not the switch processor would exist at all. Currently, the switch processor is a separately sequenced entity which connects the main processor to the static network. The processor cannot access the static network without the slave processor's cooperation.

A serious alternative to the slave-processor approach would have been to have only the main processor, with a VLIW style processor word which also specified the routes for the crossbar. The diagram "The Unified Approach" shows an example instruction encoding. Evaluating the trade-offs of the unified and slave designs is difficult.

A clear disadvantage of the slave design is that it is more complicated. It is another processor design that we have to do, with its own instruction encoding for branches, jumps, procedure calls and moves for the register file. It also requires more bits to encode a given route.

The main annoyance is that the slave processor requires constant baby-sitting by the main processor. The main processor is responsible for loading and unloading the instruction memory of the switch on cache misses, and for storing away the PCs of the switch on a procedure call (since the switch has no local storage). Whenever the processor takes a conditional branches, it needs to forward the branch condition on to the slave processor. The compiler must make sure there is a branch instruction on the slave processor which will interpret that condition.

Since the communication between the main and slave processors is statically scheduled, it is very difficult and slow to handle dynamic events. Context switches require the processor to freeze the switch, set the PC to an address which drains the register files into the processor, as well as any data outstanding on the switch ports.

The slave switch processor also makes it very difficult to use the static network to talk to the off-chip network at dynamically chosen intervals, for instance, to read a value from a DRAM that is connected to the static network. This is because the main processor will have to freeze the switch, change the switch's PC, and

then unfreeze it.

The advantages of the switch processor come in tolerating latency. It decouples the processing of networking instructions and processor instructions. Thus, if a processor takes longer to process an instruction than normal (for instance on a cache miss), the switch instructions can continue to execute, and visa versa. However, they will block when an instruction is executed that requires communication between the two. This model is reminiscent of Decoupled-Execute Access Architectures [Smith82].

The Unified approach does not give us any slack. The instruction and the route must occur at precisely the same time. If the processor code takes less time than expected, it will end up blocked waiting for the switch route to complete. If the processor code takes more time than expected, a "through-route" would be blocked up on unrelated computation. The Unified approach also has the disadvantage that through route instructions must be scheduled on both sides of an if-statement. If the two sides of the if-statement were wildly unbalanced this would create code bloat. The Slave approach would only need to have one copy of the corresponding route instructions.

In the face of a desire for this decoupling property, we have further entertained the idea of another approach, called the Decoupled-Unified approach. This would be like the Unified approach, except it would involve having a queue through which we would feed the static crossbar its route instructions. This is attractive because it would decouple the two processes. The processor would sequence, and queue up switch instructions, which would execute when ready.

With this architecture, the compiler would push the switch instructions up to pair with the processor instructions at the top of a basic block. This way through-routes could execute as soon as possible.

Switch instructions that originally ran concurrently with non-global IF-ELSE statements need some extra care. Ideally, the instructions would be propagated above the IF-ELSE statement. Otherwise, the switch instructions will have to be copied to both sides the IF-ELSE clause. This may result in code explosion, if the number of switch instructions propagated into the IF-ELSE statement is greater than the length of one of the sides of the statement.

When interrupts are taken into account, the Decoupled-Unified approach is a nightmare, because now we have situations where half of the instruction (the processor part) has executed. We can not just wait for the switch instructions to execute, because this may take an indefinite amount of time.

To really investigate the relative advantages and disadvantages of the three methods would require an extensive study, involving modifications of our compilers and simulators. To make a fair comparison, we would need to spend as much time optimizing the comparison simulators as we did the originals. In an ideal world, we might have pursued this issue more. However, given the extensive amount of infrastructure that had already been built using the Slave model, we could not justify the time investment for something which was unlikely to buy us performance, and would require such an extensive reengineering effort.

### 4.1.1 Partial Routes

One idea that our group members had was that we do not need to make sure that all routes specified in an instruction happen simultaneously. They could just fire off when they are possible, with that part of the instruction field resetting itself to the "null route." When all fields are set to null, that means we can continue onto the next instruction. This algorithm continues to preserves the static ordering property.

From a performance and circuit perspective, this is a win. It will decouple unrelated routes that are going through the processor. Additionally, the stall logic in the switch processor does not need to OR together the success of all of the routes in order to generate the "ValidOut" signal that goes to the neighboring tile.

The problem is, with partial routes, we again have an instruction atomicity problem. If we need to interrupt the switch processor, we have no clear sense of which instruction we are currently at, since parts of the instruction have already executed. We cannot wait for the instruction to fully complete, because this may take indefinite amount of time. In order to make this feature, we would have had to add special mechanisms to overcome this problem. As a result, we decided to take the simple path and stall until such a point as we can route all of the values atomically.

### 4.1.2 Virtual Switch Instruction Memory

In order to be able to run large programs, we need a mechanism to page code in and out of the various memories. The switch memory is a bit of an issue because it

is not coupled directly with the processor, and yet it does not have the means to write to its own memory. Thus, we need the processor to help out in filling in the switch memory.

There are two approaches.

In the first approach, the switch executes until it reaches a "trap" instruction. This trap instruction indicates that it needs to page in a new section of memory. The trap causes an interrupt in the processor. The processor fetches the relevant instructions and writes it into the switch processor instruction memory. It then signals the switch processor, telling it to resume.

In the second approach, we maintain a mapping between switch and processor instruction codes. When the processor reaches a junction where it needs to pull in some code, it pulls in the corresponding code for the switch. The key issue is to make sure that the switch does not execute off into the weeds while this occurs. The switch can very simply do a read from the processor's output port into its register set (or perhaps a branch target.) This way, the processor can signal the switch when it has finished writing the instructions. When the switch's read completes, it knows that the code has been put in place. Since there essentially has to be a mapping between the switch code and the processor code if they communicate, this mapping is not hard to derive. The only disadvantage is that due to the relative sizes of basic blocks in the two memories, it may be the case that one needs to page in and the other doesn't. For the most part I do not think that this will be much of a problem. If we want to save the cost of this output port read after the corresponding code has been pulled in, we can re-write that instruction.

In the end, we decided on the second option, because it was simpler. The only problem we foresee is if the tile itself is doing a completely unrelated computation (and communicating via dynamic network.) Then, the switch, presumably doing through routes, has no mechanism of telling the local tile that it needs new instructions. However, presumably the switch is synchronized with at least one tile on the chip. That tile could send a dynamic message to the switch's master, telling it to load in the appropriate instructions. We don't expect that anyone will really do this, though.

## 4.2 STATIC NETWORK BANDWIDTH

One of the questions that needs to be answered is how much bandwidth is needed in the static switch. Since a ALU operation typically has two inputs, having

only one $csti port means that one of the inputs to the instruction must reside inside the tile to not be bottlenecked. The amount of bandwidth into the tile determines very strongly the manner in which code is compiled to it. As it turns out, the RAWCC compiler optimizes the code to minimize communication, so it is not usually severely affected by this bottleneck. However, when code is compiled in a pipeline fashion across the Raw tiles, more bandwidth would be required to obtain full performance.

The network ports csti2, cNi2, cSi2, cEi2, and cWi2 have been included in this spec since it is typically easier to remove a feature than to add one at a later date. It remains to be evaluated what the speedup numbers and area (both static instruction memory, crossbar and wire area) and clock cycle costs are for this feature. As it turns out, the encoding for this fits neatly in a 64-bit switch instruction word. The inclusion of this extra routing facility emphasizes the availability of massive amounts of bandwidth inside the chip.

## 4.3 SUMMARY

The static network design makes a number of important trade-offs. The network flow control protocol contains flow-controlled buffers that allow our compiler writers some room to breath, and gives us support for events with unpredictable timing. This protocol is a distributed protocol in which the producers have imperfect information. As a result, the SIBs require a small amount of buffering to prevent delay. In this chapter, I presented a simple method for calculating how big these buffer sizes need to be in order to allow continuous streams to pass through the network bubble-free.

The static switch design also has some built-in slack for dynamic timing behaviour between the tile processor and the switch processor. This slack comes with the cost of added complexity.

Finally, we raised the issue of the static switch bandwidth.

All in all, the switch design is a success; it provides an effective low-latency network for inter-tile communication. In the next section, we will see how the static network is interfaced to the tile's processor.

# 5 DYNAMIC NET-WORK PRIMER

## 5.0 DYNAMIC NETWORK

Shortly after we developed the static network, we realized the need for the dynamic network. In order for the static network to be a high performance solution, the following must hold:

1. The destinations must be known at compile time.

2. The message sizes must be known at compile time.

3. For any two communication routes that cross, the compiler must be able generate a switch schedule which merges those two communication patterns on a cycle by cycle basis.

The static network can actually support messages which violate these conditions. However, doing this requires an expensive layer of interpretation to simulate a dynamic network.

The dynamic network was added to the architecture to provide support for messages which do not fulfill these criteria.

The primary intention of the dynamic network is to support memory accesses that cannot be statically analyzed. The dynamic network was also intended to support other dynamic activities, like interrupts, dynamic I/O accesses, speculation, synchronization, and context switches. Finally, the dynamic network was the catch-all safety net for any dynamic events that we may have missed out on.

In my opinion, the dynamic network is probably the single most complicated part of the Raw architecture. Interestingly enough, the design of the actual hardware is quite straight-forward. Its interactions with other parts of the system, and in particular, the deadlock issues, can be a nightmare if not handled correctly. For more discussions on the deadlock issues, please refer to the section entitled "Deadlock."

## 5.1 SUMMARY

The dynamic network design leveraged many of the same underlying hardware components as the static switch design. Its performance is not as good as the static network's because the route directions are not known a priori. A great deal more will be said on the dynamic network in the Deadlock section of this design document.

# 6 TILE PROCESSOR DESIGN

When we first set out to define the architecture, we chose the 5-stage MIPS R2000 as our baseline processor for the Raw tile. We did this because it has a relatively simple pipeline, and because many of us had spent hundreds of hours staring at that particular pipeline. The R2000 is the canonical pipeline studied in 6.823, the graduate computer architecture class at MIT. The discussion that follows assumes familiarity with the R2000 pipeline. For an introduction, see [Hennessey96]. (Later, because of the floating point unit, we expanded the pipeline to six stages.)



$csto, Bypass, and Writeback Networks

## 6.0 NETWORK INTERFACE

The most important part of the main processor decision is the way in which it interfaces with the networks. Minimizing the latency from tile to tile (especially on the static network) was our primary goal. The smaller the latency, the greater the number of applications that can be effectively parallelized on the Raw chip.

Because of our desire to minimize the latency from tile to tile, we decided that the static network interface should be directly attached to the processor pipeline. An alternative would have been to have explicit MOVE instructions which accessed the network ports. Instead, we wanted a single instruction to be able to read a value from the network, operate on it, and write it out in the same cycle.

We modified the instruction encodings in two ways to accomplish this magic.

For writes to the network output port SIB, $csto, we modified the encoding the MIPS instruction set to include what we call the "S" bit. The S bit is set to true if the result of the instruction should be sent out to the output port, in addition to the destination register. This allows us to send a value out of the network and keep it locally. Logically, this is useful when an operation in the program dataflow graph has a fanout greater than one. We used one of the bits from the opcode field of the original MIPS ISA to encode this.

For the input ports, we mapped the network port names into the register file name space:

| Reg | Alias | Usage |
|------|-----------|-------------------------------------------|
| $24 | $csti | Static network input port. |
| $25 | $cgn[i/o] | General ("User") Dynamic network input port. |
| $26 | $csti2 | Second static network port |
| $27 | $cmn[i/o] | Memory ("High Priority") Dynamic network port. |

This means, for instance, that when register $24 is referenced, it actually takes the result from the static network input SIB.

With the current 5-bit addressing of registers, additional register names would only be possible by adding

one more bit to the register address space. Aliasing it with an existing register name allows us to leave most of the ISA encodings unaffected. The choice of the register numbers was suggested by Ben Greenwald. He believes that we can maximize our compatibility with existing MIPS tools by reserving that particular register because it has been designated as "temporary."

## 6.1 SWITCH BYPASSING

The diagram entitled "$csto, Bypass and Writeback Networks" shows how the network SIBs are hooked up to the processor pipeline. The three muxes are essentially bypass muxes. The $csti and $cgni SIBs are logically in the decode/register fetch (RF) stage.

In order to reduce the latency of a network send, it was important that an instruction deliver its result to the $csto SIB as soon as the value was available, rather than waiting until the writeback stage. This can change the tile-to-tile communication latency from 6 cycles to 3 cycles.

The $cXXo ($cmno, $csto and $cgno) SIBs are connected to the processor pipeline in much the same way that register bypasses are connected. Values can be sent to $csto after the ALU stage, after the MEMORY stage, after the FPU stage, and at the WB stage. This gives us the minimum possible latency for all operations whose destination is the static network. The logic is very similar to the bypassing logic; however the priority of the elements is reversed: $csto wants the OLDEST value from the pipeline, rather than the newest one.

When a instruction that writes to $csto is executed, the S bit travels with it down the pipeline. A similar thing happens with a write to $cgno, except that the "D" bit is generated by the decode logic. Each cycle, the $csto bypassing logic finds the oldest instruction which has the S bit set. If that instruction is not ready, then the valid bit connecting to the output SIB is not asserted. If the oldest instruction has reached its stage of maturation (i.e., the stage at which the result of the computation is ready), then the value is muxed into the $csto port register, ready to enter into an input buffer on the next cycle. The S bit of that instruction is cleared, because the instruction has sent its value. When the instruction reaches the Writeback stage, it will also write its result into the register file.

It is interesting to note that the logic for this procedure is exactly the same as for the standard bypass logic, except that the priorities are reversed. Bypass logic favors the youngest instruction that is writing a particular value. $csto bypassing logic looks for the oldest

instruction with the S bit set because it wants to guarantee that values are sent out of the network in order that the instructions were issued.

The $cXXi ($cgni, $cmni, $csti2, and $csti) network ports are muxed in through the bypass muxes. In this case, when an instruction in the decode stage uses registers $24, $25, $26, or $27 as a source, it checks if the DataAvail signal of the SIB is set. If it is not, then the instruction stalls. This mirrors a hardware register interlock. If the decode stage decides it does not have to stall, it will acknowledge the receipt of the data value by asserting the appropriate Thanks line.

### 6.1.1 Instruction Restartability

The addition of the tightly coupled network interfaces does not come entirely for free. It imposes a number of restrictions on the operation of the pipeline.

The main issue is that of restartability. Many processor pipelines take advantage of the fact that their instruction sets are restartable. This means that the processor can squash the instruction at any point in the pipeline before the writeback stage. Unfortunately, instructions which access $cXXi modify the state of the networks. Similarly, when an instruction issues an instruction which writes to $cXXo, once the result has been sent out to the switch's SIB, it is beyond the point of no return and cannot be restarted.

Because of this, the commit point of the tile processor is right after it passes the decode stage. We have to be very careful about instructions that write to $csto or $cgno because the commit point is so early in the pipeline. If we allow the instructions to stall (because the output queues are full) in a stage beyond the decode stage, then the pipeline could be stalled indefinitely, This is because it is programmatically correct for the output queue to be full indefinitely. At that point, the processor cannot take an interrupt, because it must finish all of the "committed" instructions that passed decode.

Thus, we must also insure that if an instruction passes decode, it must not be possible for it to stall indefinitely.

To avoid these stalls, we do not let instruction pass decode unless there is guaranteed to be enough room in the appropriate SIB. As you might guess, we need to use the same analysis as we used to calculate how much buffer space we needed in the network SIBs. Having the correct number of buffers will ensure that the processor

is not too conservative. Looking at the "$csto, Bypass, and Writeback Networks", diagram, we count the number of pipeline registers in the longest cycle from the decode stage through the Thanks line, back to the decode stage. Six buffers are required.

An alternative to this subtle approach is that we could modify the behaviour of the SIBs. We can keep the values in the input SIB FIFOs until we are sure we do not need them any more. Each SIB FIFO will have three pointers: one marks the place where data should be inserted, the next marks where data should be read from, and the final one marks the position of the next element that would be committed. If instructions ever need to be squash, the "read" pointer can be reset to equal the "commit" pointer. I do not believe that this would affect the critical paths significantly, but the $csti and $cgni SIBs would require nine buffers each instead of three.

For the output SIB, creating restartability is a harder problem. We would have to defer the actual transmittal of the value through the network until the instruction has hit WRITEBACK. However, that would mean that we could not use our latency reducing bypassing optimization. This approach mirrors what some conventional microprocessors do to make store instructions restartable -- the write is deferred until we are absolute sure we need it. An alternative is to have some sort of mechanism which overrides a message that was already sent into the network. That sounded complicated.

### 6.1.2  Calculating the Tile-to-Tile Communication Latency

A useful exercise is to examine the tile-to-tile latency of the network send. The figure "Processor-Switch-Switch-Processor" path helps illustrate it. It shows the pipelines of two Raw tiles, and the path over the static network between them. The relevant path is in bold. As you can see, it takes three cycles for nearest neighbor communication.

It is possible that we could reduce the cost down to two cycles. This would involve removing the register in front of the $csti SIB, and rearranging some of the logic. We can do this because we know that the path between the switch's crossbar and the SIB is on the same tile, and thus short. However, it is not at all clear that this will not lengthen the critical path in the tile design. Whether we will be able to do this or not will become more apparent as we come closer to closing the timing issues of our verilog.

### 6.2 MORE STATIC SWITCH INTERFACE GOOK

A number of other items are required to make the static switch and main processor work together.

The first is a mechanism to write and read from the static network instruction memory. The `sload` and `sstore` operations stall the static switch for a cycle.

Another mechanism allows us to freeze the switch. This lets the processor inspect the state at its leisure. It also simplifies the process of loading in a new PC.

During context switches and booting, it is useful to be able to see how many elements are in the switch's SIBs. There is a status register in the processor which can be read to attain this information.

Finally, there is a mechanism to load in a new PC, for context switches, or if we want the static switch to do something dynamic on our behalf.

### 6.3 MECHANISM FOR READING AND WRITING INTO INSTRUCTION MEMORY

In order for us to change the stored program, we need some way of writing values into the instruction memory. Additionally, however, we want to be able to read all of the state out of the processor (which includes the instruction memory state), and we would like to support research into a sophisticated software instruction VM system. As such, we need to be able to treat the instruction memory as a true read-write memory. The basic thinking on this issue is that we will support two new instructions -- "iload" and "istore" -- which mimic the data versions but which access the instruction memory. The advantage of these instructions is that it makes it very explicit when we are doing things which are not standard, both in the hardware implementation and in debugging software. These instructions will perform their operations in the "memory" stage of the pipeline, stealing a cycle away from the "fetch" stage. This means that every read or write into instruction memory will cause a one cycle stall. Since this is not likely to be a common event, we will not concern ourselves with the performance implications.

Associated with an instruction write will be some window of time (i.e. two or three cycles unless we add in some sort of instruction prefetch, then it would be more) where an instruction write will not be reflected in the processor execution. I.E., instructions already fetched into the pipeline will not be refetched if they happen to be the ones that were changed. This is a standard caveat made by most processor architectures.

The Processor-Switch-Switch-Processor path

We also considered the alternative of using standard "load" and "store" instructions, and using a special address range, like for instance ("0xFFFFxxxx"). This approach is entirely valid and has the added benefit that standard routines ("memcpy") will be able to modify instruction memory without having special version. (If we wanted true transparency, we'd have to make sure that instruction memory was accessible by byte accesses.) We do not believe this to be a crucial requirement at this time. If needbe, the two methods could also easily co-exist.

## 6.4 RANDOM TWEAKS

Our baseline processor was the MIPS R2000. We added load interlocks into the architecture, because they aren't that costly. Instead of a single multi-cycle multiply instruction, there are three low-latency pipelined instructions, MULH, MULHU, and MULLO which place their results in a GPR instead of HI/LO. We did this because our 32-bit multiply takes only two cycles. It didn't make sense to treat it as a multi cycle instruction when it has no more delay than a load. We also removed the SWL and SWR instructions, because we didn't feel they were worth the implementation complexity.

We have a 64 bit cycle counter which lists the number of cycles since reset. There is also a watchdog timer, which is discussed in the DEADLOCK section of the

design document.

Finally, we decided on a Harvard style architecture, with separate instruction and data memories; because the design of the pipeline was more simple. See the Appendage entitled "Raw User's Manual" for a description of the instruction set of the Raw prototype. The first Appendage shows the pipeline of the main processor.

## 6.5 THE FLOATING POINT UNIT

In the beginning, we were not sure if we were going to have a floating point unit. The complexity seemed burdensome, and there were some ideas of doing it in software. One of our group members, Michael Zhang, implemented and parallelized a software floating point library [Zhang99] to evaluate the performance of a software solution. Our realization was that many of our applications made heavy use of floating point, and for that, there is no subsitute for hardware. We felt that the large dynamic range offered by floating point would further the ease of writing signal processing applications -- an important consideration for enticing other groups to make user of our prototype. This was an important consideration To simplify our task, we relaxed our compliance of the IEEE 754 standard. In particular, we do not implement gradual underflow. We decided to support only single-precision floating point operations so we would not need to worry about how to integrate a 64 bit datapath into the RAW processor. All of the network paths are 32bits, so we would have package up values and route them, reassemble them and so on. However, if we were building an industrial version, we would probably have a 64 bit datapath throughout the chip, and double precision would be easier to realize.

It was important that the FPU be as tightly integrated with the static network as the ALU. In terms of floating point, Raw had the capability of being a supercomputer even as an academic project. With only a little extra effort in getting the floating point right, we could make Raw look very exciting.

We wanted to be able to send data into the FPU in a pipelined fashion and have it stream out of the tile just as we would do with a LOAD instruction. This would yield excellent performance with signal processing codes, especially with the appropriate amount of switch bandwidth. The problem that this presented was with the $csto port. We need to make sure that values exit the $csto port in the correct order from the various floating point functional units, and from the ALU.

The other added complexity with the floating point unit is the fact that its pipeline is longer than the corresponding ALU pipeline. This means that we needed to do some extra work in order to make sure that items are stored back correctly in the writeback phase, and that they are transferred into the static network in the correct order.

The solution that we used was simple and elegant. After researching FPU designs [Oberman96], it became increasingly apparent that we could do both floating point pipelined add and multiply in three cycles. The longest operation in the integer pipeline is a load or multiply, which is two cycles. Since they are so close, we discovered that we could solve both the $csto and register file writeback problems by extending the length of the overall pipeline by one cycle. As a result, we have six pipeline stages: instruction fetch(IF), instruction decode(ID), execution(EXE), memory(MEM), floating point (FPU) and write-back(WB). See Appendix B for a diagram of the pipeline. The floating point operations execute during the Execute, Memory, and FPU stages, and write back at the same stage as the ALU instructions.

This solves the writeback and $csto issues -- once the pipelines are merged, the standard bypass and stall logic can be used to maintain sanity in the pipeline.

This solution becomes more and more expensive as the difference in actual pipeline latencies of the instructions grows. Each additional stage requires at least one more input to the bypass muxes.

As it turns out, this was also useful for implementing byte and half-word loads, which use an extra stage after the memory stage.

Finally, for floating point division, our non-pipelined 11-cycle divider uses the same decoupled HI/LO interface as the integer divide instruction.

A secondary goal we had in designing an FPU is that we make the source available for other research projects to use. Our design is constructed to be extremely portable, and will probably make its way onto the web in the near future.

## 6.6 RECONFIGURABLE LOGIC

Originally, each Raw tile was to have reconfigurable logic inside, to support bit-level and byte-level computations. Although no research can definitely say that this is a bad idea, we can say that we had a number of problems realizing this goal. The first problem is that

we had trouble finding a large number of applications that benefited enormously from this functionality. Median filter and Conway's "game of life" [Berklekamp82] were the top two contenders. Although this may seem surprising given RawLogic's impressive results on many programs, much of RawLogic's performance came from massive parallelism, which the Raw architecture leverages very capably with tile-level parallelism. Secondly, it was not clear if a reconfigurable fabric could be efficiently implemented on an ASIC. Third, interfacing the processor pipeline to the reconfigurable logic in a way that effectively used the reconfigurable logic proved difficult. Fourth, it looked as if a large area would need to be allocated to each reconfigurable logic block to attain appreciable performance gains. Finally, and probably most fundamentally for us, the complexity of the reconfigurable logic, its interface, and the software system was an added burden to the implementation of an already quite complicated chip.

