

**Precision Timed Machines**

by

Isaac Liu

A dissertation submitted in partial satisfaction of the  
requirements for the degree of  
Doctor of Philosophy

in

Electrical Engineering and Computer Science

in the

GRADUATE DIVISION  
of the  
UNIVERSITY OF CALIFORNIA, BERKELEY

Committee in charge:  
Professor Edward A. Lee, Chair  
Professor John Wawrzynek  
Professor Alice Agogino

Spring 2012

# **Precision Timed Machines**

Copyright 2012  
by  
Isaac Liu

## **Abstract**

Precision Timed Machines

by

Isaac Liu

Doctor of Philosophy in Electrical Engineering and Computer Science

University of California, Berkeley

Professor Edward A. Lee, Chair

This is my abstract

To my wife Emily Cheung, my parents Char-Shine Liu and Shu-Jen Liu, and everyone else whom I've had the privilege of running into for the first twenty-seven years of my life.

# Contents

<b>List of Figures</b>	<b>iii</b>
<b>List of Tables</b>	<b>iv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Intro Section Header 2 . . . . .	1
<b>2 Precision Timed Machine</b>	<b>2</b>
2.1 Thread-Interleaved Pipelines . . . . .	8
2.2 Memory System . . . . .	10
<b>3 Programming Models</b>	<b>11</b>
3.1 PRET Programming model Section Header . . . . .	11
3.2 Pret Programming model Section Header 2 . . . . .	11
<b>4 Related Work</b>	<b>12</b>
4.1 Architectural Modifications . . . . .	12
4.2 Related Section Header 2 . . . . .	12
<b>5 Implementation of PTARM</b>	<b>13</b>
5.1 Implementation . . . . .	13
5.1.1 PTARM Simulator . . . . .	13
5.1.2 PTARM VHDL Softcore . . . . .	13
5.1.3 Worst Case Execution Time Analysis . . . . .	13
<b>6 Applications</b>	<b>14</b>
6.1 Eliminating Side-Channel-Attacks . . . . .	14
6.2 Real Time 1D Computational Fluid Dynamics Simulator . . . . .	15
<b>7 Conclusion and Future work</b>	<b>17</b>
7.1 Summary of Results . . . . .	17
7.2 Future Work . . . . .	17
<b>Bibliography</b>	<b>18</b>

# List of Figures

1.1	Image Placeholder . . . . .	1
2.1	Handling of conditional branches in single threaded pipelines . . . . .	3
2.2	Sample code for GCD with conditional branches . . . . .	3
2.3	Sample code with data dependencies . . . . .	5
2.4	Handling of data dependencies in single threaded pipelines . . . . .	6
2.5	Simple Multithreaded Pipeline . . . . .	7
2.6	Sample execution sequence of a thread-interleaved pipeline with 5 threads and 5 pipeline stages . . . . .	8
3.1	Image Placeholder . . . . .	11
4.1	Image Placeholder . . . . .	12
6.1	Image Placeholder . . . . .	16

# List of Tables

## Acknowledgments

I want to thank my wife

I want to thank my parents

I want to thank my advisor, Edward A. Lee

I want to thank the committee members

I want to thank Hiren Patel

I want to thank Jan Rieneke

I would also like to thank the ptolemy group especially Christopher and Mary, Jia for providing me the template

I would like to thank everyone else that made this possible



# Chapter 1

## Introduction

Outline

(**Todo: make sure to add in timing anomalies**)

### 1.1 Background

- Discuss the problem
- show the difficulty in execution time analysis as background

Fig. 1.1 shows an image

### 1.2 Intro Section Header 2

Here is another header

The remaining chapters are organized as follows. Chapter 4 surveys the related research that has been done on architectures to make them more analyzable. Chapter 2 explains the architecture of PRET including the thread-interleaved pipeline and memory hierarchy, Chapter ??, Chapter 6, Chapter 7,