For reference, here is the description of the reconfigurable logic interface that we used in the first simulator:

The pipeline interface to the reconfigurable logic mimiced the connection to the dynamic network ports. There were two register mapped ports, RLO (output to RL) and RLI (input from RL to processor). These were aliased with register 30. There was a two element buffer on the RLI connection on the processor pipeline side, and a two element buffer on the reconfigurable logic input side.

## 6.7 CGNO Commit Buffer

This feature has been deprecated from the architecture. We retain the text for reference.

The processor initiates a dynamic network send by writing the destination tile number, writing the message into the $cgno commit buffer and then executing the `dlaunch` instruction [Kubiatowicz98]. $cgno is different than other SIBs because it buffers up an entire message until the message is complete. If we were to allow the messages to be injected directly into the network without queueing them up into atomic units, we could have a phenomenon we call dangling. This means that a half-constructed message is hanging out into the dynamic network. Dangling becomes a problem when interrupts occur. The interrupt handler may want to use the dynamic network output queue; however, there is a half-completed message that is blocking up the network port. The message cannot be squashed because some of the words have already been transmitted. A similar problem occurs with context switches -- to allow dan-

gling, the context switch routine would need to save and restore the internal state of the hardware Dynamic scheduler -- a prospect we do not relish. The commit buffer has to be of a fixed size. This size imposes a maximum message size constraint on the dynamic network. To reduce the complexity of the commit buffer, a write to $cgno blocks until all of the elements of the previous message have drained out.

One alternative to the commit buffer would be to require the user to enclose their dynamic network activity to constrained regions surround by interrupt enables and disables. The problem with this approach is that the tile may block indefinitely because the network queue is backed up (and potentially for a legitimate reason.) That would make the tile completely unresponsive to interrupts.

$cgni, on the other hand, operates exactly like the $csti port. However, there is a mask which, when enabled, causes a user interrupt routine to be called when the header of a message arrives at the tile.

## 6.8 SUMMARY

The Raw tile processor design descended from the MIPS R2000 pipeline design. The most interesting design decisions involved the integration of the network interfaces. It was important that these interfaces (in particular the static network interface) provide the minimal possible latency to the network so as to support as fine-grained parallelism as possible.

# 7 I/O AND MEMORY SYSTEM

## 7.0  THE I/O SYSTEM

The I/O system of a Raw processor is a crucial but up until now mostly unmentioned aspect of Raw. The Raw I/O philosophy mirrors that of the Raw parallelism philosophy. Just as we provide a simple interface for the compiler to exploit the gobs of silicon resources, we also have a simple interface for the compiler to exploit and program the gobs of pins available. Once again, the Raw architecture proves effective not because it allocates the raw pin resources to special purpose tasks, but because it exposes them to the compiler and user to meet the needs of application. The interface that we show scales with the number of pins, and works even though pin counts are not growing as fast as logic density.

An effective parallel I/O interface is especially important for a processor with so many processing resources. To support extroverted computing, a Raw architecture's I/O system must be able to interface to, at high-speed, a rich variety of input and output devices, like PCI, DRAM, SRAM, video, RF digitizers and transmitters and so on. It is likely, that in the future, a Raw device would also have direct analog connections - RF receivers and transmitters, and A/D and D/A converters, all exposed to the compiler. However, the integration of analog devices onto a silicon die is the subject of another design document.

For the Raw prototype, we will settle for being able to interface to some helper chips which can speak these dialects on our behalf.

Recently, there has been a proliferation of high speed signalling technologies like that chips SSTL, HSTL, GTL, LVTTL, and PCI. For our chip, we have been looking at SSTL and HSTL as potential candidates. We are currently leaning towards HSTL because it requires fewer external terminating resistors.

We expect to use the Xilinx Vertex parts to convert from our high-speed protocol of choice to other signaling technologies. These parts have the exciting ability to configurably communicate with almost all of the major signaling technologies. Although, in our prototype, these chips are external, I think that it is likely configurable I/O cells will find their way into the new extro-verted processors. This is because it will be so crucial for these processors to be able to communicate with all shapes and forms of devices. It may also be the case that extroverted processors will have bit-wise configurable FPGA logic near the I/O pins, for gluing together hardware protocols. After all, isn't glue logic what FPGAs were invented for? Perhaps our original conception of having fine-grained configurable logic on the chip wasn't so wrong; we just had it in the wrong place.

### 7.0.1  Raw I/O Model

I/O is a first-class software-exposed architectural entity on Raw. The pins of the Raw processor are an extension of both the mesh static and dynamic networks. For instance, when the west-most tiles on a Raw chip route a dynamic or static message to the west, the data values appear on the corresponding pins. Likewise, when an external device asserts the pins, they appear on-chip as messages on the static or dynamic network.

For the Raw prototype, the protocol spoken over the pins is the same static and dynamic handshaking network protocols spoken between tiles. If we actually had the FPGA glue logic on chip, the pins would support arbitrary handshaking protocols, including ones which require the pins to be bidirectional. Of course, for super-high speed I/O connections, there could be a fast-path straight to the pins.

The diagram "Logical View of a Raw Chip" illustrates the pin methodology. The striped lines represent the static and dynamic network pipelined buses. Some of them extend off the edge of the package, onto the pins. The number of static and dynamic network buses that are exposed off-chip is a function of the number of I/O pins that makes sense for the chip. There may only be one link, for ultra-cheap packages, or there may be total connectivity in a multi-chip module. In some cases, the number of static or dynamic buses that are exposed could be different. Or there may be a multiplex bit, which specifies whether the particular word transferred that cycle is a dynamic or static word. The compiler, given the pin image of the chip, schedules the dynamic and static communication on the chip such that it maximizes the utilization of the ports that exist on the particular Raw chip. I/O sends to non-existent ports will disappear.

The central idea is that the architecture facilitates I/O flexibility and scalability. The I/O capabilities can be scaled up or down according to the application. The I/O interface is a first-class citizen. It is not shoehorned through the memory hierarchy, and it provides an inter-

Package        Pins



▨▨▨   Static and/or Dynamic Network

**Logical View of a Raw Chip**

face which gives the compiler the access to the full bandwidth of the pins.

Originally, only the static network was exposed to the pins. The reasoning was that the static network would provide the highest bandwidth interface into the Raw tiles. Later, however, we realized that, just as the internal networks require support for both static and dynamic events, so too do the external networks. Cache line fills, external interrupts, and asynchronous devices are dynamic, and cannot be efficiently scheduled over the static network. On the other hand, the static network is the most effective method for processing a high bandwidth stream coming in at a steady rate from an outside source.

### 7.0.2 The location of the I/O ports (Perimeter versus Area I/O)

Area I/O is becoming increasingly common in today's fabrication facilities. In fact, in order to attain the pincounts that we desire on the SA-27E process, we have to use area I/O. This creates a bit of a problem, because all of our I/O connections are focused around the outside of the chip. IBM's technology allows us to simulate a peripheral I/O chip with area I/O. However, this may not be an option in the future. In that event, it is

possible to change the I/O model to match. In the Area I/O model, each switch and dynamic switch would have an extra port, which could potentially go in/out to the area I/O pads. This arrangement would create better locality between the source of the outgoing signal and the position of the actual pad on the die. Like in the peripheral case, these I/Os could be sparsely allocated.

### 7.0.3 Supporting Slow I/O Devices

In communicating with the outside world, we need to insure that we support low-speed devices in addition to the high-speed devices. For instance, it is unlikely that the off-the-shelf Virtex or DRAM parts will be able to clock as fast as the core logic of our chip. As a result, the SIB protocol needs to be re-examined to see if it still operates when connected to a client with a lesser clock speed. Ideally, the SIB protocol will support a software-settable clock speed divider feature, not unlike found on DRAM controllers for PCs. It is not enough merely to program the tiles so they do not send data words off the side of the chip too frequently; the control signals will still be switching too quickly.

The recent crop of .18 micron Xilinx parts may be sufficiently fast that we can hand-tune them to communicate with the Raw chip at a fast rate and then buffer up the words so that they can processed at a reasonable rate. An effort is currently under way to investigate this issue.

[MBT fixme: distinguish between I/O devices with low data-rates and parts which simply cannot run at the same speed as the raw chip.]

### 7.1 THE MEMORY SYSTEM

The Raw memory system has recently undergone a greater number of changes. An appendix to this document will specify the memory system in greater detail. This section will present an overview of the memory system.

A number of group members are actively researching this topic. One of the Raw group's goals is see what architectural features can actually be implemented efficiently using software and compiler technology. An effective all-software memory system would represent a significant accomplishment in the computer architecture field.

Unfortunately, the research on the memory system is not mature enough that we can omit hardware support for memory accesses and have confidence that our chip will remain useful for other research purposes.

As a result, the chip will support both hardware and software caching modes for the data memory. The instruction memory will be purely software cached.

In the software-caching mode, the SRAM will be accessed via load and store instructions. Accesses to addresses beyond the 32k size of the memory will return invalid results. It will be up to the software system to instrument the binaries for virtualization. The software caching system will use the memory network to access external DRAM.

In the hardware-caching mode, the SRAM will be used as a two way set-associative cache. We chose two way because it has low implementation cost (the tags have approximately the same area impact as a 32-bit multiplier when placed and routed), yet gives the effectiveness of a cache which is almost twice the size. These accesses will have three cycles latency; one for tag check, for memory access, and one for address generation. Misses will be handled by freezing the processor and handling the access out-of-band on the memory dynamic network. This method was chosen because it offers excellent performance and does not require extensive modification to the existing pipeline.

### 7.1.1  The Path to Copious Memory

We also need to consider the miss case. We need to have a way to reach the DRAMs residing outside of the Raw chip.This path is not as crucial as the Tag Check; however it still needs to be fairly efficient.

For this purpose, we will use a memory dynamic network to access the off-chip DRAMs. More details are given in the appendix sections.

### 7.2 SUMMARY

The strength of Raw's I/O architecture comes from the degree and simplicity with which the pins are exposed to the user as a first class resource. Just as the Raw tile expose the parallelism of the underlying silicon to the user, the Raw I/O architecture exposes the parallelism and bandwidth of the pins. It complements the key Raw goal -- to provide a simple interface to as much of the raw hardware resources to the user as possible.

# 8 DEADLOCK

In my opinion, the deadlock issues of the dynamic network is probably the single most complicated part of the Raw architecture. Finding a deadlock solution is actually not all that difficult. However, the lack of knowledge of the possible protocols we might use, and the constant pressure to use as little hardware support as possible makes this quite a challenge.

In this section, I describe some conditions which cause deadlock on Raw. I then describe some approaches that can be used to attack the deadlock problem. Finally, I present Raw's deadlock strategy.

## 8.0  DEADLOCK CONDITIONS

For the static network, it is the compiler's responsibility to ensure that the network is scheduled in a way that doesn't jam. It can do this because all of the interactions between messages on the network have been specified in the static switch instruction stream. These interactions are timing independent.

The dynamic network, however, is ripe with potential deadlock. Because we use dimension-ordered wormhole routing, deadlocks do not actually occur inside the network. Instead, they occur at the network interface to the tile. These deadlocks would not occur if the network had unlimited capacity. In every case, one of the tiles, call it tile A, has a dynamic message waiting at its input queue that is not being serviced. This message is flow controlling the network, and messages are getting backed up to a point where a second tile, B, is blocked trying to write into the dynamic network. The deadlock occurs when tile A is dependent on B's forward progress in order to get to the stage where it reads the incoming message and unblocks the network.

Below is an enumeration of the various deadlock conditions that can happen. Most of them can be extended to multiple party deadlocks. See the figure entitled "Deadlock Scenarios."

### 8.0.1  Dynamic - Dynamic

Tile A is blocked trying to send a dynamic message to Tile B. It was going to then read the message arriving from B. Tile B is blocked trying to send to Tile A. It was going to then receive from A. This forms a dependency cycle. A is waiting for B and B is waiting for A.



**Deadlock Scenarios**

### 8.0.2 Dynamic - Static

Tile A is blocked on $csto because it wants to statically communicate with processor B. It has a dynamic message waiting from B. B is blocked because it is trying to finish the message going out to A.

### 8.0.3 Static - Dynamic

Tile A is waiting on $csti because it is waiting for a static message from B. It has a dynamic message waiting from B.

Tile B is waiting because it is trying to send to tile C which is blocked by the message it sent to A. It was then going to write to processor A over the static network.

### 8.0.4 Static - Static

Processor A is waiting for a message from Processor B on $csti. It was then going to send a message.
Processor B is waiting for a message from Processor B on $csti. It was then going to send a message.
This is a compiler error on Raw.

### 8.0.5 Unrelated Dynamic-Dynamic

In this case, tile B is performing a request, and getting a long reply from D. C is performing a request, and getting a long message from A. What is interesting is that if only one or the other request was happening, there may not have been deadlock.

### 8.0.6 Deadlock Conditions - Conclusions

An accidental deadlock can exist only if at least one tile has a waiting dynamic network in-message and is blocked on either the $cgno, $csti, or $csto. Actually, technically, the tile could be polling either of those three ports. So we should rephrase that: the tile can only be deadlocked if there is a waiting dynamic message coming in and one of {$cgno is not empty, $csti does not have data available, or $csto is full}.

In all of these cases, the deadlock could be alleviated if the tile would read the dynamic message off of its input port. However, there may be some very good reasons for why the tile does not want to do this.

### 8.1 POSSIBLE DEADLOCK SOLUTIONS

The key two deadlock solutions are deadlock avoidance and deadlock recovery. These will be discussed in the next two sections.

### 8.2 DEADLOCK AVOIDANCE

Deadlock avoidance requires that the user restrict their use of the dynamic network to a certain pattern which has been proven to never deadlock.

The Deadlock avoidance disciplines that we generally arrive at are centered around two principles:

### 8.2.1 Ensuring that messages at the tail of all dependence chain are always sinkable.

In this discipline, we guarantee that the tile with the waiting dynamic message is always able to "sink" the waiting message. This means that the tile is always able to pull the waiting words off the network and break any cycles that have formed. The processor is not allowed to block while there are data words waiting.

These disciplines typically rely on an interrupt handler being fired to receive messages, which provides a high-priority receive mechanism that will interrupt the processor if it is blocked.

Alternatively, we could require that polling code be placed around every send.

Two examples disciplines which use that "always sinkable" principal are "Send Only" and "Remote Queues."

**Send Only**

For send-only protocols; like protocols which only store values, the interrupt handler can just run and process the request. This is an extremely limited model.

**Remote Queues**

For request-reply protocols, Remote Queues [Chong95], relies on an interrupt handler to dequeue arriving messages as they arrive. This handler will never send messages.

If this request was for the user process, the interrupt handler will place the message in memory, and set a flag which tells the user process that data is available. The user process then accesses the queue.

Alternatively, if the request is to be processed independently of the user process, the interrupt handler can drop down to a lower priority level, and issue a reply. While it does this will remain ready to pop up the higher priority level and receive any incoming messages.

Both of these methods have some serious disadvantages. First of all, the model is more complicated and adds software overhead. The user process must synchronize with the interrupt handler, but at the same time, make sure that it does not disable interrupts at an inopportune time. Additionally, we have lost that simple and fast pipeline-coupled interface that the network ports originally provided us with.

The Remote Queue method assumes infinite local memories, unless an additional discipline restricting the number of outstanding messages is imposed. Unfortunately, for all-to-all communication, each tile will have to reserve enough memory to handle the worst case -- all tiles sending to the same tile. This memory overhead can take up a significant portion of the on-tile SRAM.

## 8.2.2 Limit the amount and directions of data injected into the network.

The idea here is that we make sure that we never block trying to write to our output queue, making us available to read our input queue. Unless there is a huge amount of buffering in the network, this usually requires that we know a priori that there is some limit on the number of tiles that can send to us (and require replies) at any point, and that there is a limited on the amount of data in those messages. Despite this heavy restriction, this is nonetheless a useful discipline.

### The Matt Frank method

One discipline which we developed uses the effects of both principles. I called it the Matt Frank method. (It might also be called the client-server method, or the two party protocol.) In this example, there are two disjoint classes of nodes, the clients and the servers, which are connected by separate "request" and "reply" networks. The clients send a message to the servers on the request network, and then the servers send a message back on the reply network. Furthermore, each client is only allowed to have one outstanding message, which will fit entirely in its commit buffer. This guarantees that it will never be blocked sending.

Since clients and servers are disjoint, we know that when a client issues a message, it will not receive any other messages except for its response, which it will be

waiting to dequeue. Thus, the client nodes could never be responsible for jamming up the network.

The server nodes are receiving requests and sending replies. Because of this, they are not exempt from deadlock in quite the same way as the client nodes. However, we know that the outgoing messages are going to clients which will always consume their messages. The only possibility is that the responses get jammed up on their way back through the network by the requests. This is exactly what happened in the fifth dead-lock example given in the diagram "Deadlock Scenarios." However, in this case, the request and reply networks are separate, so we know that they cannot interact in this way. Thus, the Matt Frank method is deadlock free.

One simple way to build separate request-reply networks on a single dimension-ordered wormhole routed dynamic network is to have all of the server nodes on a separate side of the chip; say, the south half of the chip. With X-first dimension-ordered routing, all of the requests will use the W-E links on the top half of the chip, and then the S links on the way down to the server nodes. The replies will use the W-E links on the bottom half of the chip, and the N links back up to the clients. We have effectively created a disjoint partition of the network links between the requests and the replies.

For the Matt Frank protocol, we could lift the restriction of only one outstanding message per client if we guaranteed that we would always service all replies immediately. In particular, the client cannot block while writing a request into the network. This could be achievable via an interrupt, polling, or a dedicated piece of hardware.

### 8.2.3 Deadlock Avoidance - Summary

Deadlock avoidance is an appealing solution to handling the dynamic network deadlock issue. However, each avoidance strategy comes with a cost. Some strategies reduce the functionality of the dynamic network, by restricting the types of protocols that can be used. Others require the reservation of large amounts of storage, or cause a low utilization of the underlying network resources. Finally, deadlock avoidance can complicate and slow down the user's interface to the network. Care must be made to weigh these costs against the area and implementation cost of more brute-force hardware solutions.

## 8.3 DEADLOCK RECOVERY

An alternative approach to deadlock avoidance is deadlock recovery. In deadlock recovery, we do not restrict the way that the user employs the network ports. Instead, we have a recovery mode that rescues the program from deadlock, should one arise. This recovery mode does not have to be particularly fast, since deadlocks are not expected to be the common case. As with a program with pathological cache behaviour, a program that deadlocks frequently may need to be rewritten for performance reasons.

Before I continue, I will introduce some terminologies. These are useful in evaluating the ramifications of the various algorithms on the Raw architecture.

**Spontaneous Synchronization** is the ability of a group of Raw chips to suddenly (not scheduled by compiler) stop their current individual computations and work together. Normally, a Raw tile could broadcast a message on the dynamic network in order to synchronize everybody. However, we obviously cannot use the dynamic network if it is deadlocked. We cannot use the static network to perform this synchronization, because the tiles would have to spontaneously synchronize themselves (and clear out any existing data) in order to communicate over that network!

We could have a interrupting timer which is synchronized across all of the Raw tiles to interrupt all of the tiles simultaneously, and have them clear out the static network for communication. If we could guarantee that they would all interrupt simultaneously, then we could clear out the static network for more general communication. Unfortunately, this would mean that the interrupt timer would have to be a non maskable interrupt, which seems dangerous.

In the end, it may be that the least expensive way to achieve *spontaneous synchronization* is to have some sort of non-deadlocking synchronization network which does it for us. It could be a small as one bit. For instance, the MIT-Fugu machine had such a one bit rudimentary network [Mackenzie98].

**Non-destructive observability** requires that a tile be able to inspect the contents of the dynamic network without obstructing the computation. This mechanism could be implemented by adding some extra hardware to inspect the SIBs. Or, we could drain the dynamic network, store the data locally on the destination nodes, and have a way of virtualizing the $cgni port.

### 8.3.1  Deadlock Detection

In order to recover from deadlock, we first need to detect deadlock. In order to determine if a deadlock truly exists, we would need to analyze the status of each tile, and the network connecting them, looking for a cyclic dependency.

One deadlock detection algorithm follows:

The user would not be allowed to poll the network ports, otherwise, the detection algorithm would have no way of knowing of the program's intent to access the ports. The detection algorithm runs as follows: The tiles would synchronize up, and run a statically scheduled program (that uses the static network) which analyzes the traffic inside the dynamic network, and determines whether the each tile was stalled on a instruction accessing $csto, $csti, or $cgno. It can construct a dependency graph and determine if there is a cycle.

However, the above algorithm requires both *spontaneous synchronization* and *non-destructive observability*. Furthermore, it is extremely heavy-weight, and could not be run very often.

### 8.3.2  Deadlock Detection Approximation

In practice, a deadlock detection approximation is often sufficient. Such an approximation will never return a false negative, and ideally will not return too many false positives. The watchdog timer, used by the MIT-Alewife machine [Kubiatowicz98] for deadlock detection is one such approximation.

The operation is simple: each tile has a timer that counts up every cycle. Each cycle, if $cgni is empty, or if a successful read from $cgni is performed, then the counter is reset. If the counter hits a predefined user-specified value, then a interrupt is fired, indicating a potential deadlock.

This method requires neither *spontaneous synchronization* nor *non-destructive observability*. It also is very lightweight.

It remains to be seen what the cost of false positives is. In particular, I am concerned about the case where one tile, the aggressive producer, is sending a continuous stream of data to a tile which is consuming at a very slow rate. This is not truly a deadlock. The consumer will be falsely interrupted, and will run even slower because it will be the tile who will be running the deadlock recovery code. (Ideally, the producer would have been the one running the deadlock code.) Fugu

[MacKenzie98] dealt with these sorts of problems in more detail. At this point in time, we stop by saying that the user or compiler may have to tweak the deadlock watchdog timer value if they run into problems like this. Alternatively, if we had the spontaneous synchronization and non-destructive observability properties, we could use the expensive deadlock detection algorithm to verify if there was a true deadlock. If it was a false positive, we could bump up the counter.

### 8.3.3  Deadlock recovery

Once we have identified a deadlock, we need to recover from the deadlock. This usually involves draining the blockage from the network and storing it in memory. When the program is resumed, a mechanism is put in place so that when the user reads from the network port, he actually gets the values stored in memory.

To do this, we have a bit that is set which indicates that we are in this "dynamic refill" mode. A read from $cgni will return the value stored in the special purpose register, "DYNAMIC_REFILL." It will also cause an interrupt on the next instruction, so that a handler can transparently put a new value into the SPR. When all of the values have been read out of the memory, the mode is disabled and operation returns to normal.

An important issue is where the dynamic refill values are stored in memory. When a tile's watchdog counter goes off, it can store some of the words locally. However, it may not be expedient to allocate significant amounts of buffer space for what is a reasonably rare occurrence. Additionally, since the on-chip storage is extremely finite, in severe situations, we eventually will need to get out to a more formidable backing store. We would need spontaneous synchronization to take over the static network and attain the cooperation of other tiles, or a non-deadlocking backup network to perform this. [Mackenzie98]

### 8.3.4  More deadlock recovery problems

Most of the deadlock problems describe here have been encountered by the Alewife machine, which used a dynamic network for its memory system. However, those machines have the fortunate property that they can put large quantities of RAM next to each node. This RAM can be accessed without using the dynamic network. On Raw, we have a very tiny amount of RAM that can be accessed without travelling through the network. Unless we can access a large bank of memory deadlock-free, the deadlock avoidance and detection code must take up precious instruction SRAM space on the tile.

Ironically, a hardware deadlock avoidance mechanism may have a lesser area cost than the equivalent software ones.

### 8.3.5  Deadlock Recovery - Summary

Deadlock recovery is also an appealing solution to handling the deadlock problem. It allows the user unrestricted use of the network. However, it requires the existence of a non-deadlockable path to memory. This can be attained by using the static network and adding the ability to spontaneously synchronize. It can also be realized by adding another non-deadlocked network.