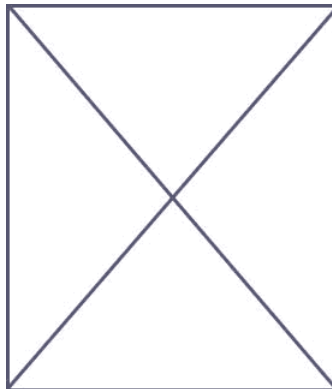


Figure 1.1: Image Placeholder

## Chapter 2

# Precision Timed Machine

In this chapter we present the PREcision Timed (PRET) Machine. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nam eu est neque. Suspendisse mollis gravida mi in blandit. Vivamus porta libero at massa sagittis pellentesque. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed nibh magna, facilisis ac dapibus vitae, tincidunt nec magna. Morbi ac neque in est porta placerat. Duis viverra blandit ante, ut scelerisque arcu sodales vel. Aenean sapien erat, tincidunt malesuada accumsan a, eleifend feugiat leo. Maecenas auctor nulla non purus fringilla nec hendrerit massa facilisis. Donec vel diam nibh. Maecenas sed massa non mauris faucibus condimentum et et metus. Fusce placerat, dolor et adipiscing suscipit, orci lectus fringilla mauris, a tincidunt dolor mauris id ligula. Vestibulum luctus, dolor in bibendum accumsan, leo turpis suscipit enim, et hendrerit odio metus eget dolor. Morbi in lectus massa.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nam eu est neque. Suspendisse mollis gravida mi in blandit. Vivamus porta libero at massa sagittis pellentesque. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed nibh magna, facilisis ac dapibus vitae, tincidunt nec magna. Morbi ac neque in est porta placerat. Duis viverra blandit ante, ut scelerisque arcu sodales vel. Aenean sapien erat, tincidunt malesuada accumsan a, eleifend feugiat leo. Maecenas auctor nulla non purus fringilla nec hendrerit massa facilisis. Donec vel diam nibh. Maecenas sed massa non mauris faucibus condimentum et et metus. Fusce placerat, dolor et adipiscing suscipit, orci lectus fringilla mauris, a tincidunt dolor mauris id ligula. Vestibulum luctus, dolor in bibendum accumsan, leo turpis suscipit enim, et hendrerit odio metus eget dolor. Morbi in lectus massa.

It is important to understand why and how current architectures fall short of timing predictability and repeatability. Thus, before we present the PRET architecture, we briefly discuss common architectural designs and their effects on execution time, and point out some key issues and trade-offs when designing architectures for predictable and repeatable timing. The introduction of pipelining vastly improved architecture average-case performance. It allowed faster clock speeds for processors, and improved the instruction throughput compared to single cycle architectures. Pipelining the datapath allowed subsequent instructions to begin execution while prior instructions were still being completed. Ideally each processor cycle one instruction completes and leaves the pipeline as another enters and begins execution. In reality, different pipeline hazards occur which reduce the throughput and create stalls in the pipeline. The handling of these hazards an important factor to the timing predictability and repeatability of the architecture design. To illustrate this point, we discuss basic hardware additions proposed to reduce performance penalty from hazards,

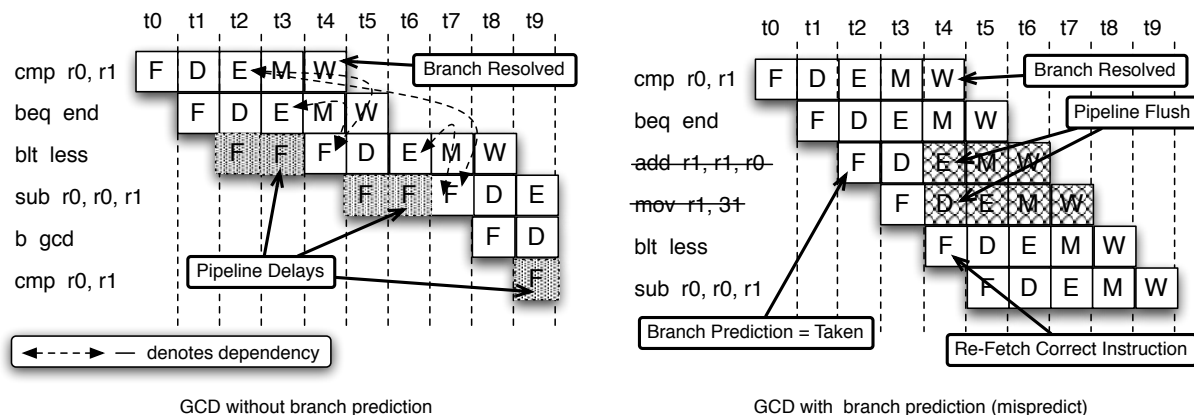


Figure 2.1: Handling of conditional branches in single threaded pipelines

and their effects on execution time and predictability.