### 8.4 DEADLOCK ANALYSIS

The issue of deadlock in the dynamic network is of serious concern. Our previous solutions (like the NEWS single bit interrupt network) have had serious disadvantages in terms of complexity, and the size of the resident code on every SRAM. For brevity, I have opted not to list them here.

In this section, I propose a new solution, which I believe offers extremely simple hardware, leverages our existing dynamic network code base, and solves the deadlock problem very solidly. It creates an abstraction which can be used to solve a variety of other outstanding issues with the Raw design. Since this is preliminary, the features described here are not described in the "User's View of Raw" section of the document.

First, let us re-examine the dynamic network manifesto:

```
The primary intention of the
dynamic network is to support memory
accesses that cannot be statically
analyzed. The dynamic network was also
intended to support other dynamic
activities, like interrupts, dynamic
I/O accesses, speculation, synchroni-
zation, and context switches.
Finally, the dynamic network was the
catch-all safety net for any dynamic
events that we may have missed out on.
```

Even now, the Raw group is very excited about utilizing deadlock avoidance for the dynamic network. We argue that we were not going to be supporting general-purpose user messaging on the Raw chip, so we could require the compiler writers and runtime system programmers to use a discipline when they use the network.

The problem is, the dynamic network is really the extension mechanism of the processor. Its strength is in its ability to support protocols that we have left out of the hardware. We are using the dynamic network for many protocols, all of which have very different properties. Modifying each protocol to be deadlock-free is hard enough. The problem comes when we attempt to run people's systems together. We then have to prove that the power set of the protocols is deadlock free!

Some of the more flexible deadlock avoidance schemes allow near-arbitrary messaging to occur. Unfortunately, these schemes often result in decreased performance, or require large buffer space.

The deadlock recovery schemes provide us with the most protocol flexibility. However, they require a deadlock-free path to outside DRAM. If this is implemented on top of the static network, then we have to leave a large program in SRAM just in case of deadlock.

## 8.5 THE RAW DEADLOCK SOLUTION

Thinking about this, I realized that the dynamic network usage falls into two major groups: memory accesses and essentially random unknown protocols. These two groups of protocols have vastly different properties.

My solution is to have two logically disjoint dynamic networks. These networks could be implemented as two separate networks, or they could be implemented as two logical networks sharing the same physical wires. In the latter case, one of the networks would be deemed the memory network and would always have priority.

The memory network would implement the Matt Frank deadlock avoidance protocol. The off-chip memory accesses will easily fit inside this framework. In this case, the processors are the "clients" and the DRAMS, hanging off the south side of the chip, are the "servers." Interrupts will be disabled during outstanding accesses. Since the network is deadlock free, and guaranteed to make forward progress, this is not a problem. This also means that we can dangle messages into the network without worry, improving memory system performance. This network will enforce a round-robin priority scheme to make sure that no tile gets starved. This network can also be used for other purposes that involve communication with remote devices and meet the requirements. For instance, this mechanism can be used to notify the tiles of external interrupts. Since the network cannot deadlock, we know that we will have a relatively fast interrupt response time. (Interrupts would be implemented as

an extra bit in the message header, and would be dequeued immediately upon arrival. This guarantees that they will not violate the deadlock avoidance protocol.)

**Note: We are no longer using the Matt Frank deadlock avoidance protocol for Raw. Instead we use the "Buffer Metering" protocol, which is described in the "Memory Network Redux" Section. This protocol relaxes the directionality constraints on memory network usage.**

The more general user protocols will use the low-priority dynamic network, which would have a commit buffer [mbt: general network no longer has a commit buffer], and will have the $cgno/$cgni that we described previously. They will use a deadlock recovery algorithm, with a watchdog deadlock detection timer. Should they deadlock, they can use the memory network to access off-chip DRAM. In fact, they can store all of the deadlock code in the DRAM, rather than in expensive SRAM. Incidentally, the DRAMs can be used to implement *spontaneous synchronization.*

One of the nice properties that comes with having the separate deadlock-avoidance network is that user codes do not have to worry about having a cache miss in the middle of sending a message. This would otherwise require loading and unloading the message queue. Additionally, since interrupt notifications come on the memory network, the user will not have to process them when they appear on the input queue.

## 8.6 THE HIGH-PRIORITY DYNAMIC NETWORK

Since the low-priority dynamic network corresponds exactly to the dynamic network described in the previous dynamic network section, it does not merit further discussion.

The use of the high-priority network needs some elaboration, especially with respect to the deadlock avoidance protocol.

The diagram "High-Priority Memory Network Protocol" helps illustrate. This picture shows a Raw chip with many tiles, connected to a number of devices (DRAM, Firewire, etc.) The protocol here uses only one logical dynamic network, but partitions it into two disjoint networks. To avoid deadlock, we restrict the selection of external devices that a given tile can communicate with. For complete connectivity, we could implement another logical network. The rule for connectivity is:

**Memory Network Protocol**

Each tile is not allowed to communicate with a device which is NORTH or WEST of it. This guarantees that all requests travel on the SOUTH and EAST links, and all replies travel on the NORTH and WEST links.

Although this is restrictive, it retains four nice properties. First, it provides high bandwidth in the common case, where the tile is merely communicating with its partner DRAM. The tile's partner DRAM is a DRAM that has been paired with the tile to allocate the network and DRAM bandwidth as effectively as possible. Most of the tile's data and instructions are placed on the tile's partner DRAM.

The second property, the *memory maintainer* property, is that the northwest tile can access all of the DRAMs. This will be extremely useful because the non-parallelizeable operating system code can run on that tile and operate on all of the other tile's memory spaces. Note that with strictly dimensioned-ordered routing, the memory maintainer cannot actually access all of the devices on the right side of the chip. This problem will be discussed in the "I/O Addressing" section.

The third property, the *memory dropbox* property, is that the southeast DRAM is accessible by all of the tiles. This means that non performance-critical synchroniza-

tion and communication can be done through a common memory space. (We would not want to do this in performance critical regions of the program, because of the limited bandwidth to a single network port.)

These last two properties are not fundamental to the operation of a Raw processor; however they make writing setup and synchronization code a lot easier.

Finally, the fourth nice property is that the system scales down. Since all of the tiles can access the southeast-most DRAMs, we can build a single DRAM system by placing the DRAM on the southeast tile.

We also can conveniently place the interrupt notification on one of the southeast links. This black box will send a message to a tile informing it that an interrupt has occurred. The tile can then communicate with the device, possibly but not necessarily in a memory-mapped fashion. Additionally, DMA ports can be created. A device would be hooked up to these ports, and would stream data through the dynamic network into the DRAMs, and vice versa. Logically, the DMA port is just like a client tile. I do not expect that we will be implementing this feature in the prototype.

Finally, this configuration does not require that the devices have their own dynamic switches. They will

merely inject their messages onto the pins, with the correct headers, and the routes will happen appropriately. This means that the edges of the network are not strictly wormhole routed. However, in terms of the wormhole routing, these I/O pins look more like another connection to the processor than an actually link to the network. Furthermore, the logical network remains partitioned because requests are on the outbound links and the replies are inbound.

## 8.7 PROBLEMS WITH I/O ADDRESSING

One of the issues with adding I/O devices to the periphery of the dynamic network is the issue of addressing. When the user sends a message, they first inject the destination tile number (the "absolute address"), which is converted into a relative X and Y distance. When we add I/O devices to the periphery, we suddenly need to include them in the absolute name space.

However, with the addition of the I/O nodes, the X and Y dimensions of the network are no longer powers of two. This means that it will be costly to convert from an absolute address to a relative X and Y distance when the message is sent.

Additionally, if we place devices on the left or top of the chip, the absolute addresses of the tiles will no longer start at 0. If we place devices on the left or right, the tile numbers will no longer be consecutive. For programs whose tiles use the dynamic network to communicate, this makes mapping a hash key to a tile costly.

Finally, I/O addressing has a problem because of dimension ordered routing. Because dimension ordered routing routes X, then Y, devices on the left and the right of the chip can only be accessed by tiles that are on the same row, unless there is an extra row of network that links all of the devices together.

## 8.8 THE FUNNY BITS

All of these problems could be solved by only placing devices on the bottom of the chip.

However, the "funny bits" solution which I propose allows us full flexibility in the placement of I/O devices, and gives us a unique name space.

The "funny bit" concept is simple. An absolute address still has a tile number. However, the three highest order bits of the address, previously unused, are reserved for the funny bits. These bits are preserved upon translation of the absolute address to relative address. These funny bits specify a final route that should be done after all dimensioned ordered routing has occurred. These funny bits can only be used to route off the side of the chip. It is a programmer error to use the funny bits in a send to a tile.

With this mechanism, the I/O devices no longer need to be mapped into the absolute address space. To route to an I/O device, one merely specifies the address of the tile that the I/O device is attached to, and sets the bit corresponding to the direction that the device is located at relative to the tile.

The funny bits mechanism is deadlock free because once again, it acts more like another processor attached to the dynamic network than a link on the network. A more rigorous proof will follow in subsequent theses.

An alternative to the funny bits solution is to provide the user with the ability to send messages with relative addresses, and to add extra network columns to the edge of the tile. This solution was used by the Alewife project [Kubiatowicz98]. Although the first half of this alternative seemed palatable, the idea of adding extra hardware (and violating the replicated uniform nature of the raw chip) was not.

## 8.9 SUMMARY

In this section, I discussed a number of ways in which the Raw chip could deadlock. I introduced two solutions, deadlock avoidance and deadlock recovery, which can be used to solve this problem.

I continued by re-examining the requirements of the dynamic network for Raw. I showed that a pair of logical dynamic networks was an elegant solution for Raw's dynamic needs.

The high-priority network uses a deadlock-avoidance scheme that I labelled the "Matt Frank protocol." Any users of this network must obey this protocol to ensure deadlock-free behaviour. This network is used for memory, interrupt, I/O, DMA and other communications that go off-chip.

The high-priority network is particularly elegant for memory accesses because, with minimal resources, it provides four properties: First, the memory system scales down. Second, the high-priority network supports *partner memories*, which means that each tile is assigned to a particular DRAM. By doing the assignments intelligently, the compiler can divide the band-

width of the high-priority network evenly among the tiles. Third, this system allows the existence of a ***memory dropbox***, a DRAM which all of the tiles can access directly. Lastly, it allows the existence of a ***memory maintainer***; which means at least one tile can access all of the memories.

The low-priority network uses deadlock recovery and has maximum protocol flexibility and places few restrictions on the user. The deadlock recovery mechanism makes use of the high-priority network to gain access to copious amounts of memory (external DRAM). This memory can be used to store both the instructions and the data of the deadlock recovery mechanism, so that precious on-chip SRAM does not need to be reserved for rare deadlock events.

This deadlock solution is effective because it prevents deadlock and provides good performance with little implementation cost. Additionally, it provides an abstraction layer on the usage of the dynamic network that allows us to ignore the interactions of the various clients of the dynamic network.

Finally, I introduced the concept of "funny bits" which provides us with some advantages in tile addressing. It also allows all of the tiles to access the I/O devices without adding extra network columns.

With an effective solution to the deadlock problem, we can breath easier.

# 9 Implementation of the DYNAMIC NET-WORKS

## 9.0  INTRODUCTION

This chapter serves as a specification for the dynamic networks. To read more about the motivation of the current dynamic networks, please refer to "The Design and Implementation of a Raw Architecture Workstation."

## 9.1 A TALE OF TWO DYNAMIC NETWORKS

The Raw chip has two logical dynamic networks (DNs), connecting the raw tiles in a two-dimensional point-to-point mesh topology.

The **general dynamic network (GDN)** is supported by a deadlock recovery mechanism which allows us to employ a diverse collection of protocols without having to prove deadlock properties about them.

The **memory dynamic network (MDN)** is restricted to protocols which follow a specific deadlock avoidance pattern, the Matt Frank protocol. These protocols all communicate with the periphery of the chip. Interrupt, DRAM and device messages consitute the primary traffic for this network.

## 9.2 THE DYNAMIC NETWORK HARDWARE

In this section, we will discuss the various components of the dynamic network. We will discuss the wires, the routers, and the processor interface to the network.

### 9.2.1  The Dynamic Network Hardware

The two DNs have separate routers and input buffers, and separate physical data wires. They are essentially two copies of the exact same hardware.

Since the DNs use the SIB protocol, the inter-tile wiring of the static networks and the DNs is identical. The diagram "The Dynamic Network Wires" shows the signals that run between dynamic routers on separate tiles.



**The Dynamic Network Wires**

### 9.2.1.1 Alternative: Multiplexing the Wires

Since we do not anticipate heavy traffic on both networks simultaneously, we could multiplex the physical data wires that connect the two logical networks. The multiplexing of the data wires would be performed transparently without user intervention. Although, the high-priority network is labelled "high-priority," the physical wires would be round-robin shared. Thus, from the perspective of the hardware, the networks have equal priority.

The round-robin sharing would work as follows:
If neither network has a data word to send, nothing occurs. If one network has a data word to send, it sends it. If both networks have a data word to send, the network which sent the last data word will stall a cycle and allow the other network to procede.

Although we suspect that this configuration will offer good performance even in the event of heavy load, it remains to be proven.

However, multiplexing the dynamic networks adds an extra mux on the tile-to-tile path. This would make the tile-to-tile delay greater on the dynamic network than on the static network, a property we would like to avoid. For the time being we will assume that the dynamic networks are not multiplexed.

### 9.2.2  The Dynamic Network Router

The dynamic network is a dimension-ordered, wormhole routed flow-controlled network [Dally86]. The router routes in the X direction, then in the Y direc-

The Dynamic Network Router

tion. We implemented the protocol on top of the SIB protocol that was used for the static network. The figure entitled "The Dynamic Network Router" illustrates this. The dynamic network router is identical to the static network router, except it has a dynamic scheduler instead of the actual switch processor.

### 9.2.2.1 Dynamic Network Message Format

Each dynamic network message that travels through the network has a header, followed by the message data.

This header is described in the figure "Dynamic Message Header." The header contains a route encoded by an absolute X position and an absolute Y position, a length and F, the funny bits. The middle fourteen bits (non-bold font) are ignored by the dynamic routers. The absX and absY positions encode the position on the Raw mesh that the message is destined for. Each tile has a SPR which indicates the tile's x and y position, which



Dynamic Message Header

is used to determine if a message has arrived or not. The length field details the number of words in the message, not counting the header. The final route (F) field encodes specifies the direction of the final route to be performed after all X-Y routing is performed. The header restricts the maximum size of a raw array to a 32x32 array, and restricts message transfer sizes to 31 words, not including the header. The header word is dropped upon delivery to the destination tile, so the recipient of the message never sees it. However, when the message is routed off the sides of the pins, the header is retained, and is visible to outside devices. Thus the settings of the "usr" and "origY/origX" fields are completely irrelevant to a receiving tiles.

### 9.2.3  Dynamic Scheduler

The dynamic sheduler examines the headers of incoming message. If the X position is not equal, then it initializes a state machine which will transfer one word per cycle of the header and message out of the west or east port, based on the sign of the difference between the tile's position and the message's destination.

If the X position is equal, and the Y position is not equal, the message is sent out of the south or north ports. If both X and Y are equal, then the final route bits encode the final route to be performed. The final route

**Table 3: Final Route Field**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|------|-----|---|---|---|---|-----|-----|
| none | inv | W | S | E | N | inv | inv |

and top bits fields are never cleared as the message is passed along.

It is invalid to specify a final route which does not go off of the side of the chip. In fact, with a naive implementation, a final route which does not go off the side of the chip will actually end up be infinitely routed back and forth inside the network. We may want to modify the switches slightly so that messages that come from the same direction as we are to route to are squashed. This would involve supressing the "valid out" signal on a send, and updating the "spaceAvail" counter appropriately.

Because multiple input messages can contend for the same output port, there needs to be a priority scheme. We use a simple clock-wise round robin scheduler to select between contenders. This means that an aggressive producer of data will not be able to block

44

other users of the network from getting their messages through.

Because the scheduler must parse the message header, and then modify it to forward it along, the original design took two cycles for the header to pass through the network. Each word after that would only take one cycle. We have arrived at slightly different solution which allows through-routes to occur in a single cycle latency, but turns occur with two cycle latency. This design appears to have good critical path properties and gives us reasonable performance, especially, on large Raw mesh configurations.

### 9.2.4 Dynamic Network Interfaces

The two networks appear very similar from the programmer's perspective. However, the semantics of each network are slightly different. In this section, we describe the common interface, and then describe the particulars for each network.

#### 9.2.4.1 The Common Network Interface

The dynamic network ports appear to the user as $CGNI and $CGNO (for the general network) and $CMNI and $CMNO (for the memory network).

In either case, the user sends a message by writing a message header into the network, followed by the datawords. The message header was described in the section "Dynamic Network Message Format."

The user can use the "IHDR" instruction to construct a message header. (Typically, the target for the IHDR instruction will be the $cgno port.) Since there is nothing magical about the "IHDR" instruction, the user could also use choose to construct his/her own header, or even store commonly used headers in memory.

```
# low N bits of R3 = tile number

IHDR $cgno, 0x400($3)
```

The IHDR instruction converts the low N bits to a location on the mesh, filling in the appropriate X and Y destination. The top 12 bits of the imm field are blindly copied into the top 12 bits of the header. The origin of the message is inserted, mostly for debugging reasons. We do not anticipate that F or usr fields will be set to anything except zero. The IHDR instruction is used for applications which use low-order interleaving type communication, within the array of tiles.

```
IHDR result, R3, 0x0400
```

```
(R3[4:0]
  & GDN_XMASK[4:0]) + GDN_XADJ[4:0]
```

```
((R3[11:0] >> GDN_YSHIFT[2:0])
  & GDN_YMASK[4:0]) + GDN_YADJ[4:0]
```

IHDR Header Construction

The DN_XPOS, DN_YPOS, GDN_XMASK, GDN_YMASK, GDN_XADJ, GDN_YADJ, and GDN_YSHIFT values are located in a special purpose register. These provide the relevant data required to translate from an address to a position within the process's virtual configuration (say 4x4) to an absolute address on the mesh (which could be say 32x32.)

Alternatively, the OHDR instruction specifies a different sort of translation, (used by the hardware caching mechanism), which is intended for memory and device accesses external to the tile.

The DRAMs are mapped into a common 32bit address space, which is doled out according to the dimensionality of the Raw processor. A 4x4 Raw chip has 8 ports each of which own 512 MB of the address space, while a 32x32 chip has 64 ports, each of which own 64 MB of DRAM. It is possible to use some skanky tricks to increase the amount of addressible memory, but I think that a commercial Raw implementation would use a 64 bit dataword anyways, making the address space issue a moot point. The job of the OHDR instruction is to take an address and generate the correct header which will transport data to the port that that address has been assigned to.

45

The OHDR instruction uses R3 to fill in the F, X, and Y fields. The origX and origY fields are also updated with the DN_XPOS and DN_YPOS values. This constitutes a return tile address for the instruction. The bits 12..4 of the immediate are copied into bits 28..20 of the result. This immediate thus specifies the length of the message for the DN, as well as an optional request type. This somewhat convoluted mechanism allows the user to convey much of the standard information (destination, source, request type, and length) with a single instruction. This saves network latency and bandwidth. Note that the X dimension and Y dimension of the tile array do not have to be the same.

The MDN_EXTEND bit is a later-arriving feature which modifies the behaviour of the OHDR instruction, mapping addresses on all four sides of the chip, instead of just on two. This requires a modification in the deadlock avoidance protocols of the machine, but enables the tiles to use a greater proportion of the pin bandwidth for memory traffic.

```
OHDR $cmno, 0x1FF0($3)

side = R3[30] & MDN_EXTEND;
bits = MDN_EXTEND ?
         R3[29..25] : R3[30..26]

switch (R3[31])

0: ABS_X = side ? 0    : MDN_XMAX;
   FBITS = side ? west : east;
   ABS_Y = (bits >> MDN_YSHIFT)


1: ABS_Y = side ? 0    : MDN_YMAX;
   FBITS = side ? north : south;
   ABS_X = (bits >> MDN_XSHIFT)

  (3 >= MDN_xSHIFT >= 0)

        3 --> 4x4
        2 --> 8x8
        1 --> 16x16
        0 --> 32x32
```

OHDR Header Construction

Address Ranges Assigned to Ports,

(for 4x4, 8x8 MDN_EXTEND = 1, 8x8,
4x4 MDN_EXTEND = 1, 8x4, and 16x16 Raw)

46

### 9.2.4.2 Memory Network Interface

When a word is written into the MDN, the value goes directly out to the MDN switch for interpretation. This offers the greatest performance. However, it means that it is the programmers' responsibility to ensure the atomicity of the message send. For the most part, this means that interrupts should be masked during an outstanding transaction. It also means that users of the MDN must gaurantee that no cache miss occurs during MDN accesses. More detail is given on this subject in elsewhere in this document.

During a hardware cache miss, the pipeline is frozen, so interrupts are implicitly masked.

### 9.2.4.3 Memory Interrupts

External devices can signal interrupts to user tiles using the Memory network. The interrupt controller must be positioned on the SE-most tile so that it can communicate with all of the tiles, and maintain our communication discipline. In order to avoid deadlock, each tile's dynamic network hardware must be prepared to remove interrupt messages at all times. This will occur transparent to the user. However, the corresponding "External Interrupt" bit will be set in the appropriate user-accessible field.

This interrupt notification will be used by both the O/S and devices to communicate with tiles in a asynchronous, dead-lock free method.

The interrupt controller will be implemented on an FPGA and will interact with the devices and the tile in an appropriate manner. Only one outstanding interrupt per tile will be allowed. Once the tile branches to its handler, it will acknowledge the interrupt controller (which will give it an interrupt number or a vector to jump to) and further external interrupts may occur.

| 31 29 | 28    24 | 20 19 | 15 14 | 10 9 | 5 4 | 0 |
|---|---|---|---|---|---|---|
| **F** | **00000** | **1111** | | | **abs Y** | **abs X** |
| 3 | 5 | 4 | 5 | 5 | 5 | 5 |

Interrupting Message Header

An interrupt message will be defined as messages whose length field is set to zero and whose user field is set to 1111. Interrupt messages only make sense inside the tile. There is no fixed interpretation of interrupt messages outside of the chip.

### 9.2.4.4 General Network Interface

This section is deprecated. The GDN interface is identical to the MDN interface. Ignore the following two paragraphs:

"Because the general network can deadlock, a more rigorous mechanism must be put in place than for the MDN. We call this device a commit buffer. When the user writes to $cgno, the value is queued up in the commit buffer. When the last word is written (the message length is encoded by the header word), the commit buffer commences the streaming of data into the GDN. A write to $cgno during the time between which the last word is written to the commit buffer and the time at which the last word leaves the commit buffer will cause the processor to stall.

In the event of an interrupt or context switch, it may be necessary to drain uncommited elements from the commit buffer. There is a SPR which shows how many elements are left in the buffer, and another register which a read from will cause the youngest message word to be dequeued. The programmer can check the element count register, and then read the other values out, including the header word. "

### 9.2.4.5 General Network - Deadlock Recovery

The general network has a watchdog timer which is incremented every cycle that data is available on the general network but the network port is not read. If data is read, or nothing is available, the counter will be reset to zero. If this counter reaches the watchdog value, an interrupt will be caused. The interrupt handler will pull words out of the network and store them locally (which may possibly go out to DRAM). It will then set a flag, and the DRVAL register, which will enable virtualization of the network.

**Table 4: User Fields, MDN**

| Value | meaning |
|---|---|
| 0000 | Cache-line Read (Addr) |
| 0001 | DMA Read - North |
| 0010 | DMA Read - West |
| 0011 | |
| 0100 | Cache-line Write (Addr, 8 Words) |
| 0101 | Escape (next word is command) |

**Table 4: User Fields, MDN**

| Value | meaning |
|-------|---------|
| 0110 | User defined |
| 0111 | DMA Write (Address, Mask, 8 Words) |
| 1000 | DMA tagged Read Reply (multi session) |
| 1001 | User defined |
| 1010 | User defined |
| 1011 | User defined |
| 1100 | User defined |
| 1101 | MDN Relay |
| 1110 | System Monitor Service (external) Store Acknowledgement (internal to chip) (The Store_Ack message type looks at the sender field to determine if the "Partner" or "Non-partner" counter should be decre-mented.) |
| 1111 | reserved |

## 9.3 SUMMARY

The dynamic network design leveraged many of the same underlying hardware components as the static switch design. Its performance is not quite as good as the static network's because the route directions are not known a priori. A great deal more was said in the dead-lock section of this document.