We began by looking at how control-flow changes are handled in pipelines. Branches cause control-flow hazards in the pipeline; the instruction after the branch, which should be fetched the next cycle, is unknown until after the branch instruction is completed. Conditional branches adds more complexity (Todo: ?), as whether or not the branch is taken depends on an additional condition bit. The code segment in figure 2.2 shows assembly instructions from the ARM instruction set architecture (ISA) that implement the Greatest Common Divisor (GCD) algorithm using conditional branch instructions *beq* (branch equal) and *blt* (branch less than). Conditional branch instructions in ARM branch based on conditional flags that are set with special compare instructions (Todo: citation).

The *cmp* instruction is one such compare instruction that subtracts two registers and updates the conditional flags according to the results. The GCD implementation shown in the code uses this mechanism to determine whether to continue or end the algorithm. Figure 2.1 show two ways branches are commonly handled in a single-threaded pipeline. In the figure, time progresses horizontally towards the right, each time step, or column, represents a processor cycle. Each row represents an instruction that is fetched and executed within the pipeline. Each block represents the instruction entering the different stages of the pipeline – fetch (F), decode (D), execute (E), memory (M) and writeback (W).

A simple but effective way of handling control-flow hazards is by simply stalling the pipeline until the branch instruction completes. This is shown on the left of figure 2.1. Two pipeline delays (or bubbles) are inserted after each branch instruction to wait until address calculation is com-

```

gcd:
    cmp r0, r1      # compare r0 and r1
    beq end        # branch if r0 == r1
    blt less       # branch if r0 < r1
    sub r0, r0, r1  # r0 = r0 - r1
    b gcd          # branch to label gcd
less:
    sub r1, r1, r0  # r1 = r1 - r0
    b gcd          # branch to label gcd
end:
    add r1, r1, r0  # r1 = r1 + r0
    mov r3, r1     # r3 = r1

```

Figure 2.2: Sample code for GCD with conditional branches

pleted. The dependencies between instructions are also drawn out to make clear why the pipeline bubbles are necessary. In order for the *blt* instruction to be fetched, its address must be calculated during the execution stage of the *beq* instruction. At the same time, because *beq* is a conditional branch, whether or not the branch is taken depends on the *cmp* instruction. The architecture used contain forwarding circuitry, so the addresses calculated by the branch instructions and the results of the *cmp* instruction can be used before the instructions are committed. The performance penalty incurred is the pipeline delays inserted to wait for the branch address calculation to complete. Conditional branches will also incur extra delays for deeper pipelines if the branch condition cannot be resolved in time. Some architectures enforce the compiler to insert one or more non-dependent instructions after a branch that is always executed before the change in control-flow of the program. These are called branch delay slots and can mitigate the branch penalty, but become less effective as pipelines grow deeper by design.

In attempt to remove the need of inserting pipeline bubbles, branch predictors were invented to predict the results of a branch before it is resolved (Todo: citation). Many clever branch predictors have been proposed, and they can accurately predict branches up to 93.5% (Todo: citation). Branch predictors predict the condition and target addresses of branches, so pipelines can speculatively continue execution based upon the prediction. If the prediction was correct, no penalty occurs for the branch, and execution simply continues. However, when a mispredict occurs, then the speculatively executed instructions need to be flushed and the correct instructions need to be refetched into the pipeline for execution. The right of figure 2.1 shows the execution of GCD in the case of a branch misprediction. After the *beq* instruction, the branch is predicted to be taken, and the *add* and *mov* instructions from the label *end* is directly fetched into execution. When the *cmp* instruction is completed, a misprediction is detected, so the *add* and *mov* instruction are flushed out of the pipeline while the correct instruction *blt* is immediately re-fetched and execution continues. The misprediction penalty is typically the number of stages between fetch and execute, as those cycles are wasted executing instructions from an incorrect execution path. This penalty only occurs on a mispredict, thus branch prediction typically yields better average performance and is preferred for modern architectures. Nonetheless, it is important to understand the effects of branch prediction on execution time.