# 10 Memory Network Redux

## 10.0 INTRODUCTION

The original deadlock avoidance discipline for the memory network created restrictions on the directionality of network requests and replies. Essentially, tiles could only send requests to the east and south. This guaranteed that requests and replies never conflicted, and thus that the I/O devices could always make forward progress. Unfortunately, this strategy restricts the flexibility of DRAM (only half of the I/O ports could be used for DRAM) and device placement on the chip and makes communication with the devices on the left and top side of the chip particularly cumbersome. Additionally, it means that tiles are inherently unequal and non-uniform; the upper-left hand tile has far superior powers than the lower right-hand tile. This complicates the programming module and reduces the Raw processor's applicability as a computing fabric.

Additionally, the original deadlock avoidance scheme relies on the flow-control within the mesh network to prevent deadlock. In other words, if one device is running behind, its input will back up into the network, causing unrelated I/O transactions to block as well. Conventional thought predicts that dynamic networks tend to break down as congestion approaches %30. It is preferable that the I/O port have enough buffering to sink unprocessed requests so that it will not have a more global effect on performance.

Finally, even in the case that all of the tiles are targetting the same I/O port, and they would need to wait anyways, the local round-robin policy of the routers does not translate into a globally fair policy. This is the so called "Parking Lot" pathology. Imagine four tiles on a row all repeatedly transmitting to a DRAM on the east side of the chip on that same row. The tile closest to the DRAM will get %50 of the bandwidth (because the bandwidth will be equally shared between a tile's compute processor and all of the tiles west of it) , the next tile will get %25, and the last two %12.5. On the other hand, if the port can sink the messages and keep the network clear, the tiles will get a near-equal alotment of the bandwidth.

## 10.1 Buffer Metering

The "buffering metering" protocol is summarized with the following rules:

1. Each network node is not allowed to block on a memory network send unless it can guarantee that its input buffers can sink all of memory network messages that it may be receiving.

The data caches trivially obey this protocol; $cmno has 16 elements of buffering. The worst possible case occurs when a read request (2 words) followed by an evict message (10 words) is sent. If the last word of the read request has escaped the buffer, it is gauranteed that the buffer has 16 elements of buffering to receive the 12 elements of the evict message. So the processor will not be blocked when the read request response arrives, and thus it will always be able to sink all of the memory network words it is receiving.

The I/O nodes have a more complex analysis. Each I/O node can receive requests from every tile. Because it pulls values off of the network, and then places responses back onto the network, it is possible that it may block on output due to other traffic in the network. In that case, it needs to gaurantee that it can continue to pull off any words that are coming on the network. If no disciplines were imposed, the tiles would have the potential to stream an infinite amount of data into the I/O node and cause it to deadlock.

The solution that we employ allocates enough RAM at the I/O node to sink a finite number (>=1) of messages from each tile. Nominally, each tile is allocated 12 words of space on the input buffer. This amounts to 12 N words of buffer space for every I/O port. Additionally, each tile is allocated an additional quantity of buffers on one I/O ports that is typically (but doesn't strictly have to be) directly N,E,W or S of it. This allows for high bandwidth transfers between the tile and its "**partner**" port. Finally, there is a pool of free buffers, maintained by the OS, which can be allocated to I/O devices or software running on the tiles.

The existance of **partner** ports derives from two important factors: first, a desire to partition memory requests from tiles in order to load balance the ports. The stacks and other local data structures of the tiles can be assigned by the runtime system to the partner DRAM of a tile. Secondly, the FBITs protocol degenerates to some degree on I/O port that are not in a direct line with the originating tiles -- because of dimension-ordered routes, they end up flattening (see figure "Fbits Related Flattening") to the links on the edges of the chips for either the request or the reply side of the transaction. Note that the two straight routes do not observe any flattening. Alternative, non-dimension-order routing protocols could also be used to alleviate the flattening problem.

We supplement the cache with a counter (the "store counter") that counts the number of outstanding stores. When a store request is issued, it is decremented. When a store request reply is received, it is incremented. Should the counter value be zero at the time of a cache



0xC* 0xD* 0xE* 0xF*

0x4*          0x0*
0x5*          0x1*
0x6*          0x2*
0x7*          0x3*

0x8* 0x9* 0xA* 0xB*

Example of a good full-bandwith Raw partner  memory allocation.



FBits Related Flattening

transaction, the cache state machine will stall until it receives the message. The store counter is actually implemented as two counters that disinguish between the partner port, and the non-partner port. Because of timing issues, the cache state machine actually will stall if either of the counters is zero.

The software can also access the I/O ports and must maintain the discipline as well. The software can query the operating system for additional buffers on a particular I/O port, and ensure that it does not excede the maximum number of words in flight at any time.

Alternatively it can carefully interact with the store counter and borrow accesses from the hardware (more complex). It issues this by issuing an mtsri/mtsr instruction to STORE_ACK with the DECREMENT_MODE set. It also needs to check the appropriate store counter and make sure that it is non-zero.  Alternatively, the software can maintain the counter separately, use non-acking stores from the DRAMS, and use pings to ensure that data has arrived.

The buffer metering scheme is somewhat of an experiment. It is likely that for larger or commercial Raw arrays, separate request-reply networks combined with an ample but moderate amount of buffering at the I/O ports will lower the minimum IQ required for programming Raw!

# 11 INTERRUPTS

## 11.0  Types of interrupts supported

There are two classes of interrupts that we support, synronous and asynchronous.

The asynchronous interrupts have some degree of flexibility in when they can issue relative to when the processor knows that the interrupting condition exists. There are three asynchronous interrupts supported: the external interrupt, the general network availability interrupt, and the timer interrupt.

The synchronous interrupts, on the other hand, must be initiated at certain place in the code stream. Deferring these interrupts can be disastrous. There are three synchronous interrupts supported; the GDN_COMPLETE ("message finished") interrupt, the GDN_REFILL ("dynamic refill") interrupt, and the TRACE interrupt.

### 11.0.1  Two Interrupt Levels

The TIMER and EXTERNAL interrupts need to be operational even in the middle of a GDN_AVAIL interrupt. This means that we need to be able to take interrupts with a interrupt handler. We could achieve this by requiring the GDN_AVAIL handler to mask the appropriate bits and return to non-interrupt mode. However, doing so requires that the user save registers, and perform a number of operations which increases the message response latency and occupancy.

Instead, we have two interrupt types, "SYSTEM" and "USER", which use different sets of SPRs to save off the PC. The "USER" interrupt does not disable SYSTEM interrupts from happening. However, the SYSTEM interrupts do disable the "USER" interrupt. The GDN_AVAIL interrupt is the only interrupt which operates at the USER level.

## 11.1 Priority of the Interrupts

0. GDN_REFILL (highest)
1. GDN_COMPLETE
2. TRACE
3. EXTERNAL
4. TIMER
5. GDN_AVAIL
6. EVENT_COUNTER  (lowest)

## 11.2 Interrupt Procedure

**First Cycle:**

An interrupt occurs if the following value is non-zero (a find last one operation is executed on this value to find the highest priority interrupt):

(EX_MASK & EX_BITS[5:0]) &
(EX_BITS[USER] 1 1 1 1 1)
& (EX_BITS[SYSTEM] replicated 6 times)

The interrupt is converted to a vector by shifting. This address is fed to the fetch unit (so it can be fetched on the next cycle.) Meanwhile, the processor decode stage stalls. If the interrupt is [0..4], it is a system interrupt, EX_PC is saved, and EX_BITS[SYSTEM] is set to 0. If the interrupt is the GDN_AVAIL interrupt (i.e. a USER interrupt), EX_UPC is saved, and EX_BITS[USER] is set to 0 (system interrupts can still occur).

**Second cycle:**

The decode stage stalls again. The first instruction is being fetched. The address of the next instruction at (VEC+4)  is setup for the IMEM.

**Third cycle:**

The first instruction of the handler is now being executed in the decode stage.

**If a instruction memory store or load instruction is present in the pipeline, the interrupt procedure must be stalled appropriately.  The instruction memory operation has priority.**

**Nth Cycle:**

When a system handler is done, it will execute an ERET. This will restart the processor at address EX_PC and sets  EX_BITS[SYSTEM] to 1.

When a user handler is done, it will execute DRET, which returns to EX_UPC, and sets EX_BITS[USER] to 1.

## 11.3 Synchronous Interrupts

In order to get the behaviour of synchronous interrupts correct, it is important the interrupts be enabled on

the cycle immediately following the DRET. If somehow two synchronous interrupts occur on the same cycle, the highest priority interrupt will be serviced, will return, and immediately the next highest priority interrupt should be serviced before any instructions are executed. If interrupts are not immediately enabled. then an instruction may slip in in-between the interrupts, proving disastrous.

Synchronous interrupts must be blocked by the SYSTEM interrupt mask because more than one of them may trigger at the same time. Furthermore, they must disable interrupts upon activation for the same reason.

### 11.3.1 GDN Refill Interrupt

The GDN refill interrupt is enabled by setting the GDN_REFILL bit in the EX_MASK register.

When a read to cgni occurs in the decode stage, it will be substituted with the value in the GDN_RF_VAL register. If the instruction does not stall, the GDN_REFILL bit of the EX_BITS register will be asserted for the next cycle.

Since the synchronous interrupts have priority over the asynchronous resets, this interrupt is gauranteed to occur on the next cycle (since the GDN_REFILL bit is presumably enabled.) This is the way it must be, the next instruction might read $cgni and thus will have need of the refill services.

The firing of the GDN_REFILL interrupt will reset the GDN_REFILL bit of the EX_BITS register.

### 11.3.2 GDN Complete Interrupt

When the OS wants to context switch, it needs to ensure that the user finishes any incomplete messages on the GDN network. It does this by checking the GDN_PENDING field of the GDN_BUF SPR. This register reports how many words are required to complete the current message. The OS should make sure to read this register at least 5 or 6 instructions after the last write to $CGNO, since it will not be up-to-date until the results leave the pipeline and enter the $CGNO buffer.

(A special IB_WATCHER watches the values heading to CGNO from the processor pipeline. It keeps track of each message that the processor issues, maintaining a count of how many words are left to go. This value is reflected in the so-called PENDING field of the GDN_BUF SPR.)

To enable the GDN_COMPLETE interrupt, the OS sets the GDN_REMAIN SPR to the value of the PENDING field of the GDN_BUF SPR. The OS also enables the GDN_COMPLETE bit of the EX_MASK. GDN_REMAIN decrements each time an instruction issues which writes to CGNO.

When the GDN_REMAIN counter reaches zero, and the GDN_COMPLETE bit of EX_MASK is set, the GDN_COMPLETE bit of the EX_BITS is set, meaning that the message is complete. An interrupt will occur before the next instruction is issued.

* Note that a direct set of the GDN_REMAIN counter has a refractory period of up to one cycle; i.e., it will not be decremented if a instruction that writes to CGNO is placed in the cycle immediately after the MTSR to GDN_REMAIN. *

In the likely case, the operating system will also set the watchdog timer and give the user some reasonable amount of time to complete the send. If the user fails to do so, the operating system will fill write in data elements until the message has been completed, and then kill the process.

It is important to get the timing right of this interrupt, so that interrupt is asserted before the user can issue another send back-to-back with the last one.

### 11.3.3 Trace Interrupt

The trace interrupt is used for debugging support of the Raw processor. There are both TRACE bits in both the EX_BIT and the EX_MASK SPRs. When the processor wants to single step, it will (with interrupts off) set the EX_MASK trace interrupt bit and perform the ERET operation. The processor will execute exactly one instruction at user level, at which point the EX_BIT bit will be set to one, enabling an interrupt. Although other synchronous interrupts may fire before this one, the trace interrupt will fire before the second instruction has executed in user mode.

The EX_BIT bit is reset automatically when the processor vectors to the TRACE handler.

The TRACE bit is enabled in the EX_BIT (and remains enabled until the interrupt fires) when:

( (!stalled last cycle & EX_BITS[SYSTEM]
  & EX_MASK[TRACE]))

The trace EX_BIT can be cleared by the OS (for instance, for a context switch) by ERET to a nop instruction. The trace bit EX_BIT can be set (for instance, to context switch back to a program) in a subtle way by the OS through the GDN_COMPLETE interrupt mechanism. However, since the trace operation is an OS facility, it is not expected that this would be necessary. Most likely the OS would dispatch directly to the TRACE interrupt vector on its own in order to emulate this functionality.

The trace bit can be used to emulate the GDN complete interrupt, theoretically.

The primary purpose of the trace bit is to allow single-stepping, and to allow stepping over breakpoints with out implementing a partial machine simulator.

The trace bit can be found in the motorola 68K and x86 processors.

## 11.4 Asynchronous Interrupts

### 11.4.1 External (MDN) Interrupts

External interrupts are triggered by one-word messages sent by the off-chip interrupt controller. These interrupts indicate that some device requires the processor's attention. The operating system may also trigger the interrupt controller to perform some operation on its behalf.

The delivery mechanism for external interrupts are described in more detail in the "Dynamic Network section." The memory network is used because its discipline gaurantee deadlock free and relatively timely delivery of interrupts.

### 11.4.2 Timer Interrupts

The watchdog timer is configured to initiate an interrupt if WATCH_VAL == WATCH_MAX. There are various conditions which can reset the value of WATCH_VAL.

### 11.4.3 GDN Avail Interrupt

The GDN_AVAIL interrupt is asserted if data is available on the general network input queue. This allows the processor to perform receive unexpected dynamic messages without polling. This interrupt is assumed to be under user control, and has the lowest pri-

ority. It is important that this interrupt not disable the deadlock detection mechanism, else we will not be able to recover from deadlock on the general network.

### 11.4.4 Interrupt Levels

The TIMER and EXTERNAL interrupts need to be operational even in the middle of a GDN_AVAIL event. We could achieve this by requiring the user to mask the appropriate bits and return to non-interrupt mode. However, doing so requires that the user save registers, and perform a number of operations which increases the message response latency and occupancy.

In the interests of decreasing these times, we place the GDN_AVAIL interrupt at a lower level (let's call it "user level.") This allows other interrupts to continue to occur even in the user's general network handler.

The other four interrupt types do not have such hefty requirements for latency, and can be placed in the same priority group (let's call it "system level.")

With two interrupt priority levels, we need two registers to save the exception PCs at both levels. Hence, we have both the EX_PC and EX_UPC special purpose registers.

### 11.4.5 Masks and enables

The EX_ENABLE read-only register has two bit positions, one indicating if system level interrupts are disabled, and the other indicating if user level interrupts are disabled.

The EX_MASK register allows the user to mask out various interrupts.

The EX_BITS register registers the presence of interrupt requests. These bits are implicitly cleared when:

| GDN_COMPLETE | when the interrupt fires. |
|---|---|
| GDN_AVAIL | data is removed from general network |
| TIMER | when the interrupt fires |
| TRACE | when the interrupt fires |
| EXTERNAL | when the interrupt fires |
| GDN_REFILL | when the interrupt fires |

| EVENT_COUNTERS | when EVENT_BITS is zero |
|---|---|

This subtlely allows us to avoid having special instructions to set individual bits, which is required for atomicity reasons. If we read the value of EX_BITS, changed it and wrote it back, an interrupt might have come in in the meantime, causing us to loose an interrupt.

The GDN_AVAIL bit is wired more or less directly (perhaps with one intermittant register) to the "validout" wire of the CGNI SIB.

The GDN_COMPLETE bit of the EX_BITS register is set by a special IB_WATCHER which watches values heading to $cgno. It can be hooked up in a clever fashion to count the number of words that are left in the message.

The EXTERNAL bit is an actual register, which is set to a '1' when an external interrupt message arrives on the MDN.

The TIMER bit is set to '1' anytime that the timer reaches "MAXVAL", and the timer interrupt is enabled.

The DRET, ERET, INTOFF and INTON instructions enable and disable all interrupts (user and system). The UINTOFF and UINTON instructions enable and disables user interrupts. These must be effective immediately in the beginning of the EXECUTE stage (so that an interrupt does not occur on the next cycle in the DECODE stage), so that interrupts do not happen after they are executed. Note the diagram "SPR Update Logic" which shows a subtle difference in how these are implemented.

When an interrupt fires, it disables interrupts for the next cycle, after the effects of the DRET, ERET, INTOFF, and INTON. See "SPR UPDATE LOGIC."

Circuitry for Changes to SYSTEM/USER
(ie xINTON/xINTOFF, DRET/ERET
instrs in Execute stage)

Circuitry for changes to all other SPRs

**SPR UPDATE LOGIC**

54

# 12 MULTITASKING

## 12.0  MULTITASKING

One of the many big headaches in processor design is enabling multitasking -- the running of several processes at the same time. This is not a major goal of the Raw project. For instance, we do not provide a method to protect errant processes from modifying memory or abusing I/O devices. It is nonetheless important to make sure that our architectural constructs are not creating any intractable problems. Raw could support both spatial and temporal multitasking.

In spatial multitasking, two tiles could be running separate processes at the same time. However, a mechanism would have to be put in place to prevent spurious dynamic messages from obstructing or confusing unrelated processes. A special operating system tile could be used to facilitate communication between processes.

## 12.1 CONTEXT SWITCHING

Temporal multitasking creates problems because it requires that we be able to snapshot the state of a Raw processor at an unknown location in the program and restore it back later. Such a context switch would presumably be initiated by a dynamic message on the Memory network. Saving the state in the main processor would be much like saving the state of a typical microprocessor. Saving the state of the switch involves freezing the switch, and loading in a new program which drain all of the switch's state into the processor.

The dynamic and static networks present more of a challenge. In the case of the static network, we can freeze the switches, and then inspect the count of values in the input buffers. We can change the PC of the switch to a program which routes all of the values into the processor, and then out to the southeast shared DRAM over the high-priority dynamic network. Upon return from interrupt, that tile's neighbor can route the elements back into the SIBs. Unfortunately, this leaves no recourse for tiles on the edges of the chip, which do not have neighbor tiles. This issue will be dealt with later in the section.

The dynamic network is somewhat easier. In this case, we can assume command of all of the tiles so that we know that no new messages are being sent. Then we can have all of the tiles poll and drain the messages out of the network. The tiles can examine the buffer counts on the dynamic network SIBs to know when they are

done. Since they can't use the dynamic network to indicate when they are done (they're trying to drain the network!) they can use the common DRAM, or the static network to do so. Upon return, it will be as if the tile was recovering from deadlock; the DYNAMIC REFILL mechanism would be used. For messages that are in the commit buffer, but have not been LAUNCHed, we provide a mechanism to drain the commit buffer.

### 12.1.1  Context switches and I/O Atomicity

One of the major issues with exposing the hardware I/O devices to the compiler and user is I/O atomicity. This is a problem that occurs any time resources are multiplexed between clients. For the most part, we assume that a higher-order process (like the operating system) is ensuring that two processes don't try to write the same file or program the same sound card.

However, since we are exposing the hardware to the software, there is another problem. Actions which were once performed in hardware atomically are now in software, and are suddenly not atomic. For instance, on a request to a DRAM, getting interrupted before one has read the last word of the reply could be disastrous.

The user may be in the middle of issuing a message, but suddenly get swapped out due to some sort of context switch or program exit. The next program that is running may initiate a new request with the device. The hardware device will now be thoroughly confused. Even if we are fortunate enough that it just resets and ignores the message, the programs will probably blithely continue, having lost (or gained) some bogus message words. I call this the I/O Message Atomicity problem.

There is also the issue that a device may succeed in issuing a request on one of the networks, but context switch before it gets the reply. The new program may then receive mysterious messages that were not intended for it. I call this the I/O Request Atomicity problem.

The solution to this problem is to impose a discipline upon the users of the I/O devices.

### 12.1.1.1 Message atomicity on the static network

To issue a message, enclose the request in an interrupt disable/enable pair. The user must guarantee that this action will cause the tile to stall with interrupts disabled for at most a small, bounded period of time.

This may entail that the tile synchronize with the switches to make sure that they are not blocked because

they are waiting for an unrelated word to come through.

It also means that the message size must not overflow the buffer capacity on the way to the I/O node, or if it does, the I/O device must have the property that it sinks all messages after a small period of time.

### 12.1.1.2 Message atomicity on the dynamic network

If the commit buffer method is used for the high-or-low priority dynamic networks, then the message send is atomic. If the commit buffer method is not used, then again, interrupts must be disabled, as for the static network. Again, the compiler must guarantee that it will not block indefinitely with interrupts turned off. It must also guarantee that sending the message will not result in a deadlock.

### 12.1.1.3 Request Atomicity

Request atomicity is more difficult, because it may not feasible to disable interrupts, especially if the time between a request and a reply is long.

However, for memory accesses, it is reasonable to turn off interrupts until the reply is received, because we know this will occur in a relatively small amount of time. After all, standard microprocessors ignore interrupts when they are stalled on a memory access.

For devices with longer latencies (like disk drives!), it is not appropriate to turn off interrupts. In this case, we really are in the domain of the operating system. One or more tiles should be dedicated to the operating system. These tiles will never be context switched. The disk request can then be proxied through this OS tile. Thus, the reply will go to the OS tile, instead of the potentially swapped out user tile. The OS tile can then arrange to have the data transferred to the user's DRAM space (possibly through the DMA port), and potentially wake up the user tile so it can operate on the data.

### 12.2 SUMMARY

In this section, I showed a strategy which enables us to expose the raw hardware devices of the machine to the user and still support multi-tasking context switches. This method is deadlock free, and allows the user to keep the hardware in a consistent state in the face of context switches.

# 13 THE MULTICHIP PROTOTYPE

## 13.0 THE RAW FABRIC / SUPERCOMPUTER

The implementation of the larger Raw prototype creates a number of interesting challenges, mostly having due to with the I/O requirements of such a system. Ideally, we would be able to expose all of the networks of the peripheral tiles to the pins, so that they could connect to an identical neighbor chip, creating the image of a larger Raw chip. Just as we tiled Raw tiles, we will tile Raw chips! To the programmer, the machine would look exactly like a 256 tile Raw chip. However, some of the network hops may have an extra cycle of latency.

## 13.1 PIN COUNT PROBLEMS AND SOLUTIONS

Our package has a whopping 1124 signal pins. This in itself is a bit of a problem, because building a board with 16 such chips is non-trivial. Fortunately, our mesh topology makes building such a board easier. Additionally, the possibility of ground bounce due to simultaneously switching pins is sobering.

For the ground bounce problem, we have a potential solution which reduces the number of pins that switch simultaneously. It involves sending the negation of a signal vector in the event that more than half of the pins would change values. Unfortunately, this technique requires an extra pin for every thirty-two pins, exacerbating our pin count problem. Alternatively, we can switch to a differential signalling system (two pins per signal), clock the pins at twice the rate, so that we use the same number of pins. The signal gurus inform me that this creates a net cancellation of supply strain.

Unfortunately, 1124 pins is also not enough to expose all of the peripheral networks to the edges of the chip so that the chips can be composed to create the illusion of one large tile. The table entitled "Pin Count - ideal" shows the required number of pins. In order to build the Raw Fabric, we needed to find a way to reduce the pin usage.

### TABLE 5. Pin Count - ideal

| Purpose | Count |
| --- | --- |
| Testing, Clocks, Resets, PSROs | 10 |
| Dynamic Network Data | 32x2x16 |
| Dynamic Network Thanks Pins | 2x2x16 |
| Dynamic Network Valid Pins | 1x2x16 |
| Dynamic Network Mux Pins | 1x2x16 |
| Static Network Data | 32x2x16 |
| Static Network Thanks Pins | 1x2x16 |
| Static Network Valid Pins | 1x2x16 |
| Total | 70*32+10 |
| | = 2250 |

We explored a number of options:

### 13.1.1 Expose only the static network

One option was to expose only the static network. Originally, we had opted for this alternative. However, over time, we became more and more aware of the importance of having a dynamic I/O interface to the external world. This is particularly important for supporting caching. Additionally, not supporting the dynamic network means that many of our software systems would not work on the larger system.

### 13.1.2 Remove a subset of the network links

For the static network, this is not a problem -- the compiler can route the elements accordingly through network to avoid the dead links.

For a dimension ordered wormhole routed network, a sparse mesh created excruciating problems. Suddenly, we have to route around the "holes", which means that the sophistication of the dynamic network would have to increase drastically. It would be increasingly hard to remain deadlock free.

### TABLE 6. Pin Count - with muxing

| Purpose | Count |
| --- | --- |
| Testing, Clocks, Resets, PSROs | 10 |
| Network Data | 32x2x16 |

**TABLE 6. Pin Count - with muxing**

| Purpose | Count |
| --- | --- |
| Dynamic Network Thanks | 2x2x16 |
| Dynamic Network Valid | 1x2x16 |
| Mux Pins | 2x2x16 |
| Static Network Thanks | 1x2x16 |
| Static Network Valid Pins | 1x2x16 |
| Total | 39*32+10 |
| | = 1258 |

### 13.1.3 Do some more muxing

The alternative is to retain all of the logical links and mux the data pins. Essentially, the static, dynamic and high-priority dynamic networks all become logical channels. We must add some control pins which select between the static, dynamic and high-priority dynamic networks. See the Table entitled "Pin Count - with muxing."

### 13.1.4 Do some encoding

The next option is to encoding the control signals:

**TABLE 7. States -- encoded**

| State | Value |
| --- | --- |
| No value | 0 |
| Static Value | 1 |
| Memory Dynamic | 2 |
| Low Priority Dynamic | 3 |

This encoding combines the mux and valid bits. Individual thanks lines are still required.

**TABLE 8. Pin Count - with muxing and encoding**

| Purpose | Count |
| --- | --- |
| Testing, Clocks, Resets, PSROs | 10 |
| Network Data | 32x2x16 |
| Dynamic Network Thanks | 2x2x16 |
| Encoded Mux Pins | 2x2x16 |
| Static Network Thanks | 1x2x16 |
| Total | 37*32+10 |
| | = 1194 |

At this point, we are only 70 pins over budget. At this point, we can:

### 13.1.5 Pray for more pins

The fates at IBM may smile upon us and provide us with a package with even better pin counts. We're not too far off.

### 13.1.6 Find a practical but ugly solution

As a last resort, there are some skanky but effective techniques that we can use. We can multiplex the pins of two adjacent tiles, creating a lower bandwidth stripe across the Raw chip. Since these signals will not be coming from the same area of the chip, the latency will probably increase (and thus, the corresponding SIB buffers). Or, we can reduce the data sizes of some of the paths to 16 bits and take two cycles to send a word.

More cleverly, we can send the value over as a 16 bit signed number, along with a bit which indicates if the value fit entirely within the 16 bit range. If it did not, the other 16 bits of the number would be transmitted on the next cycle.

### 13.2 SUMMARY

Because of the architectural headaches involved with exposing only parts of the on-chip networks, we have decided to use a variety of muxing, encoding and praying to solve our pin limitations. These problems are however, just the beginning of the problems that the MULTI-CHIP Raw system of 2007 would encounter. At that time, barring advances in optical interconnects, there will have an even smaller ratio of pins to tiles. At that time, the developers will have to derive more clever dynamic networks [Glass92], or will have to make heavy use of the techniques described in the "skanky solution" category.

# 14 CONCLUSIONS

## 14.0 CURRENT PROGRESS ON THE PROTOTYPE

We are fully in the midst of the implementation effort of the Raw prototype. I have written a C++ simulator named btl, which corresponds exactly to the prototype processor that we are building. It accurately models the processor on a cycle-by-cycle basis, at a rate of about 8000 cycles per second for a 16 tile machine. My pet multi-threaded, bytecode compiled extension language, bC, allows the user to quickly prototype external hardware devices with cycle accurate behaviour. The bC environment provides a full-featured programmable debugger which has proven very useful in finding bugs in the compiler and architecture. I have also written a variety of graphic visualization tools in bC which allow the user to gain a qualitative feel of the behaviour of a computation across the Raw chip. See the Appendages entitled "Graphical Instruction Trace Example" and "Graphical Switch Animation Example." Running wordcount reveals that the simulator, extension language, debugger and user interface code total 30,029 lines of.s,.cc,.c,.bc, and.h files. This does not include the 20,000 lines of external code that I integrated in.

Rajeev Barua and Walter Lee's parallelizing compiler, RawCC, has been in development for about two years. It compiles a variety of benchmarks to the Raw simulators. There are several ISCA and ASPLOS papers that describe these efforts.

Matt Frank and I have ported a version of GCC for use on serial and operating system code. It uses inline macros to access the network ports.

Ben Greenwald has ported the GNU binutils to support Raw binaries.

Jason Kim, Sam Larsen, Albert Ma, and I have written synthesizeable verilog for the static and dynamic networks, and the processors. It runs our current code base, but does not yet implement all of the interrupt handling and deadlock recovery schemes.

Our testing effort is just beginning. We have Krste Asanovic's automatic test vector generator, called Torture, which generates random test programs for MIPS processors. We intend to extend it to exert the added functionality of the Raw tile.

We also have plans to emulate the Raw verilog. We have a IKOS logic emulator for this purpose.

Jason Kim and I have attended IBM's ASIC training class in Burlington, VT. We expect to attend the Static Timing classes later in the year.

A board for the Raw handheld device is being developed by Jason Miller.

This document will form the kernel of the design specification for the Raw prototype.

## 14.1 PRELIMINARY RESULTS

We have used the Raw compiler to compile a variety of applications to the Raw simulator, which is accurate to within %10 of the actual Raw hardware. However, in both the base and parallel case, the tile has unlimited local SRAM. Results are summarized below.

**TABLE 9. Preliminary Results - 16 tiles**

| Benchmark | Speedup versus one tile |
|---|---|
| Cholesky | 10.30 |
| Matrix Mul | 12.20 |
| Tomcatv | 9.91 |
| Vpenta | 10.59 |
| Adpcm-encode | 1.26 |
| SHA | 1.44 |
| MPEG-kernel | 4.48 |
| Moldyn | 4.48 |
| Unstructured | 5.34 |

More information on these results is given in [Barua99].

Mark Stephenson, Albert Ma, Sam Larsen, and I have all written a variety of hand-coded applications to gain an idea of the upper bound on performance for a Raw architecture. Our applications have included median filter, DES, software radio, and MPEG encode. My hand-coded application, median filter, has 9 separate interlocking pipeline programs, running on 128 tiles, and attains a 57x speedup over a single issue processor, compared to the 4x speedup that a hand-coded dual-issue Pentium with MMX attains. Our hope is that the Raw supercomputer, with 256 MIPS tiles, will enable us to attain similarly outrageous speedup numbers.

## 14.2 EXIT

In this design document, I have traced the design decisions that we have made along the journey to creating the first Raw prototype. I detail how the architecture was born from our experience with FPGA computing. I familiarize the reader with Raw by summarizing the programmer's viewpoint of the current design. I motivate our decision to build a prototype. I explain the design decisions we made in the implementation of the static and dynamic networks, the processor, and the prototype systems. I finalize by showing some results that were generated by our compiler and run on our simulator.

The Raw prototype is well on its way to becoming a reality. With many of the key design decisions determined, we now have a solid basis for finalizing the implementation of the chip. The fabrication of the chip and the two systems will aid us in exploring the application space for which Raw processors are well suited. It will also allow us to evaluate our design and prove that Raw is, indeed, a realizable architecture.

## 14.3 REFERENCES

J.L. Hennessey, **"The Future of Systems Research,"** IEEE Computer Magazine, August 1999. pp. 27-33.

D. L. Tennenhouse and V. G. Bose, **"SpectrumWare - A Software-Oriented Approach to Wireless Signal Processing,"** ACM Mobile Computing and Networking 95, Berkeley, CA, November 1995.

R. Lee, **"Subword Parallelism with MAX- 2"**, IEEE Micro, Volume 16 Number 4, August 1996, pp. 51-59.

J. Babb et al. **"The RAW Benchmark Suite: Computation Structures for General Purpose Computing,"** IEEE Symposium on Field-Programmable Custom Computing Machines, Napa Valley, CA, April 1997.

Agarwal et al. "**The MIT Alewife Machine: Architecture and Performance,"** Proceedings of ISCA '95, Italy, June, 1995.

Waingold et al. **"Baring it all to Software: Raw Machines,"** IEEE Computer, September 1997, pp. 86-93.

Waingold et al. **"Baring it all to Software: Raw Machines,"** MIT/LCS Technical Report TR-709, March 1997.

Walter Lee et al. **"Space-Time Scheduling of Instruction-Level Parallelism on a Raw Machine,"** Proceedings of ASPLOS-VIII, San Jose, CA, October 1998.

R. Barua et al. **"Maps: A Compiler Managed Memory System for Raw Machines,"** Proceedings of the Twenty-Sixth International Symposium on Computer Architecture (ISCA), Atlanta, GA, June, 1999.

T. Gross. **"A Retrospective on the Warp Machines,"** 25 Years of the International Symposia on Computer Architecture, Selected Papers. 25th Anniversary Issue. 1998. pp 45-47.

J. Smith. **"Decoupled Access/Execute Computer Architectures,"** 25 Years of the International Symposia on Computer Architecture, Selected Papers. 25th Anniversary Issue. 1998. pp 231-238. (Originally in ISCA 9)

W. J. Dally. **"The torus routing chip,"** Journal of Distributed Computing, vol. 1, no. 3, pp. 187-196, 1986.

J. Hennessey, and D. Patterson **"Computer Architecture: a Quantitative Approach (2nd Ed.)"**, Morgan Kauffman Publishers, San Francisco, CA, 1996.

M. Zhang. **"Software Floating-Point Computation on Parallel Machines,"** Master's Thesis, Massachusetts Institute of Technology, 1999.

S. Oberman. **"Design Issues in High Performance Floating Point Arithmetic Units,"** Ph.D. Dissertation, Stanford University, December 1996.

E. Berlekamp, J. Conway, R. Guy, **"Winning Ways for Your Mathematical Plays,"** vol. 2, chapter 25, Academic Press, New York, 1982.

John D. Kubiatowicz. "**Integrated Shared-Memory and Message-Passing Communication in the Alewife Multiprocessor,"** Ph.D. thesis, Massachusetts Institute of Technology, Department of Electrical Engineering and Computer Science, February 1998.

C. Moritz et al. "**Hot Pages: Software Caching for Raw Microprocessors,"** MIT CAG Technical Report, Aug 1999.

Fred Chong et al. **"Remote Queues: Exposing Message Queues for Optimization and Atomicity,"** Symposium on Parallel Algorithms and Architecture (SPAA) Santa Barbara, July 1995.

K. Mackenzie et al. "**Exploiting Two-Case Delivery for Fast Protected Messaging."** Proceedings of 4th International Symposium on High Performance Computer Architecture Feb. 1998.

C. J. Glass et al. **"The Turn Model for Adaptive Routing,"** 25 Years of the International Symposia on Computer Architecture, Selected Papers. 25th Anniversary Issue. 1998. pp 441-450. (Originally in ISCA 19)

# **15** APPENDAGES

Raw user's manual

Massachusetts Institute of Technology
Laboratory of Computer Science

# RAW 823 and 824
# ISA
# User's Manual

# 1.0  Foreword

This document is the ISA manual for the Raw prototype processor. Unlike other Raw documents, it does not contain any information on design decisions, rather it is intended to provide all of the information that a software person would need in order to program a Raw processor. This document assumes a familiarity with the MIPS architecture. If something is unspecified, one should assume that it is exactly the same as a MIPS R4000.
(See http://www.mips.com/publications/index.html, "R4000 Microprocessor User's Manual".)

# 2.0  Processor

Each Raw Compute Processor looks very much like a MIPS R4000.

The follow items are different:

0. Registers 24, 25,  26, and 27 are used to address network ports and are not available as GPRs.
1. Floating point operations use the same register file as integer operations.
2. Floating point compares have a destination register instead of setting a flag.
3. The floating point branches, BC1T and BC1F are removed, since the integer versions have equivalent functionality.
4. Instead of a single multiply instruction, there are three low-latency instructions, MULH, MULHU, and MULLO which place their results in a GPR instead of HI/LO.
5. The pipeline is eight stage pipeline, with FETCH, DECODE, RF/STALL, EXE, MUL, MEM, FPU and WB stages.
6. Floating point divide uses the FD register instead of a destination register.
7. The instruction set, the timings and the encodings are slightly different. The following section lists all of the instructions available in the processor. There are some omissions and some additions. For actual descriptions of the standard computation instructions, please refer to the MIPS manual. The non-standard raw instructions (marked with **823** and **824**) will be described later in this document.
7a. A tile has  8k words (32bit word) of  local instruction memory, and 8k words (64 bit word) of switch memory. The Raw 824 ISA also supports a mode with a cached instruction memory.
8. A tile has a 8k word (32 bit word) , 2-way set-associative, 3 cycle latency data cache with 32 byte lines.  For replacement, it has the following replacement policy: use set 0 if it is invalid, otherwise use set 1 if it is invalid, otherwise use the least recently used of the two sets. The cache uses a 2 element bypassing write buffer to defer stores until after the tag has been checked. See the TAGSW instruction for more information.
9. cvt.w does round-to-nearest even rounding (instead of a "current rounding mode"). the trunc operation (which is the only one used by GCC) can be used to round-to-zero.
10. All floating point operations are single precision.
11. The Raw prototype is a LITTLE ENDIAN processor. In other words, if there is a word stored at address P, then the low order byte is stored at address P, and the most significant byte is stored at address P+3. (Sparc, for reference, is big endian.)
12. Each non-branching instruction instruction has one bit reserved in the encoding, called the S-bit. The S-bit determines if the result of the instruction is written to $csto port, in addition to the register file. If the instruction has no output, the behaviour of the S-bit is undefined. The S-bit is set by using an exclamation point with the instruction, as follows:

and!          $3,$2,$0                              # writes to static switch and r3

13. All multi-cycle non-branch operations (loads, multiplies, divides) on the raw processor are fully interlocked.
14. We use static branch prediction, and no delay slots. A "+" appended to the end of the conditional branch/jmp indi-

cates that the branch is likely taken, a "-" indicates that the branch is likely not taken. If it is unspecified, then the assembler applies a simple heuristic (backwards branches taken) and sets the appropriate bit. The mispredict penalty is 3 cycles. For instance   beq+ $csto, loop   indicates a branch likely to "loop." The same holds true for the switch processor.

# 3.0  Register Conventions

 The following register convention map has been modified for Raw from page D-2 of the MIPS manual). Various software systems by the raw group may have more restrictions on the registers.

**Table 1: Register Conventions**

| reg | alias | Use |
| --- | --- | --- |
| $0 | | Always has value zero. |
| $1 | $at | Reserved for assembler |
| $2..$3 | | Used for expression evaluation and to hold procedure return values. |
| $4..$7 | | Used to pass first 4 words of actual arguments. Not preserved across procedure calls. |
| $8..$15 | | Temporaries. Not preserved across procedure calls |
| $16..$23 | | Callee saved registers /<br>Must be restored by interrupt handlers for i-caching to function. |
| $24 | $cst[i/o] | Static network input port. |
| $25 | $cgn[i/o] | User Dynamic Network input/output port. |
| $26 | $csti2 | Second static network input port. |
| $27 | $cmn[i/o] | Memory Dynamic network input/output port. |
| $28 | $gp | Global pointer. |
| $29 | $sp | Stack pointer. |
| $30 | | A callee saved register. |
| $31 | | The link register. |

The switch processor has the following register conventions:

$0,$1,$2 -- caller saved.   $3 -- link register

# 4.0 Instruction Set

Sample Instruction Listing:



The Raw Chip implements the 823 ISA.
The 824 ISA is available only in simulation, and only when "824 mode" is enabled.

# 5.0 Integer Computation Instructions

**ADDIU**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|
| s | ADDIU | rs | | rt | | signed immediate | |
| 1 | 5 | 5 | | 5 | | 16 | |

ADDIU rt, rs, imm    1

**ADDU**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | | rt | | rd | | 0 0 0 0 0 | | ADDU 1 0 0 0 0 1 | |
| 1 | 5 | 5 | | 5 | | 5 | | 5 | | 6 | |

ADDU rd, rs, rt    1

**AND**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | | rt | | rd | | 0 0 0 0 0 | | AND 1 0 0 1 0 0 | |
| 1 | 5 | 5 | | 5 | | 5 | | 5 | | 6 | |

AND rd, rs, rt    1

**ANDI**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|
| s | ANDI 0 0 1 0 0 | rs | | rt | | unsigned immediate | |
| 1 | 5 | 5 | | 5 | | 16 | |

ANDI rt, rs, imm    1

**AUI**
**823**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|
| s | AUI 0 0 1 1 1 | rs | | rt | | unsigned immediate | |
| 1 | 5 | 5 | | 5 | | 16 | |

AUI rt, rs, imm    1

Replaces LUI. [rt] ← [rs] + (imm << 16)

**BL**
**824**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|---|
| 1 | B 1 1 0 0 0 | | signed offs4 | | | | | |
| 1 | 5 | | 26 | | | | | |

BL  offs4    1

Enables full-address space hardware instruction caching. Relative branch reachs  +/− 128 MB of code.

**BLAL**
**824**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|---|
| 1 | BAL 1 1 0 0 1 | | signed offs4 | | | | | |
| 1 | 5 | | 26 | | | | | |

BLAL  offs4    1

Enables full-address space hardware instruction caching. Relative function call.
[r31] <- PC+4
branch by offs4

**BEQ**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|---|
| p | BEQ 1 1 0 1 1 | rs | | rt | | signed offs4 | | |
| 1 | 5 | 5 | | 5 | | 16 | | |

BEQ rs, rt, offs    1m

**BGEZ**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | | BGEZ 0 0 0 1 0 | | signed offs4 | |
| 1 | 5 | 5 | | 5 | | 16 | |

BGEZ rs, offs    1m

**BGEZAL**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | BGEZAL 1 0 0 1 0 | signed offs4 | BGEZAL rs, offs | 1m |
| 1 | 5 | 5 | 5 | 16 | | |

**BGTZ**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | BGTZ 0 0 0 1 1 | signed offs4 | BGTZ rs, offs | 1m |
| 1 | 5 | 5 | 5 | 16 | | |

**BLEZ**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | BLEZ 0 0 0 0 1 | signed offs4 | BLEZ rs, offs | 1m |
| 1 | 5 | 5 | 5 | 16 | | |

**BLTZ**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | BLTZ 0 0 0 0 0 | signed offs4 | BLTZ rs, offs | 1m |
| 1 | 5 | 5 | 5 | 16 | | |

**BLTZAL**

| 31 | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | BLTZAL 1 0 0 0 0 | signed offs4 | BLTZAL rs, offs | 1m |
| 1 | 5 | 5 | 5 | 16 | | |

**BNE**

| 31 | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|
| p | BNE 1 1 0 1 0 | rs | rt | signed offs4 | BNE rs, rt, offs | 1m |
| 1 | 5 | 5 | 5 | 16 | | |

**BNEA**

| 31 | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|
| p | BNEA 1 1 0 0 1 | rs | rt | signed offs4 | BNE rs, rt, offs | 1m |
| 1 | 5 | 5 | 5 | 16 | | |

Branch and add.
tmp = [rs]; [rs] = tmp + SPR[incr];
if (rt != tmp) branch offs4;

**CLZ**
**823**

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | 0 0 0 0 0 | rd | 0 0 0 0 0 | CLZ 1 1 1 0 0 1 | CLZ rd, rs | 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | |

Count leading zero. Counts the number of leading '0' bits.

**DIV**

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | rt | 0 0 0 0 0 | 0 0 0 0 0 | DIV 0 1 1 0 1 0 | DIV rs, rt | 42 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | |

**DIVU**

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | rt | 0 0 0 0 0 | 0 0 0 0 0 | DIVU 0 1 1 0 1 1 | DIVU rs, rt | 42 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | |

Note: the DIVU/DIV instruction has on cycle of occupancy on the fetch unit, but occupies the divide unit for the whole duration of the instruction.

## J

| 31 | | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| 1 | REGIMM 11000 | 0 0 0 0 0 | J 0 1 1 0 0 | targ4 | | J  targ | 1 |
| 1 | 5 | 5 | 5 | 16 | | | |

---

## JEQL
**823**

| 31 | | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| p | JEQL 1 1 1 1 1 | rs | rt | targ4 | | JEQL rs, rt, targ4 | 1m |
| 1 | 5 | 5 | 5 | 16 | | | |

$[31] \leftarrow$ PC+4;  if ([rs] == [rt]) { PC $\leftarrow$ (targ4 << 2) }

---

## JGEZL
**823**

| 31 | | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | JGEZL 0 0 1 1 0 | targ4 | | JGEZL rs, targ4 | 1m |
| 1 | 5 | 5 | 5 | 16 | | | |

$[31] \leftarrow$ PC+4;  if ([rs] >= 0) { PC $\leftarrow$ (targ4 << 2)}

---

## JGTZL
**823**

| 31 | | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | JGTZL 0 0 1 1 1 | targ4 | | JGTZL rs, targ4 | 1m |
| 1 | 5 | 5 | 5 | 16 | | | |

$[31] \leftarrow$ PC+4;  if ([rs] > 0) { PC $\leftarrow$ (targ4 << 2)}

---

## JLEZL
**823**

| 31 | | 27 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | JLEZL 1 0 1 0 1 | targ4 | | JLEZL rs, targ4 | 1m |
| 1 | 5 | 5 | 5 | 16 | | | |

$[31] \leftarrow$ PC+4;  if ([rs] <= 0) { PC $\leftarrow$ (targ4 << 2)}

---

## JLTZL
**823**

| 31 | | 27 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | rs | JLTZL 1 0 1 0 0 | targ4 | | JLTZL rs, targ4 | 1m |
| 1 | 5 | 5 | 5 | 16 | | | |

$[31] \leftarrow$ PC+4;  if ([rs] < 0) { PC $\leftarrow$ (targ4 << 2)}

---

## JNEL
**823**

| 31 | | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| p | JNEL 1 1 1 1 0 | rs | rt | targ4 | | JNEL rs, rt, targ4 | 1m |
| 1 | 5 | 5 | 5 | 16 | | | |

$[31] \leftarrow$ PC+4;  if ([rs] != [rt]) { PC $\leftarrow$ (targ4 << 2) }

---

## JAL

| 31 | | 26 25 | 21 20 | 16 15 | 0 | | |
|---|---|---|---|---|---|---|---|
| p | REGIMM 1 1 0 0 0 | 0 0 0 0 0 | JAL 1 1 1 0 0 | targ4 | | JAL targ | 1 |
| 1 | 5 | 5 | 5 | 16 | | | |

## JALR

| 31 | | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | 0 0 0 0 0 | 0 0 0 0 0 | 0 0 0 0 0 | JALR 0 0 1 0 0 1 | | JALR  rs | 4 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | | |

**JR**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | | 0 0 0 0 0 | | 0 0 0 0 0 | | 0 0 0 0 0 | | JR 0 0 1 0 0 0 | | JR rs | 4 |
| 1 | 5 | 5 | | 5 | | 5 | | 5 | | 6 | | | |

**JRHOFF**
**824**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | | 0 0 0 0 0 | | 0 0 0 0 0 | | 0 0 0 0 0 | | JRHOFF 0 0 1 0 1 1 | | JRHOFF rs | 4 |
| 1 | 5 | 5 | | 5 | | 5 | | 5 | | 6 | | | |

Jump register, disabling hardware caching. Used to transition from hardware to software cached mode.

**JRHON**
**824**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | | 0 0 0 0 0 | | 0 0 0 0 0 | | 0 0 0 0 0 | | JRHON 0 0 1 0 1 0 | | JRHON rs | 4 |
| 1 | 5 | 5 | | 5 | | 5 | | 5 | | 6 | | | |

Jump register, enabling hardware caching. Used to transition from software to hardware cached mode.

**LB**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|
| s | LB 1 0 0 0 0 | base | rt | signed offset | | LB rt, base(offs) | 3 1 |
| 1 | 5 | 5 | 5 | 16 | | | |

**LBU**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|
| s | LBU 1 0 0 0 1 | base | rt | signed offset | | LBU rt, base(offs) | 3 1 |
| 1 | 5 | 5 | 5 | 16 | | | |

**LH**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| s | LH 1 0 0 1 0 | base | rt | signed offset | | LH rt, base(offs) | 3 1 |
| 1 | 5 | 5 | 5 | 16 | | | |

**LHU**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| s | LHU 1 0 0 1 1 | base | rt | signed offset | | LHU rt, base(offs) | 3 1 |
| 5 | 1 | 5 | 5 | 16 | | | |

**LW**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| s | LW 1 0 1 0 0 | base | rt | signed offset | | LW rt, base(offs) | 3 1 |
| 1 | 5 | 5 | 5 | 16 | | | |

The low two bits of the addression are ignored. However, during a cache miss, they are actually sent off of the edge of the chip, for debugging purposes.

**MAGIC**
**SIM**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | rt | code | | MAGIC 0 0 0 0 0 1 | | magic rt, rs,code | 1 |
| 1 | 5 | 5 | 5 | 5 | | 6 | | | |

This instruction allows the user to extend the instruction set, and does not exist on the real chip. However, the opcode has been selected such that, on the real Raw chip, the instruction merely trashes the output register. So, as long as a program containing MAGIC instructions does not rely on them for correctness, programs containing these instructions will run on the real chip.

**MFFD**
**823**

| 31 | 27 26 | 25 | 16 15 | 11 10 | 6 5 | 0 | |
|----|-------|-----|-------|-------|-----|---|---|
| s | SPECIAL 0 0 0 0 0 | 0 0 0 0 0 0 0 0 0 0 | rd | 0 0 0 0 0 | MFFD 0 1 0 1 0 0 | | MFFD rd |
| 1 | 5 | 10 | 5 | 5 | 6 | | |

FD holds result of floating point divider. Blocks until FP divider is ready.

**MFHI**

| 31 | 27 26 | 25 | 16 15 | 11 10 | 6 5 | 0 | |
|----|-------|-----|-------|-------|-----|---|---|
| s | SPECIAL 0 0 0 0 0 | 0 0 0 0 0 0 0 0 0 0 | rd | 0 0 0 0 0 | MFHI 0 1 0 0 0 0 | | MFHI rd |
| 1 | 5 | 10 | 5 | 5 | 6 | | |

**MFLO**

| 31 | 27 26 | 25 | 16 15 | 11 10 | 6 5 | 0 | |
|----|-------|-----|-------|-------|-----|---|---|
| s | SPECIAL 0 0 0 0 0 | 0 0 0 0 0 0 0 0 0 0 | rd | 0 0 0 0 0 | MFLO 0 1 0 0 1 0 | | MFLO rd |
| 1 | 5 | 10 | 5 | 5 | 6 | | |

**MTFD**
**823**

| 31 | 27 26 | 25 | 21 20 | 6 5 | 0 | |
|----|-------|-----|-------|-----|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | MTFD 0 1 0 1 0 1 | | MTFD rs |
| 1 | 5 | 5 | 15 | 6 | | |

FD holds result of floating point divider. Blocks until FP divider is ready.

**MTHI**
**823**

| 31 | 27 26 | 25 | 21 20 | 6 5 | 0 | |
|----|-------|-----|-------|-----|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | MTHI 0 1 0 0 0 1 | | MTHI rs |
| 1 | 5 | 5 | 15 | 6 | | |

Like MIPS , but for the integer divider only.

**MTLO**
**823**

| 31 | 27 26 | 25 | 21 20 | 6 5 | 0 | |
|----|-------|-----|-------|-----|---|---|
| 0 | SPECIAL 0 0 0 0 0 | rs | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | MTLO 0 1 0 0 1 1 | | MTLO rs |
| 1 | 5 | 5 | 15 | 6 | | |

Like MIPS , but for the integer divider only.

**MULH**
**823**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 | |
|----|-------|-----|-------|-------|-------|-----|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | MULH 1 0 1 0 0 0 | | MULH rd, rs, rt |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | |

The contents of register *rs* and *rt* are multiplied as signed values to obtain a 64-bit result.
The high 32 bits of this result is stored into register *rd*.

Operation:        $[rd] \leftarrow ([rs] *_s [rt])_{63..32}$

**MULHU**
**823**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 | |
|----|-------|-----|-------|-------|-------|-----|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | MULHU 1 0 1 0 0 1 | | MULHU rd, rs, rt |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | |

The contents of register *rs* and *rt* are multiplied as unsigned values to obtain a 64-bit result.
The high 32 bits of this result is stored into register *rd*.
Operation:        $[rd] \leftarrow ([rs] *_u [rt])_{63..32}$

## MULLO
**823**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | | rs | | rt | | rd | | 0 0 0 0 0 | | MULLO 0 1 1 0 0 0 |
| 1 | 5 | | 5 | | 5 | | 5 | | 5 | | 6 |

MULLO rd, rs, rt   2 1

The contents of register *rs* and *rt* are multiplied as signed values to obtain a 64-bit result.
The low 32 bits of this result is stored into register *rd*.

Operation: $[rd] \leftarrow ([rs]*[rt])_{31..0}$

## MULLU
**823**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | | rs | | rt | | rd | | 0 0 0 0 0 | | MULLU 0 1 1 0 0 1 |
| 1 | 5 | | 5 | | 5 | | 5 | | 5 | | 6 |

MULLO rd, rs, rt   2 1

The contents of register *rs* and *rt* are multiplied as unsigned values to obtain a 64-bit result.
The low 32 bits of this result is stored into register *rd*.

Operation:   $[rd] \leftarrow ([rs]*_u[rt])_{31..0}$

## NOR

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | | | rs | | rt | | rd | | 0 0 0 0 0 | | NOR 1 0 0 1 1 1 |
| 1 | 5 | | | 5 | | 5 | | 5 | | 5 | | 6 |

NOR rd, rs, rt   1

## OR

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | | rs | | rt | | rd | | 0 0 0 0 0 | | OR 1 0 0 1 0 1 |
| 1 | 5 | | 5 | | 5 | | 5 | | 5 | | 6 |

OR rd, rs, rt   1

## ORI

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 0 |
|---|---|---|---|---|---|---|---|
| s | ORI 0 0 0 0 0 | | rs | | rt | | unsigned immediate |
| 1 | 5 | | 5 | | 5 | | 16 |

ORI  rt, rs, imm   1

## POPC
**823**

| 31 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | | rs | | 0 0 0 0 0 | | rd | | 0 0 0 0 0 | | POPC 1 1 1 0 0 0 |
| 1 | 5 | | 5 | | 5 | | 5 | | 5 | | 6 |

POPC rd, rs   1

Population count -- sums all "1" bits in the input.

## RLM
**823**

| 31 | 27 | 26 | 25 | 21 | 20 | 16 | 15 | 11 | 10 | 6 | 5 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RLM 1 0 1 X 0 0 | | | rs | | rt | | SA | | MB | | ME | Z |
| 6 | | | 5 | | 5 | | 5 | | 5 | | 5 | |

RLM rt, rs, SA, M   1

Rotate Left And Mask (or, Rotate Right and Mask)
Operation:   [rt] = ([rs] lrot SA) & MASK (MB,ME,Z)

Note that the assembler takes the full mask (e.g., 0xf0f00FFF) and encodes it for you. If it can't be encoded, it will tell you. Also, the user can use RRM assembly alias.

MASK:  case { ME[1], ME[0], Z }
            001:  MASK = { MB[4] MB[3] MB[2] MB[1] MB[0]  ME[4] ME[3] ME[2] }replicated 4 times
            111:  MASK = nibble mask;  i.e. { MB[4] MB[4] MB[4] MB[4] MB[3] MB[3] MB[3] MB[3]
                          MB[2] MB[2] MB[2] MB[2] MB[1] MB[1] MB[1] MB[1]
                          MB[0] MB[0] MB[0] MB[0] ME[4] ME[4]  ME[4] ME[4]
                          ME[3] ME[3]  ME[3] ME[3] ME[2] ME[2]  ME[2] ME[2] }

xx0:  MASK = bit range; if (MB <= ME) bits MB..ME set to 1, else all bits except ME..MB set to 1.

## RLMI
**823**

| 31        27 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| RLMI<br>1 0 1 X 0 1 | rs | rt | SA | MB | ME | Z |
| 6 | 5 | 5 | 5 | 5 | 5 | |

RLMI rt,rs, SA,M    | 1 |

Rotate Left and Mask With Insert
Operation: [rt] = (([rs] lrot SA) & MASK(MB,ME,Z)) | ([rt] & ~MASK(MB,ME,Z))

## RLVM
**823**

| 31        27 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| RLVM<br>1 0 1 X 1 0 | rs | rt | rd | MB | ME | Z |
| 6 | 5 | 5 | 5 | 5 | 5 | |

RLVM rd,rs,rt,M    | 1 |

Rotate Left Variable And Mask
Operation: [rd] =  ([rs] lrot rt) & MASK(MB,ME,Z)
e.g,   RLVM $3,$4,$5,0xffff0000

## SLL

| 31 | 27 26 | 25  21 | 20  16 | 15  11 | 10  6 | 5  0 |
|---|---|---|---|---|---|---|
| s | SPECIAL<br>0 0 0 0 0 | rs | rt | SA | 0 0 0 0 0 | SLL<br>0 0 0 0 0 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SLL  rt, rs,sa    | 1 |

## SLLV

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL<br>0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | SLLV<br>0 0 0 1 0 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SLLV rd, rs, rt    | 1 |

## SLT

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL<br>0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | SLT<br>1 0 1 0 1 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SLT rd, rs, rt    | 1 |

## SLTI

| 31 | 26 25 | 21 20 | 16 15 | 0 |
|---|---|---|---|---|
| s | SLTI<br>1 0 1 1 0 | rs | rt | signed immediate |
| 1 | 5 | 5 | 5 | 16 |

SLTI rt, rs, simm    | 1 |

## SLTIU

| 31 | 27 26 25 | 21 20 | 16 15 | 0 |
|---|---|---|---|---|
| s | SLTIU<br>1 0 1 0 1 | rs | rt | signed immediate |
| 1 | 5 | 5 | 5 | 16 |

SLTIU rt, rs,simm    | 1 |

## SLTU

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL<br>0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | SLTU<br>1 0 1 0 1 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SLTU rd, rs, rt    | 1 |

## SRA

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL<br>0 0 0 0 0 | rs | rt | SA | 0 0 0 0 0 | SRA<br>0 0 0 0 1 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SRA  rt, rs,sa    | 1 |

## SRAV

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL<br>0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | SRAV<br>0 0 0 1 1 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SRAV rd, rs, rt    | 1 |

## SRL

| 31 27 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | rt | SA | 0 0 0 0 0 | SRL 0 0 0 0 1 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SRL  rt, rs,sa    1

## SRLV

| 31 27 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | SRLV 0 0 0 1 1 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SRLV rd, rs,rt    1

## SUBU

| 31 27 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | SUBU 1 0 0 0 1 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

SUBU  rd, rs, rt    1

## SB

| 31 27 | 26 25 | 21 20 | 16 15 | 0 |
|---|---|---|---|---|
| 0 | SB 0 1 0 0 0 | base | rt | signed offset |
| 1 | 5 | 5 | 5 | 16 |

SB  rt, offset(base)    1

## SH

| 31 27 | 26 25 | 21 20 | 16 15 | 0 |
|---|---|---|---|---|
| 0 | SH 0 1 0 1 0 | base | rt | signed offset |
| 1 | 5 | 5 | 5 | 16 |

SH  rt, offset (base)    1

## SW

| 31 27 | 26 25 | 21 20 | 16 15 | 0 |
|---|---|---|---|---|
| 0 | SW 0 1 1 0 0 | base | rt | signed offset |
| 1 | 5 | 5 | 5 | 16 |

SW  rt, offset(base)    1

The low two bits of the address are ignored. However, during a cache miss, they are actually sent off of the edge of the chip, for debugging purposes.

## XOR

| 31 27 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | SPECIAL 0 0 0 0 0 | rs | rt | rd | 0 0 0 0 0 | XOR 1 0 0 1 1 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

XOR  rd, rs, rt    1

## XORI

| 31 27 | 26 25 | 21 20 | 16 15 | 0 |
|---|---|---|---|---|
| s | XORI 0 0 0 1 1 | rs | rt | unsigned immediate |
| 1 | 5 | 5 | 5 | 16 |

XORI rt, rs, imm    1

# 6.0  Floating Point Computation Instructions

## ABS.s

| 31 27 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | 0 0 0 0 0 | rd | fmt | ABS 0 0 0 1 0 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

ABS.s  rd, rs    4
1

**ADD.s**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | rt | rd | fmt | ADD 0 0 0 0 0 0 | ADD.s rd, rs, rt |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

**C.xx.s**

**823**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | rt | rd | fmt | cond 1 1 x x x x | C.xx.s  rd, rs, rt |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

Description:    Precisely like MIPS but the result is stored in *rt*, instead of a flags register.

**CVT.s**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | 0 0 0 0 0 | rd | fmt | CVT.S 1 0 0 0 0 0 | CVT.s rd,  rs |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

**CVT.w**

**823**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | 0 0 0 0 0 | rd | fmt | CVT.W.s 1 0 0 1 0 0 | CVT.w  rd,  rs |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

Description:  Precisely like MIPS but always uses round to nearest even rounding mode.
Converts to fixed point.

**DIV.s**

**823**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | rt | 0 0 0 0 0 | fmt | DIV.s 0 0 0 0 1 1 | DIV.s  rs, rt |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 12 |
|---|
| 1 |

Description:   Precisely like MIPS but the result is stored in the FD register, instead of a FPR.
Note: the DIV.s instruction has one cycle of occupancy on the fetch unit, but occupies the divide
unit for the whole latency of the instruction (12 cycles).

**MUL.s**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | rt | rd | fmt | MULT.s 0 0 0 0 1 0 | MUL.s  rd, rs, rt |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

**NEG.s**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | 0 0 0 0 0 | rd | fmt | NEG.s 0 0 0 1 1 1 | NEG.s  rd, rs |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

**SUB.s**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | rt | rd | fmt | SUB.s 0 0 0 0 0 1 | SUB.s  rd, rs, rt |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

**TRUNC.w**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| s | FPU 0 0 1 1 0 | rs | 0 0 0 0 0 | rd | fmt | TRUNC.w.s 0 0 1 1 0 1 | TRUNC.w rd, rs |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

| 4 |
|---|
| 1 |

Converts to fixed point. Does not use round to nearest even rounding mode.

# 7.0 Floating Point Compare Options

**Table 2: Floating Point Comparison Condition (for c.xxx.s)**

| Predicate | | | Relations(Results) | | | | Invalid operation exception if unordered |
|---|---|---|---|---|---|---|---|
| Cond | Mnemonic | Definition | Greater Than | Less Than | Equal | Unordered | |
| 0 | F | False | F | F | F | F | No |
| 1 | UN | Unordered | F | F | F | T | No |
| 2 | EQ | Equal | F | F | T | F | No |
| 3 | UEQ | Unordered or Equal | F | F | T | T | No |
| 4 | OLT | Ordered Less Than | F | T | F | F | No |
| 5 | ULT | Onordered or Less Than | F | T | F | T | No |
| 6 | OLE | Ordered Less Than or Equal | F | T | T | F | No |
| 7 | ULE | Unordered or Less Than or Equal | F | T | T | T | No |
| 8 | SF | Signaling False | F | F | F | F | Yes |
| 9 | NGLE | Not Greater Than or Less Than or Equal | F | F | F | T | Yes |
| 10 | SEQ | Signaling Equal | F | F | T | F | Yes |
| 11 | NGL | Not Greater Than or Less Than | F | F | T | T | Yes |
| 12 | LT | Less Than | F | T | F | F | Yes |
| 13 | NGE | Not Greater Than or Equal | F | T | F | T | Yes |
| 14 | LE | Less Than or Equal | F | T | T | F | Yes |
| 15 | NGT | Not Greater Than | F | T | T | T | Yes |

# 8.0 Administrative Instructions

**DRET**
**823**

| 0 | COMM 0 1 0 1 1 | 0 0 0 0 0 0 0 0 0 0 | 0 0 0 0 0 | 0 0 0 0 0 | DRET 0 0 0 0 0 0 | DRET | 1 |
|---|---|---|---|---|---|---|---|
| 1 | 5 | 10 | 5 | 5 | 6 | | |

Bit positions: 31  26 25  16 15  11 10  6 5  0

Returns from an interrupt, JUMPs through EX_UPC, enables USER interrupts for next instruction that executes.

**ERET**
**823**

| 0 | COMM 0 1 0 1 1 | 0 0 0 0 0 0 0 0 0 0 | 0 0 0 0 0 | 0 0 0 0 0 | ERET 0 0 0 0 1 1 | ERET | 1 |
|---|---|---|---|---|---|---|---|
| 1 | 5 | 10 | 5 | 5 | 6 | | |

Bit positions: 31  26 25  16 15  11 10  6 5  0

Returns from an interrupt, JUMPs through EX_PC, enables SYSTEM interrupts for next instruction that executes.

**IHDR**
**823**

| 0 | IHDR 1 1 1 0 1 | rs | rt | unsigned imm | IHDR rt, rs, imm | 1 |
|---|---|---|---|---|---|---|
| 1 | 5 | 5 | 5 | 16 | | |

Bit positions: 31  26 25  21 20  16 15  0

See dynamic network section.

**ILW**
**823**

| s | ILW 0 0 0 0 1 | base | rt | signed offset | ILW rt, offs(base) | 5 2 |
|---|---|---|---|---|---|---|
| 1 | 5 | 5 | 5 | 16 | | |

Bit positions: 31  26 25  21 20  16 15  0

The 16-bit *offset* is sign-extended and added to the contents of *base* to form the effective address. The word at that effective address in the instruction memory is loaded into register *rt*. Last two bits of the effective address must be zero.

Operation:   $Addr \leftarrow ( (offset_{15})^{16} \parallel offset_{15..0}) + [base]$
   $[rt] \leftarrow IMEM[Addr]$

**INTOFF**
**823**

| 0 | COMM 0 1 0 1 1 | 00000 | 0 0 0 0 0 | 00000 | 0 0 0 0 0 | INTOFF 0 0 0 0 0 1 | INTOFF | 1 |
|---|---|---|---|---|---|---|---|---|
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | |

Bit positions: 31  26 25  21 20  16 15  11 10  6 5  0

Clears the interrupt enable bit.

**INTON**
**823**

| 0 | COMM 0 1 0 1 1 | 00000 | 0 0 0 0 0 | 00000 | 0 0 0 0 0 | INTON 0 0 1 0 0 1 | INTON | 1 |
|---|---|---|---|---|---|---|---|---|
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | | |

Bit positions: 31  26 25  21 20  16 15  11 10  6 5  0

Sets the interrupt enable bits.

**ISW**
**823**

| 0 | ISW 0 1 0 0 1 | base | rt | signed offset | ISW rt, offs(base) | 2 |
|---|---|---|---|---|---|---|
| 1 | 5 | 5 | 5 | 16 | | |

Bit positions: 31  27 26 25  21 20  16 15  0

The 16-bit offset is sign-extended and added to the contents of base to form the effective address. The contents of *rt* are stored at the effective address in the instruction memory.

Operation:   $Addr \leftarrow ( (offset_{15})^{16} \parallel offset_{15..0} ) + [base]$
   $IMEM[Addr] \leftarrow [rt]$

Move from event register.

## MFEC

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | rd | 0 0 0 0 0 | MFEC 0 1 0 0 1 0 | |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

MFEC rd,rs — 1

**823**

Move from event counter. Note: MFEC actually captures its value in the RF stage.  [rd] = EC[rs]
Thus a sequence like:
```
lw $0,4($0)                    # cache miss in TV
nop                            # TL
mfec  $4, EC_CACHE_MISS        # Execute -- will not register the cache miss
mfec  $4, EC_CACHE_MISS    # RF       -- will register the cache miss
```
Also, there is one cycle of lag between when the event actually occurs and when
the event counter is actually updated.
 For example, you need a nop between a MTEC and a MFEC to the same register.

## MFSR

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | rd | 0 0 0 0 0 | MFSR 0 1 0 0 0 0 | |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

MFSR rd,rs — 1

**823**

Loads a word from a status register. See "status and control register" table.

Operation:    [rd] = SR[rs]

## MUNLK

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | 00000 | 0 0 0 0 0 | 00000 | 0 0 0 0 1 | INTON 0 0 1 0 0 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

MUNLK — 1

**823/824**

"MDN unlock"
Marks end of mdn locked region. Enables interrupts.

## MLK

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | 00000 | 0 0 0 0 0 | imm5 | 0 0 0 0 1 | INTOFF 0 0 0 0 0 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

MLK <imm5> — 1

**823/824**

"MDN lock"
Signals to hardware and software caching system that the following imm5 * 8 instructions need to be resident
in the cache in order for correct execution to occur. Also, disables interrupts.

## MTEC

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | rs | rt | 0 0 0 0 0 | 0 0 0 0 0 | MTEC 0 1 0 0 1 1 | |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

MTEC rt ,rs — 1

**823**

Move to event counter.  EC[rt] = [rs]

## MTSR

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | rs | rt | 0 0 0 0 0 | 0 0 0 0 0 | MTSR 0 1 0 0 0 1 | |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

MTSR rt ,rs — 1

**823**

Loads a word into a control register, changing the behaviour of the Raw tile. See "status and control register" page.

Operation:    SR[rt] = [rs]

78

**MTSRi**

**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|----|----------|-------|-------|-------|-----|---|
| 0 | MTSRi 1 1 1 0 0 | 0 0 0 0 0 | rt | unsigned immediate | | |
| 1 | 5 | 5 | 5 | 16 | | |

MTSRi rt, imm    | 1 |

Loads a word into a control register, changing the behaviour of the Raw tile. See "status and control register" page.

Operation:    $SR[rt] = 0^{16} \parallel imm$

**OHDR**

**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 |
|----|----------|-------|-------|---|
| 0 | OHDR 0 1 1 1 0 | rs | rt | signed offs |
| 1 | 5 | 5 | 5 | 16 |

OHDR rt, imm(rs)    | 1 |

See dynamic network section.

**OHDRX**

**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 |
|----|----------|-------|-------|---|
| 0 | OHDRX 0 1 1 1 1 | rs | rt | signed offs |
| 5 | 1 | 5 | 5 | 16 |

OHDRX rt, imm(rs)    | 1 |

Like OHDR, but disables interrupts.

**PWRBLK**

**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|----|----------|-------|-------|-------|-----|---|
| 0 | COMM 0 1 0 1 1 | 0 0 0 0 0 | 0 0 0 0 0 | rd | 0 0 0 0 0 | PWRBLK 1 0 0 0 0 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

PWRBLK rd    | 1 |

Stalls in RF stage until: interrupt fires, a word comes into the NESW  buffers of st1/st2 or CGNI or CMNI.
returns bit vector:   G000 0000 0000 0000 0000 0neswNESWMT
G =gdn avail, T = timer, M = cmni, W = cWi ... w = cWi2

**SWLW**

**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 |
|----|----------|-------|-------|---|
| s | SWLW 0 0 1 0 1 | base | rt | signed offset |
| 1 | 5 | 5 | 5 | 16 |

SWLW  rt, offs(base)    | 5 / 1 |
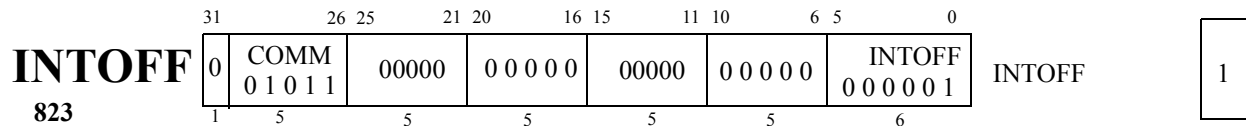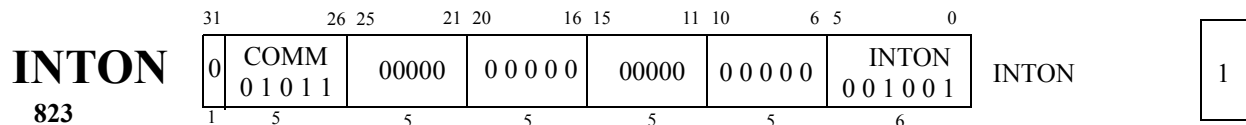
The 16-bit *offset* is sign-extended and added to the contents of *base* to form the effective address.  The word at that effective address in the switch memory is loaded into register *rt*.  Last two bits of the effective address must be zero.
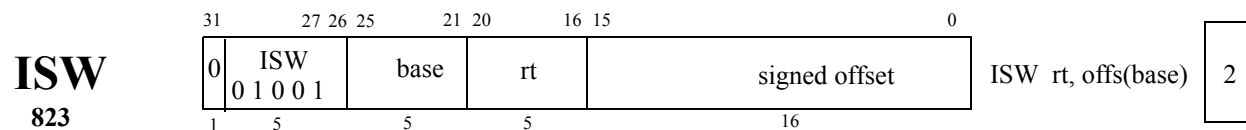
Operation:    Addr ← ( $(offset_{15})^{16} \parallel offset_{15..0}$ ) + [base]
           [rt] ← SWMEM[Addr]
Also occupies one cycle on switch fetch unit.

**SWSW**

**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 0 |
|----|----------|-------|-------|---|
| 0 | SWSW 0 1 1 0 1 | base | rt | signed offset |
| 1 | 5 | 5 | 5 | 16 |

SWSW  rt, offs(base)    | 3 / 1 |
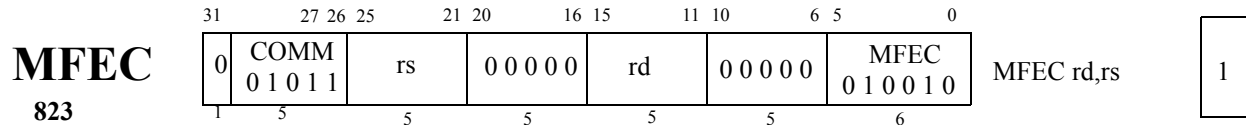
The 16-bit offset is sign-extended and added to the contents of base to form the effective address.  The contents of *rt* are stored at the effective address in the switch memory.

Operation:    Addr ← ( $(offset_{15})^{16} \parallel offset_{15..0}$ ) + [base]
           SWMEM[Addr] ← [rt]
Also occupies one cycle on switch fetch unit.

**UINTON**

**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|----|----------|-------|-------|-------|-----|---|
| 0 | COMM 0 1 0 1 1 | 00000 | 0 0 0 0 0 | 00000 | 0 0 0 0 0 | UINTON 0 0 1 0 1 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

UINTON    | 1 |

**UINTOFF**

**823**

| 31 | 27 26 | 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|---|
| 0 | COMM<br>0 1 0 1 1 | 00000 | 0 0 0 0 0 | 00000 | 0 0 0 0 0 | UINTOFF<br>0 0 0 0 10 | |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 | |

UINTOFF

1

## 8.1 Cache Administrative Instructions - Address Based

| 31 | | 26 25 | | 21 20 | | 16 15 | | 11 10 | | 6 5 | | 0 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**AINV**
**823**

| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | 0 0 0 0 0 | saaaaa | AINV 0 1 1 1 1 0 | AINV rs, saaaa | 1,5 |
|---|---|---|---|---|---|---|---|---|
| 1 | 5 | | | 5 | 5 | 6 | | |

Checks if   rs + (s*16384+aaaa << 5) is in the cache. If it is, it invalidates it. An invalidation takes 4 additional cycles.

**AFL**
**823**

| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | 0 0 0 0 0 | saaaaa | AFL 0 1 1 1 0 0 | AFL rs, saaaa | 1,5,13 |
|---|---|---|---|---|---|---|---|---|
| 1 | 5 | | | 5 | 5 | 6 | | |

Checks if   rs + (s*16384+aaaa << 5)  in the cache. If it is and the line is dirty, the corresponding cache line is written
the dirty bit is cleared. The line is NOT invalidated. If the address is not in the cache, this instruction
takes 1 cycles, if it is in the cache but not dirty, it takes 5 cycles, otherwise, it takes at least 13 cycles.

**AFLINV**
**823**

| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | 0 0 0 0 0 | saaaaa | AFLINV 0 1 1 1 0 1 | AFLINV rs, saaaa | 1,5,13 |
|---|---|---|---|---|---|---|---|---|
| 1 | 5 | | | 5 | 5 | 6 | | |

Checks if   rs + (s*16384+aaaa << 5)  in the cache. If it is and the line is dirty, the corresponding cache line is written
the dirty bit is cleared. The line is then invalidated. If the address is not in the cache, this instruction
takes 1 cycles, if it is in the cache but not dirty, it takes 5 cycles, otherwise, it takes at least 13 cycles.

## 8.2 Cache Administrative Instructions - Tag based

The input parameters of the tag based cache instructions do not specify memory addresses, but rather positions in the tag array of the cache. Bit 15 of the rs value specifies the set (since the cache is two-way set associative.) and bits 14..6 specify the line (there are 512 lines of 32 bytes each), using the convention that bit 1 is the first bit of the word. These instructions should be employed only if you know what you are doing. If you are depending on the state of the cache not changing, make sure to disable interrupts during this period (and later re-enable them!) The additional 5-bit "saaaa" field allows one to specify an immediate offset to the input register. The S bit in the field can be used to toggle the set, the four aaaa bits allow an offset of between 0 and 15 cache lines to be specified. To use the saaaa mode, append an "I" to the end of the mnuemonic, as in the examples.

**TAGFL**
**823**

| 31 | 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | 00000 | saaaaa | TAGFL 0 1 1 0 0 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

TAGFLI rs, saaaa    3

Writes the line located at <set + s, line + a a a a, 5b00000 > pair back to memory. Does not flush if the line is not dirty or not valid. Does not invalidate the line. Sets the MRU bit to the other set.

**TAGLA**
**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | rd | s a a a a | TAGLA 0 1 1 0 1 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

TAGLAI rd, rs, saaaa    3 1

Gets the line address for a <set + s, line + aaaa, 5b00000> pair.

**TAGLV**
**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | rs | 0 0 0 0 0 | rd | s a a a a | TAGLV 0 1 1 0 1 1 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

TAGLVI rd, rs, saaaa    3 1

Returns a 1 if the given {set + s/line + a a a a/5b00000} pair is valid, otherwise a 0.
Not to be issued cycle after LOAD/STORE because of RAW hazards on tag memory.

**TAGSW**
**823**

| 31 | 27 26 25 | 21 20 | 16 15 | 11 10 | 6 5 | 0 |
|---|---|---|---|---|---|---|
| 0 | COMM 0 1 0 1 1 | base | rt | 00000 | s a a a a | TAGSW 0 1 1 0 0 0 |
| 1 | 5 | 5 | 5 | 5 | 5 | 6 |

TAGSWI rt, base, saaaa    1

Writes to the tags. base specifies a { set + s, line + aaaa, 5b00000 } pair. The 19th bit is the valid bit, the other bits are the physical tag. The dirty bit is cleared. The MRU bit is set to the other set. Not to be issued cycle after LOAD/STORE because of WAW hazards on tag memory. The set bit is the 15th bit. (bit numbers start at 1.)
Cache replacement policy: If set0 is invalid, picks that set, else if set1 is invalid, pick that set, else pick the LRU set.

# 9.0 Opcode Map

**OPCODE Map**

This map is for the first six bits of the instruction (the "opcode" field.)
The branch predictor can just look at the top 3 bits of the instruction for a predicted branch.

| | | instruction[28..26] | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
| 111 | | **REGIMM+** | BNEA+ | BNE+ | BEQ+ | BL | BLAL | JNEL+ | JEQL+ |
| 110 | | LB! | LBU! | LH! | LHU! | LW! | SLTIU! | SLTI! | ADDIU! |
| 101 | | RLM | RLMI | RLVM | | RLM! | RLMI! | RLVM! | ! |
| 100 | | **SPECIAL!** | ILW! | ORI! | XORI! | ANDI! | SWLW! | **FPU!** | AUI! |
| 011 | | **REGIMM-** | BNEA- | BNE- | BEQ- | MTSRI | IHDR | JNEL- | JEQL- |
| 010 | | LB | LBU | LH | LHU | LW | SLTIU | SLTI | ADDIU |
| 001 | | SB | ISW | SH | **COM** | SW | SWSW | OHDR | OHDRX |
| 000 | | **SPECIAL** | ILW | ORI | XORI | ANDI | SWLW | **FPU** | AUI |

**SPECIAL Map**

This map is for the last six bits of the instruction when opcode == "SPECIAL".

| | | instruction[2..0] | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
| 000 | | SLL | MAGIC | SRL | SRA | SLLV | | SRLV | SRAV |
| 001 | | JR | JALR | JRHON | JRHOFF | | | | |
| 010 | | MFHI | MTHI | MFLO | MTLO | MFFD | MTFD | | |
| 011 | | MULL | MULLU | DIV | DIVU | | | | |
| 100 | | | ADDU | | SUBU | AND | OR | XOR | NOR |
| 101 | | MULH | MULHU | SLT | SLTU | | | | |
| 110 | | | | | | | | | |
| 111 | | POPC | CLZ | | | | | | |

**FPU Function map**

This opcode map is for the last six bits of the instruction when the opcode field is FPU.

| | | instruction[2..0] | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
| 000 | | ADD.s | SUB.s | MUL.s | DIV.s | | ABS.s | | NEG.s |
| 001 | | | | | | | TRUNC.s | | |
| 010 | | | | | | | | | |
| 011 | | | | | | | | | |
| 100 | | CVT.S | | | | CVT.W | | | |
| 101 | | | | | | | | | |
| 110 | | C.F | C.UN | C.EQ | C.UEQ | C.OLT | C.ULT | C.OLE | C.ULE |
| 111 | | C.SF | C.NGLE | C.SEQ | C.NGL | C.LT | C.NGE | C.LE | C.NGT |

**COM Function map**

This opcode map is for the last six bits of the instruction when the opcode field is COM.

| | | instruction[2..0] | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
| 000 | | DRET | INTOFF | UINTOFF | ERET | | | | |
| 001 | | | INTON | UINTON | | | | | |
| 010 | | MFSR | MTSR | MFEC | MTEC | | | | |
| 011 | | TAGSW | FLUSH | TAGLA | TAGLV | AFL | AFLINV | AINV | |
| 100 | | PWRBLK | | | | | | | |

**REGIMM Map**

This map is for the rt field of the instruction when opcode == "REGIMM."

| | | instruction[18..16] | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
| 00 | | BLTZ | BLEZ | BGEZ | BGTZ | | | | |
| 01 | | | | | | J | | | |
| 10 | | BLTZAL | | BGEZAL | | JLTZL | JLEZL | JGEZL | JGTZL |
| 11 | | | | | | JAL | | | |

# 10.0  Status and Control Registers

| | Status Reg Name | | Purpose |
|---|---|---|---|
| 0 | SW_FREEZE | RW | Switch is frozen ( 1, 0) |
| 1 | SW_BUF1 | R | # els in static switch buffers ( sss NNN EEE SSS WWW III OOOOO) (s,N,E,S,W,I <= 4, O <= 8) |
| 2 | SW_BUF2 | R | # els in switch buffers pair 2 ( sss NNN EEE SSS WWW III OOOOO) (s,N,E,S,W,I <= 4, O <= 8) |
| 3 | MDN_BUF | R | # els in mdn switch buffers   ( NNN EEE SSS WWW III OOOOO) (N,E,S,W,I <= 4, O <= 16) |
| 4 | SW_PC | RW | Switch's PC, byte address (ie low three bits are zero.) When a write occurs to this register, it causes a branch mispredict in the switch. Allow at least three cycles for the corresponding instruction at that switch PC to be executed. |
| 5 | BR_INCR | RW | Increment value for BNEA instruction. (32 bits) |
| 6 | EC_DYN_CFG | RW | Event counter - Dynamic Network events - Config [30:16] = memory network  NNN EEE SSS WWW  PPP [14:0]  = general network   NNN EEE SSS  WWW PPP See event counter section for the meaning of the possible values. |
| 7 | WATCH_VAL | RW | 32 bit Timer count up 1 per cycle |
| 8 | WATCH_MAX | RW | value to reset/interrupt at |
| 9 | WATCH_SET | RW | mode for watchdog counter (Stall Dynamic) (2 bits) |
| 10 | CYCLE_HI | RW | number of cycles from bootup (hi 32 bits) (read first) (writeable for test) |
| 11 | CYCLE_LO | RW | number of cycles from bootup (low 32 bits) (read second, subtract 1) (writeable for test) |
| 12 | EVENT_CFG2 | RW | Event counter configuration, part 2. See event counter section. |
| 13 | GDN_RF_VAL | RW | Dynamic refill value (32 bits) |
| 14 | GDN_REMAIN | RW | GDN_COMPLETE countdown register |
| 16 | GDN_BUF | R | # els in gdn switch buffers  (PPPPP NNN EEE SSS WWW III OOOOO) (PPPPP = GDN_PENDING) (N,E,S,W,I <= 4, O <= 16, P < 32) |
| 17 | GDN_CFG | RW | General dynamic network configuration [31:27] GDN_XMASK - Masks X bits from an address (5 bits) [26:22] GDN_YMASK - Mask Y bits from an address (5 bits) [21:17] GDN_XADJ    - Adjust from local to global X address (5 bits) [16:12] GDN_YADJ    - Adjust from local to global Y address (5 bits) [11:09] GDN_YSHIFT - Gets Y bits from an address (3 bits 0..5) |
| 18 | STORE_METER | RW | STORE_ACK COUNTERS [31:27] PARTNER_Y                - Y location of partner port [26:22] PARTNER_X                - X location of partner port [21]    ENABLE                      - enable store meter-based stalls [10]    DECREMENT_MODE    (reads always 0) [9:5]  COUNT_PARTNER          - counter, num partner accesses left [4:0]  COUNT_NON_PARTNER - counter num non-partner accesses left <br><br> When writing without Bit 10 set, the user must make sure that all store acks have been received, otherwise the counter's value may change spontaneously later. With Bit 10 set, only bits 5 and 0 are used; they specify respectively, whether COUNT_PARTNER or COUNT_NON_PARTNER should be decremented. Bit 10 mode mtsris should be executed at least two cycles after any instruction that can cache miss. |

| | Status Reg Name | | Purpose |
|---|---|---|---|
| 19 | MDN_CFG | RW | Memory Network configuation<br>[31:27] DN_XPOS - Absolute X position of tile in array (5 bits)<br>[26:22] DN_YPOS - Absolute Y position of tile in array (5 bits)<br>[21:17] MDN_XMAX - X Coord of East-Most Tiles  (5 bits)<br>[16:12] MDN_YMAX - Y Coord of South-Most Tiles  (5 bits)<br>[10:09] MDN_XSHIFT - Shift amount (2 bits) (Sim uses addt'l bit 11)<br>[07:06] MDN_YSHIFT - Shift amount (2 bits) (Sim uses addt'l bit 08)<br>[00:00] MDN_EXTEND - Use all four edges of the chip |
| 20 | EX_PC | RW | PC where system-level exception occurred |
| 21 | EX_UPC | RW | PC where user exception occurred. |
| 22 | FPSR | RW | Floating Point  Status Register (E V Z O U I)<br>(Unimplemented, Invalid, Div by Zero, Overflow,<br> underflow, Inexact Operation)<br>These bits are sticky, ie a floating point operation can only set the bits, never clear. However, the user can both set and clear all of the bits.<br>These flags are set the cycle after the floating point instruction has finished its result, i.e. to get a valid value, you need to insert 3 nops to read the correct value. |
| 23 | EVENT_BITS | R | [15:0] the list of events that has triggered |
| 24 | EX_BITS | R | [6:3] =  EVENT_COUNTER / GDN_AVAIL / TIMER / EXTERNAL<br>[2:0] =  TRACE / GDN_COMPLETE / GDN_REFILL<br><br>"1" indicates a request for a given interrupt occured at some point<br>Other bits are implicitly reset.<br><br>[31:30] = USER / SYSTEM<br>Indicates whether exceptions are enabled.<br>SYSTEM = 0 --> all interrupts disabled<br>USER = 0 --> user interrupts disabled |
| 25 | EX_MASK | RW | [6:3] =  EVENT_COUNTER / GDN_AVAIL / TIMER / EXTERNAL<br>[2:0] =  TRACE / GDN_COMPLETE / GDN_REFILL<br><br>A "1" indicates that that exception is enabled. |
| 26 | EVENT_CFG | RW | [0]    1,0 -->  single / global  instr mode (single mode uses bits 15:1)<br>[15:1] PC to profile (omit low two bits) for single mode<br>each bit enables an event counter (C = cache, B = branch)<br>[31:16 ] Enables for all of the events |

| | Status Reg Name | | Purpose |
|---|---|---|---|
| 27 | POWER_CFG | RW | bit 0 = disable power saving in comparator<br>bit 1 = disable power saving in ALU<br>bit 2 = disable power saving in FPU<br>bit 3 = disable power saving in multiplier<br>bit 4 = disable power saving in divider<br>bit 5 = disable power saving in data cache<br>bit 6 = enable power saving in instr memory<br>bit 7 = enable power saving in data memory<br>bit 8 = enable power saving in switch memory<br>bit 9 =   disable pwrblk wake up after timer wakeup<br>bit 10 = disable pwrblk wake up after external wakeup<br>(the following can be set and cleared by both the chip and mtsr/mtsri...)<br>bit 11 =  timer wakeup pending on return to pwrblk instruction<br>bit 12 = external wakeup pending on return to pwrblk instruction<br>at reset, set to 0. |
| 28 | TN_CFG | RW | |
| 29 | TN_DONE | W | |
| 30 | TN_PASS | W | |
| 31 | TN_FAIL | W | |

These status and control registers are accessed by the MTSR and MFSR instructions, and by the UINTOFF, UINTON instructions.

**Status Registers - Bank 2**  (Raw 824 only)

| | Status Reg Name | | Purpose |
|---|---|---|---|
| 0 | HW_ICACHE | R | bit 0 = hw icaching enabled |
| 1 | | | |

# 11.0 Exception Vectors

| | Vector Name | Imem Addr >> 4 | Purpose |
|---|---|---|---|
| | | | |
| 0 | VEC_GDN_REFILL | 0 | Dynamic Refill Exception |
| 1 | VEC_GDN_COMPLETE | 1 | GDN send is complete |
| 2 | VEC_TRACE | 2 | Trace interrupt |
| 3 | VEC_EXTERN | 3 | External Exception, over MDN |
| 4 | VEC_TIMER | 4 | Timer Exception |
| 5 | VEC_GDN_AVAIL | 5 | Data avail on GDN |
| 6 | VEC_EVENT_COUNTERS | 6 | Event Counter Interrupt |
| 7 | | | |

The exceptions vectors are stored in IMEM. One of the main concerns with storing vectors in unprotected memory is that they can be easily overwritten by data accesses, resulting in an unstable machine. Since we are a Princeton architecture, however, the separation of the instruction memory from the data memory affords us a small amount of protection. Another alternative is use a register file for this purposes. Given the number of vectors we support, this is not so exciting. The ram requirements of these vectors is 4 words per vector: this lets us put a 32bit value into a register and jump somewhere (useful if the vectors are going to be cached by the i-caching system).

```
sw save($gp), $3
lw $3, vec($gp)
jmp icache.vec
```

# 12.0 Switch Processor

The switch processor is responsible for routing values between the Raw tiles. One might view it as a VLIW processor which can execute a tremendous number of moves in parallel. The assembly language of the switch is designed to minimize the knowledge of the switch microarchitecture needed to program it while maintaining the full functionality.

The switch processor has three structural components:

1. A 1 read port, 1 write port, 4-element register file.
2. A pair of crossbars, which is responsible for routing values to neighboring switches.
3. A sequencer which executes a very basic instruction set.

A switch instruction consists of a processor instruction and a list of routes for the two crossbars. All combinations of processor instructions and routes are allowed subject to the following restrictions:

1. The source of a processor instruction can be a register or a switch port but the destination must be a register.
2. The source of a route can be register or a switch port but the destination must always be a switch port.
3. Two values can not be routed to the same location.
4. If there are multiple reads to the register file, they must use the same register number. This is because there is only one read port.
5. routes between [cN, cS, cE, cW, cst1] and [cN2, cS2, cE2, cW2, cst2] are forbidden.
To switch domains, one should route the word through swo2 or swo1.
It will appear on swi2 or swi1 on the next cycle.

For instance,

```
MOVE     $3,$2        ROUTE   $2->$csti, $2->$cNo, $2->$cSo, $cSi->$cEo
MOVE     $3,$csto     ROUTE   $2->$csti, $2->$cNo, $2->$cSo, $cSi->$cEo
```

are legal because they read exactly one register (r2) and write one register (r3).

```
JAL      $3, myAddr   ROUTE $csto->$2
```

is illegal because the ROUTE instruction is trying to use r2 as a destination.

```
JALR     $2,$3        ROUTE $2->$csti
```

is illegal because two different reads are being initiated to the register file (r2,r3).

```
JALR     $2,$3        ROUTE $2->$csti, $cNi->$csti
```

is illegal because two different writes are occurring to the same port.

# 13.0 Switch Processor Instruction Set

## BEQZ/JEQZ \<rp\>, ofs16                    beqz $cEi, myRelTarget

**823**                    R = 1/0, pr = 0 or 1,  op = 7, imm = ( ofs16 >> 3), rego = \<rp\>

## BEQZD/JEQZD \<r1\>,\<r2\> ofs16                    beqzd $1,$2 myRelTarg

**823**                    R = 1/0, pr = 0 or 1,  op = E, imm = ( ofs16 >> 3), rego = \<regi\> rsrc = \<r1\>, rdst = \<r2\>

Branch equal to zero, then decrement. MUST BE REGISTER-REGISTER.
  r2_old = r2;
  r1 = r2 - 1;
  if (r2_old == 0) goto ofs16

## BLEZ/JLEZ \<rp\>, ofs16                    blez $cNi, myRelTarget

**823**                    R = 1/0, pr = 0 or 1, op = C, imm = ( ofs16 >> 3), rego = \<rp\>

## BLTZ/JLTZ \<rp\>, ofs16                    bltz $cNi, myRelTarget

**823**                    R = 1/0, pr = 0 or 1, op = 1, imm = ( ofs16 >> 3), rego = \<rp\>

## BNEZ/JNEZ \<rp\>, ofs16                    bnez $2, myRelTarget

**823**                    R = 1/0, pr = 0 or 1, op = 2, imm = ( ofs16 >> 3), rego = \<rp\>

## BNEZD/JNEZD \<r1\>,\<r2\> ofs16                    bnezd $1,$2 myRelTarg

**823**                    R = 1/0, pr = 0 or 1, op = D, imm = ( ofs16 >> 3), rego = \<regi\> rsrc = \<r1\>, rdst = \<r2\>

Branch not equal to zero, then decrement. MUST BE REGISTER-REGISTER.
  r2_old = r2;
  r1 = r2 - 1;
  if (r2_old != 0) goto ofs16

## BGEZ/JGEZ \<rp\>, ofs16                    bgez $cSti, myRelTarget

**823**                    R = 1/0, pr = 0or 1, op = 3, imm = ( ofs16 >> 3), rego = \<rp\>

## BGTZ/JGTZ \<rp\>, ofs16                    bgtz $cNi, myRelTarget

**823**                    R = 1/0, pr = 0 or 1, op = 6, imm = ( ofs16 >> 3), rego = \<rp\>

## DEBUG imm13

debug imm13

**823**

R = 0, pr = 0, op =15, imm = (imm13) , rego = \<none\>

## BAL/JAL  \<rd\>, ofs16

jal $2, myAbsTarget

**823**

R = 1/0, pr = 1, op = 4, imm = ( ofs16 >> 3), rego = "none", rdst = \<rd\>

## JALR \<rd\>, \<rp\>

jalr $2, $cWi

**823**

R = 0, pr = 0, op = 0 (other), ext_op  = 3, rego = \<rp\>, rdst = \<rd\>

## B/J  ofs16

j myAbsTarget

**823**

R = 1/0, pr = 1, op = 5, imm = (ofs16 >> 3), rego = "none"

## JR \<rp\>

jr $cWi

**823**

R = 0, pr = 0, op =0 (other), ext_op = 2, rego = \<rp\>

## MOVE  \<rd\>, \<rp\>

move $1, $cNi

**823**

R = 0, pr = 0, op =0 (other), ext_op = 1, rego = \<rp\>, rdst = \<rd\>

## NOP

nop

**823**

R = 0, pr = 0, op =0 (other), ext_op = 0, rego = \<none\>

Switch instruction formats

**Default instruction format**

register number, if a register is read

route instruction sources

| pr | op | R | 0 | imm | rdst | rsrc | w | rego | cNo | cEo | cSo | cWo | csti | swo2 | cNo2 | cEo2 | cSo2 | cWo2 | csti2 | swo1 |
|----|----|---|---|-----|------|------|---|------|-----|-----|-----|-----|------|------|------|------|------|------|-------|------|
| 1 | 4 | 1 | 1 | 13 | 2 | 2 | 1 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |

bit positions: 63  59 58 57  44  42  39  36  33  30  27  24  21  18  15  12  9  6  3  0

pr = (whether we predict a branch/jump or not)
R = (relative bit) - add PC to imm (for branches)
w = (which bit) whether rego is from 1st or 2nd xbar)

op =7, imm = (0), rego = <rp>, rdst = <rd>

**ExOp instruction format (and conceptual wiring of switch crossbars)**

| 0 | op | 00 | 9'b0 | Exop | ... | rego | cNo | cEo | cSo | cWo | csti | swo2 | cNo2 | cEo2 | cSo2 | cWo2 | csti2 | swo1 |
|---|----|----|------|------|-----|------|-----|-----|-----|-----|------|------|------|------|------|------|-------|------|
| 1 | 4 | 2 | | 4 | | | | | | | | | | | | | | |

bit positions: 63  59  57  48  44  27  24

Crossbar Data from buffers

regi
cNi
cEi
cSi
cWi
csto

Crossbar control

swi1

regi
cNi2
cEi2
cSi2
cWi2
csto

Crossbar control

swi2

Crossbar Data from buffers

cNo
cEo
cSo
cWo
csti
rego

cNo2
cEo2
cSo2
cWo2
csti2
rego

swo1

swo2

Crossbar output (to neighboring tile or proc input port)

# Opcode Map (bits 62..58)

| | instruction[60..58] | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 000<br>(R=0) | 010<br>(R=0) | 100<br>(R=0) | 110<br>(R=0) | 001<br>(R=1) | 011<br>(R=1) | 101<br>(R=1) | 111<br>(R=1) |
| 00 | ExOp | JLTZ | JNEZ | JGEZ | | BLTZ | BNEZ | BGEZ |
| 01 | JAL | J | JGTZ | JEQZ | BAL | B | BGTZ | BEQZ |
| 10 | | | | | | | | |
| 11 | JLEZ | JNEZD | JEQZD | DEBUG | BLEZ | BNEZD | BEQZD | |

Note: for this implementation, the BXX instructions are automatically rewritten into the equivalent JUMP versions on swsw instructions. The switch RTL code (except in implementation of the swlw/swsw instructions) does not need to treat the branches and jumps separately.

# ExOp Map (bits 47..44)

| | instruction[45..44] | | | |
|---|---|---|---|---|
| | 00 | 01 | 10 | 11 |
| 00 | NOP | MOVE | JR | JALR |
| 01 | | | | |
| 10 | | | | |
| 11 | | | | |

# Port Name Map (primary static switch)

| Port | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| | none | csto | cWi | cSi | cEi | cNi | swi1 | regi |

# Port Name Map (secondary static switch)

| Port | 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| | none | csto | cWi2 | cSi2 | cEi2 | cNi2 | swi2 | regi |

# 14.0  Event Counter Support

The event counters are built as an array of 16 c_trigger modules. Each c_trigger has a 32 bit counter. These counters count down every time a particular event occurs. When the counter transitions from 0 to -1, it will assert a line (the "trigger") which will hold steady until the user writes new value into the counter. These triggers are visible in the EVENT_BITS register, and are OR'd together to form the EX_BITS EVENT_COUNTER bit, which can cause an interrupt. The counter latches the PC (without the low zero bits) of the instruction that caused the event into bits [31:16] of the counter (use the rlm instruction to extract them!). The counter will continue to count down regardless of the setting of the trigger. Because the PC is stored in the high bits, there is a window of time in which subsequent events will not corrupt the captured pc. Note that if the event is not instruction related, the setting of the PC is undefined. The event counters can be both read and written by the user. There is typically a one cycle delay between when an event occurs and when a mfec instruction will observe it, there is also a number of cycles (2) before an event trigger interrupt will fire.

**Table 3: Event counters**

headings:  # = counter number, s = stage associated with this event.
(M = memory, E = execute, F = fpu, S = switch, R = register fetch, @ = ignores single instruction mode)

| # | s | Function | Notes | EVENT_CFG2 |
|---|---|----------|-------|------------|
| 0 | @ | Cycles | so handler can bound acquisition time | bit 25 = 0 |
| 0 | F | Write Over Read | detects when a resident cache line is marked dirty by a sw to an odd address for the first time. note: if the sw instruction is preceded by a lw/sw/flush instruction, this mechanism does not have the bandwidth to verify the previous state of the bits. It will conservatively count it as an event. | bit 25 = 1 |
| 1 | M | Cache Writebacks | includes flushes | |
| 2 | M | Cache Fills | | |
| 3 | M | Cache Stall Cycles | Total number of cycles that backend of pipeline is frozen by cache state machine. Includes writeback and fill time, as well as timing stolen by non-dirty flush instruction executions. | |
| 4 | E | Cache Miss Ops | # of flush, lw, sw instructions | Bit 0 = 0 |
| 4 | E | FPU Ops | number of instructions issued to FPU (includes *.s, *.w) | Bit 0 = 1 |
| 5 | E | Possible Mispredicts | Conditional Jumps and Branches, ERET, DRET, JR, JALR | Bit 1 = 0 |
| 5 | E | Possible Mispredicts | Possible Mispredicts due to wrong SBIT (i.e, only conditional jumps and branches) | Bit 1 = 1 |
| 6 | E | Actual Mispredicts | Mispredicts | Bit 2 = 0 |

**Table 3: Event counters**
headings:  # = counter number, s = stage associated with this event.
(M = memory, E = execute, F = fpu, S = switch, R = register fetch, @ = ignores single instruction mode)

| # | s | Function | Notes | EVENT_CFG2 |
|---|---|----------|-------|------------|
| 6 | E | Actual Mispredicts | Mispredicts due to wrong SBIT | Bit 2 = 1 |
| 7 | @ | Switch Stalls | On switch  (PC captured is switch pc) | |
| 8 | @ | Possible Mispredicts | On switch (PC captured is switch pc) | |
| 9 | @ | Actual Mispredicts | On switch  (PC captured is switch pc) | |
| 10 | @ | Pseudo Rand LFSR | X_next =  (X >> 1) \| (xor(X[31,30,10,0]) << 31) note: sampling this more than once per 32 cycles will produce highly correlated numbers. | |
| 11 | R | Func Unit Stalls | Stalls due to bypassing (e.g., the output of the preceding instruction is not available yet) or because of interlocks on the fp/int dividers. | Bit 3 = 1/0 |
| 11 | @ | GP | General Network Processor Port Counting | Bit 4 = 1/0 |
| 11 | @ | MP | Memory Network Processor Port Counting | Bit 5 = 1/0 |
| 11 | @ | Instrs Issued1 | Number of instructions that enter Execute stage. | Bit 23 = 1/0 |
| 12 | R | Non-cache stalls | # stalls not due to cache misses (includes isw/ilw stalls, if trigger fires on isw/ilw, the pc will be the pc of the instruction in the rf stage, rather than the isw/ilw instruction) | Bit 6 = 1/0 |
| 12 | @ | GW | General Network West Port Counting | Bit 7 = 1/0 |
| 12 | @ | MW | Memory Network West Port Counting | Bit 8 = 1/0 |
| 13 | R | ISW/ILW issued | # ilw/isw instructions issued | Bit 9 = 1/0 |
| 13 | @ | GS | General Network South Port Counting | Bit 10 = 1/0 |
| 13 | @ | MS | Memory Network South Port Counting | Bit 11 = 1/0 |
| 13 | @ | Instrs Issued2 | Number of instructions that enter execute stage | Bit 24 = 1/0 |
| 14 | R | CSTO stalls | Blocked on CSTO | Bit 12 = 1/0 |
| 14 | R | CGNO stalls | Blocked on CGNO | Bit 13 = 1/0 |
| 14 | R | CMNO stalls | Blocked on CMNO | Bit 14 = 1/0 |
| 14 | @ | GE | General Network East Port Counting | Bit 15 = 1/0 |
| 14 | @ | ME | Memory Network East Port Counting | Bit 16 = 1/0 |
| 15 | R | CSTI stalls | No data available on CSTI | Bit 17 = 1/0 |
| 15 | R | CSTI2 stalls | No data available on CSTI2 | Bit 18 = 1/0 |
| 15 | R | CGNI stalls | No data available on CGNI | Bit 19 = 1/0 |
| 15 | R | CMNI stalls | No data available on CMNI | Bit 20 = 1/0 |

**Table 3: Event counters**

headings:  # = counter number, s = stage associated with this event.
(M = memory, E = execute, F = fpu, S = switch, R = register fetch, @ = ignores single instruction mode)

| # | s | Function | Notes | EVENT_CFG2 |
|---|---|----------|-------|------------|
| 15 | @ | GN | General Network North Port Counting | Bit 21 = 1/0 |
| 15 | @ | MN | Memory Network North Port Counting | Bit 22 = 1/0 |

The low bits of the EVENT_CFG register allow the user to only count events that occur on a single instruction at a particular main processor PC instead of across all PCs. For single instruction mode, set the low bit of EVENT_CFG, and place the PC to sample (excluding the low two zero bits) into bits [15:1]. In cases where the event does not have an associated main processor PC (marked with the "@" in the table), the EVENT_CFG single instruction mode setting is ignored. The high bits of EVENT_CFG selectively enable counting on a per event basis, but do not suppress existing triggers.

The EVENT_CFG2 SPR allows the user to configure the events that a particular event counter counts. The Event Counter table shows the setting of these bits. In the cases where two events are listed in the same box, the counter is configured through EVENT_CFG2 to count one or the other. In cases where multiple events share the same event counter but are not listed in the same box, the corresponding bits EVENT_CFG2 select the subset of events that will trigger an increment of the event counter. In some cases, there are nonsensical combinations (say GE and csto stalls).

The meaning of the GN GE GS GW GP MN ME MS MW MP events are configured by the EC_DYN_CFG SPR. Each event corresponds to a network N (G = general, M = memory) and a direction D (N=north, E=east, ...) The meanings are as follows:

| val | ec_dyn_cfg fields - Description |
|-----|--------------------------------|
| 0 | # cycles the output port D wanted to transmit but could not because of blockage on a neighboring tile |
| 1 | # of words transmitted from input port D to output port P |
| 2 | # of words transmitted from input port D to output port W |
| 3 | # of words transmitted from input port D to output port S |
| 4 | # of words transmitted from input port D to output port E |
| 5 | # of words transmitted from input port D to output port N |
| 6 | # of words transmitted from input port D |
| 7 | # cycles that input port D had data to transmit was not able to |

# Administrative Procedures

**Warning: This stuff is somewhat outdated.**

## Interrupt masking

To be discussed at a later date.

## Processor thread switch  (does not include switch processor)

EPC must be saved off and new values put in place. A ERET
will cause an atomic return and interrupt enable.

```
mfsr      $29, EPC
sw        $29, EPC_VAL($0)
lw        $29, NEW_EPC_VAL($0)
mtsr      EPC, $29
lw        $29, OLD_R29($0)
eret                                    # return and enable interrupt bits
```

## Freezing the Switch

The switch may be frozen and unfrozen at will by the processor.
This is useful for a variety of purposes. When the switch is frozen,
it ceases to sequence the PC, and no routes are performed.

## Reading or Write the Switch's PC

The switch processor's PC can be written at any time. However, it is often the case that one will
want to freeze the switch before doing so.

```
mtsr SW_PC, $2                     # set new switch PC to $2
```

The PC of the switch may be read at any time, in any order. However, we imagine that this opera-
tion will be most useful when the switch is frozen.

```
mtsri FREEZE, 1# freeze the switch
mfsr $2, SW_PC# get PC
mtsri FREEZE, 0                         # unfreeze the switch
```

## Reading or Writing the Processors's IMEM

This will stall the processor for one cycle per access. The read or write will cause
the processor to stall for one cycle. Addresses are multiples of 4. Any low bits will
be ignored.

```
ilw $3, 0x160($2)# load a value from the proc imem
```

```
isw $5, 0x168($2)# store a value into the proc imem
```

## Reading or Writing the Switch's IMEM

The switch can be frozen or unfrozen. The read or write will cause
the switch processor to stall for one cycle. Addresses are multiples of 4. Any low bits will
be ignored. Note that instructions must be aligned to 8 byte boundaries. The switch imem
instructions are 64 bits. They are store in XINU form; i.e., the word containing the opcode is
stored last, at an address ending in xxx100b.

When looking at a hex-dump of the instructions, the last byte of the instruction contains the
opcode.

```
swlw    $3, 0x160($2)                    # load a value from the switch imem
swlw    $4, 0x164($2)                    # load a value from the switch imem
                                         # (this word contains the opcode)

swsw    $5, 0x168($2)                    # store a value into the switch imem
swsw    $6, 0x168($2)                    # store a value into the switch imem
```

**Determining how many elements are in a given switch buffer**

At any point in time, it is useful to determine how many elements are waiting in the buffer of a
given switch. There are two SRs used for this purpose, SWBUF1, which is for the first set of input
an output ports, and SWBUF2, which is for double-bandwidth switch implementations. The for-
mat of these status words is as follows:

```
# to discover how many elements are waiting in csto queue

mfsr    $2, SWBUF1                    # load buffer element counts
andi $2, $2, 0x1F                     # get $csto count
```
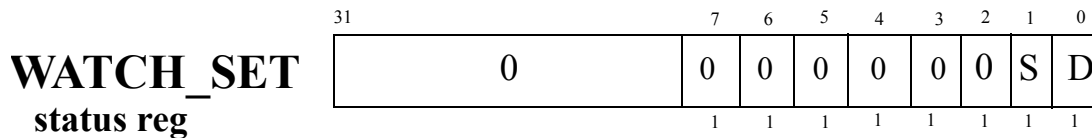
## Using the watchdog timer

The watchdog timer can be used to monitor the dynamic network and determine if a deadlock condition may have occurred. WATCH_VAL is the current value of the timer, incremented every cycle, regardless of what is going on in the processor.
WATCH_MAX is the value of the timer which will cause a watch event to occur:

There are several bits in WATCH_SET which determine when WATCH_VAL is reset and if an interrupt fires (by default, these values are all zero):

| 31 | | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|----|---|---|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | S | D |

**WATCH_SET**
**status reg**

bit widths: 1 1 1 1 1 1 1 1

| | Bit | Name | effect |
|---|-----|------|--------|
| | 0 | DYN_MOVE | reset WATCH_VAL when a data element is removed from dynamic network (or refill buffer), or if no data is available on dynamic network ? |
| | 1 | NOT_STALLED | reset WATCH_VAL if the processor was not stalled ? |
| | 2 | | |
| | 3 | | |
| | 4 | | |
| | 5 | | |

```
# code to enable watch dog timer for dynamic network deadlock

mtsr WATCH_MAX, 0xFFFF          # 65000 cycles
mtsr WATCH_VAL, 0x0            # start at zero
mtsr WATCH_SET, 0x3           # interrupt on stall and no
                             # dynamic network activity
jr 31

# watchdog timer interrupt handler
# pulls as much data off of the dynamic network as
# possible, sets the DYNREFILL bit and then
# continues


sw      $2, SAVE1($0)          # save a reg
                             # (not needed
                             #  if reserved regs for handlers)
sw      $3, SAVE2($1)          # save a reg
lw      $2, HEAD($0)           # get the head index
```

```
lw      $3, TAIL($0)              # get the tail index
add     $3, $3,1
and     $3, $3, 0x1F              # thirty-one element queue
beq     $2, $3, dead             # if queue full, we need some serious work
blop:
lw      $2, TAIL($0)
sw      $3, TAIL($0)              # save off new tail value
sw      $cgni, $2(BUFFER)        # pull something out of the network
mfsr    $2, D_AVAIL              # stuff on the dynamic network still?
beqz    $2, out                  # nothing on, let's progress
lw      $2, SAVE1($0)            # restore register (delay slot)

# otherwise, let's try to save more
move    $2, $3
add     $3, $2, 1
and     $3, $3, 0x1F             # thirty-one el queue
bne     $2, $3, blop             # if queue not full, we process another
lw      $2, SAVE1($0)            # restore register (delay slot)

out:
mtsr    DYNREFILL, 1             # enable dynamic refill
dret
lw      $3, SAVE2($1)            # restore register
```

## Setting or Reading an Exception Vector

Exception vectors are instructions located at predefined locations in memory to which the processor should branch when an exceptional case occurs. They are typically branches followed by delay slots. See the Exceptions sections for more information on this.

```
ILW     $2, ExceptionVectorAddress($0)    # save old interrupt instruction
ISW     $3, ExceptionVectorAddress(40)    # set new interrupt instruction
```

## Using Dynamic Refill (DYNREFILL/EX_DYN_REF/DR_VAL)

Dynamic refill mode allows us to virtualize the dynamic network input port. This functionality is useful if we find ourselves attempt to perform deadlock recovery on the dynamic network. When DYNREFILL is enabled, a dynamic read will take its value from the "DR_VAL" register and cause a EX_DYN_REF immediately after. The deadlock countdown timer (if enabled) will be reset as with an dynamic read. This will give the runtime system the opportunity to either insert another value into the refill register, or to turn off the DYNREFILL mode.

```
# enable dynamic refill

mtsri   DYNREFILL, 1                # enable dynamic refill
mtsr    DR_VAL, $2                  # set refill value
dret                                # return to user

# drefill exception vector
# removes an element off of a circular fifo and places it in DR_VAL
# if the circular fifo is empty, disable DYNREFILL
# if (HEAD==TAIL), fifo is empty
# if ((TAIL + 1) % size == HEAD), fifo is full

sw      $2, SAVE1($0)               # save a reg (not needed if
                                    #         reserved regs for handlers)
sw      $3, SAVE2($1)               # save a reg
lw      $2, HEAD($0)                # get the head index
lw      $3, $2(BUFFER)              # get next word
mtsr    DR_VAL, $3                  # set DR_VAL
add     $2, $2, 1                   # increment head index
and     $2, $2, 0xF                 # buffer is 32 (31 effective) entries big
lw      $3, TAIL($0)                # load tail
sw      $2, HEAD($0)                # save new head
bne     $2,$3, out                  # if head == tail buffer is empty
lw      $2, SAVE1($0)               # restore register (delay slot)
mstri   DYNREFILL, 0                # buffer is empty, turn off DYNREFILL

out:
dret
lw      $3, SAVE2($1)               # restore register
```

# Schedules

| | Design | Sim | RTL |
|---|---|---|---|
| **Initial Run** | | | |
| Muxing / Pin Problems | 3 jm | 2 jm | 2 sl |
| Caching - Integration | 5 mt | 4 mt | 6 jk |
| Dynamic Network | 3 mt | 2 mt | 8 jk |
| SPRs/INT | 1 mt | 1 mt | 7 jk |
| Misc | | 4 mt | |
| **21 days = 5 weeks** | | | Mar 10 |

MBT: 20    JK: 21    SL: 2   JM: 5

| **Pipelining, Timing** | Design | Sim | RTL |
|---|---|---|---|
| Planning | 6 mt/jk | | |
| FPU | 5 alb | - | 5 alb |
| Static Net | 5 mt | 5 jm | 9 sl |
| Mustang | 0 | 0 | 15 alb |
| Control | 5 mt | 8 mt | 14 jk |
| Branch | 5 mt | 2 mt | 5 jk |
| Timing Iteration | 30 al/jk/mt | | |
| **65 days = 13 weeks** | | | May 1 |

MBT: 31+30    JK: 25+30    ALB:  25+30    SL:  9   JM: 5

Resources:  Anant, Andras, Albert, Ben, JasonM, Saman,
            Elliot, Mark, JasonK, Walt, Sam, Omar, Michael, Matt, Kevin

These quantities are in work-weeks, not wall-clock weeks. They
do not include times for FIXES, SETBACKS, CLASSES,
CONFERENCES, PAPERS, VACATIONS, or IBM INTERACTIONS.

| Design Validation | |
|---|---|
| RawCC | (1 month) Walt |
| SUDS | (1 month) Matt |
| I-Caching | (1 month) Jason Miller |
| BenOS | (1 month) Ben |
| ElliotCC | (1 month) Elliot |
| **5 weeks** | |

| Design Verification | |
|---|---|
| Testing Infrastructure<br>- Torture<br>- Test integration framework<br>- Basher infrastructure<br>- Testing Strategy | 8 weeks |
| White box testers | All    8 weeks |
| Bashers | All    20 weeks |
| **20 Weeks** | |

| Design Finalization | |
|---|---|
| Pins PLLs, Clocks, IBM tools, Design Rule Fixes | 7 weeks  (Albert, Omar) |
| Boards | JasonM, Kevin |
| **7 weeks** | |

Hardware Path: 5+13+7 weeks = 25 weeks
Testing Path: 5+5+20 weeks = 30 weeks

Massachusetts Institute of Technology
Laboratory of Computer Science

# RAW Prototype Chip
# Verification Strategy

Michael Taylor

Version 1.0

# 1.0  Foreword

This subdocument describes the verification strategy of the Raw processor. Verification seems to be an uncertain science. This document is a first stab at trying to solidify our plans. I recommend that readers of this document also read "www.cs.princeton.edu/~doug/papers/cmu.ps," which is the reference that I have used for this document.

# 2.0  Introduction

The thesis of  "Large-Scale Hardware Simulation: Modeling and Verification Strategies" is that the purpose of verification is to find bugs. This may seem like a simple premise, but an alternative view might presuppose that the purpose of verification is to attain assurance that the design is sound.
If we believe the thesis of that paper, then our goal should be to find bugs in the design rather than to seek assurance. That paper proposes that a more dynamic approach to bug finding should be applied: explore approaches that maximize payoff, and change those approaches as they yield lesser results. This implies that a basic idealogy and initial plan for bug finding should be applied, but as more information about the types and frequencies of bugs becomes available, we should adjust our plan. This amounts more or less to an open season on bugs, where the hunters rely on previous experience and their own cleverness to bring the bugs in, rather than creating a grid across the forest and checking each point.

# 3.0  Targetted (unit) Tests

Initial, we should have unit targetted tests for every major component of the chip. This includes:

Integer Datapath,
Floating Point Datapath,
Switch Processor,
Static Networks,
Dynamic Networks,
Muxing Logic,
Data Cache (including misses),
Fetch Unit (including mispredicts),
Interrupts and interrupt logic,
Status Registers,
Stall Logic,
Bypass Logic,
Off-chip I/O support

Some of these tests will require coordination of events between different systems in order to test thoroughly. These tests will be hand-written by raw group members, and in such a way that we can chain them together and test interaction conditions. These tests will take the form of programs

rather than as test-benchs for the verilog so that the results of the simulator and the rtl can be cross-verified.

## 4.0  Unit test Methodology

To facilitate the creation of tests that operate in a common framework, we need to come up with a set of primitives which test writers can use in their test code.

Certainly, two obvious ones are:

PASS();
FAIL();

It might also be productive to have some other hints that describe what resources this test uses so that it can be automatically mixed with other tests.

## 5.0  Small Scale Test

There might also be a use in collecting together a set of smallish codes which exercise the abilities of the chip but are still easy to debug. This allows us to rely on the compiler and external resources for creating useful tests.

## 6.0  Random Tests

These are tests that are automatically generated by TORTURE, krste's test generator program.

## 7.0  Bashing / Daemon testing

These tests are intended to target interactions of systems. They combined together tests with interrupt handlers and other extra-ordinary behaviour. One idea is to have a stall signal in the tiles which generates random stalls for random periods of time.

## 8.0  Code reviews

Careful reading of verilog code by relevant parties.