Typical branch predictors predict branches based upon the history of previous branches encountered. As each branch instruction is resolved, the internal state of the predictor, which stores the branch histories, is updated and used to predict the next branch. This implicitly creates a dependency between branch instructions and their execution history, as the prediction is affected by its history. In other words, the execution time of a branch instruction will depend on the branch results of previous branch instructions. During static execution timing analysis, the state of the branch predictor is unknown because it is often infeasible to keep track of execution history so far back. There has been work on explicitly modeling branch predictors for execution time analysis (Todo: citation), but the results are (Todo: the results of branch predictor modeling for execution time analysis). The analysis needs to conservatively account for the potential branch mispredict penalty for each branch, which leads to overestimated execution times. To make matters worse, as architectures grow in complexity, more internal states exist in architectures that could be affected by the speculative execution. For example, cache lines could be evicted when speculatively executing instructions from a mispredicted path, changing the state of the cache. This makes a tight static execution time analysis extremely difficult, if not impossible; explicitly modeling all hardware states and their ef-

fects together often lead to an infeasible explosion in state space. On the other hand, although the simple method of inserting pipeline bubbles for branches could lead to more branch penalties, the static timing analysis is precise and straight forward, as no prediction and speculative execution occur. In the case of the ARM ISA, the analysis simply accounts for the branch penalty after every branch. Additional penalties from a conditional branch can be accounted for by simply checking for instructions that modify the conditional flag above the conditional branch. We explicitly showed this simple method of handling branches to point out an important trade-off between speculative execution for better average performance and consistent stalling for better predictability. Average-case performance can be improved by speculation at the cost of predictability and potentially prolonging the worst-case performance. The challenge remains to maintain predictability while improving worst-case performance, and how pipeline hazards are handled play an integral part of tackling this challenge.

Data hazards occur when instructions depend on the results of previous instructions that have yet to commit. The code segment shown in figure 2.3 contains instructions that each depend on the result of its previous instruction. The top of figure 2.4 shows an execution of the code segment on a naive pipeline without handling of data hazards. Data hazards in this case are handled by inserting pipeline de-

add r0, r1, r2	# r0 = r1 + r2
sub r1, r0, r1	# r1 = r0 - r1
ldr r2, [r1]	# r2 = mem[r1]
sub r0, r2, r1	# r0 = r2 - r1
cmp r0, r3	# compare r0 and r3

Figure 2.3: Sample code with data dependencies

lays to ensure the completion of all dependent instructions. Similar to inserting pipeline delays of control-flow hazards, this method allows for predictable static execution time analysis, but at a slight cost of performance. Pipeline forwarding is the most common way of handling data hazards that occur from pipelining. A pipeline forwarding circuitry consists of backwards paths for data from different pipeline stages to the inputs of arithmetic units, and multiplexers to select amongst them. It provides a way to directly access computation results from the previous instruction before it commits. The pipeline controller dynamically detects whether a data-dependency exists, and changes the selection bits to the multiplexers accordingly so the correct operands are selected. The bottom of figure 2.4 shows the execution with forwarding in the pipeline. No pipeline bubbles are needed for the first *sub* instruction and *ld* instruction, as the results they depend on can be computed in one cycle by the ALU, and forwarded through the forwarding paths. (Todo: Needs some reworking here) Notice that although there is dynamic execution in the pipeline forwarding circuitry, it is actually possible to statically predict the execution time accurately. The logic in the pipeline controller that enables and selects the correct forwarding bits only needs to check a small set of previous instructions to detect data-dependencies. Thus, static execution time analysis can detect forwarding by simply checking a short window of previous instructions to account for stalls accordingly. (Todo: find papers to back this up) The internal state of the branch predictor on the other hand is dependent on branch histories which may have happened arbitrarily long ago in the execution sequence, (Todo: more complicated ways of handling data hazards include out-of-order execution...)

(Todo: Talk about caches here) We assume the *ld* instruction accesses results in a cache miss, requiring access to main memory. Stalls are still inserted for the the second *sub* instruction, as it is waiting upon the results of a memory operation. The memory access latency in the figure is arbitrarily chosen to be 5 cycles so the figure is not too long and instructions after the *ld* instruction

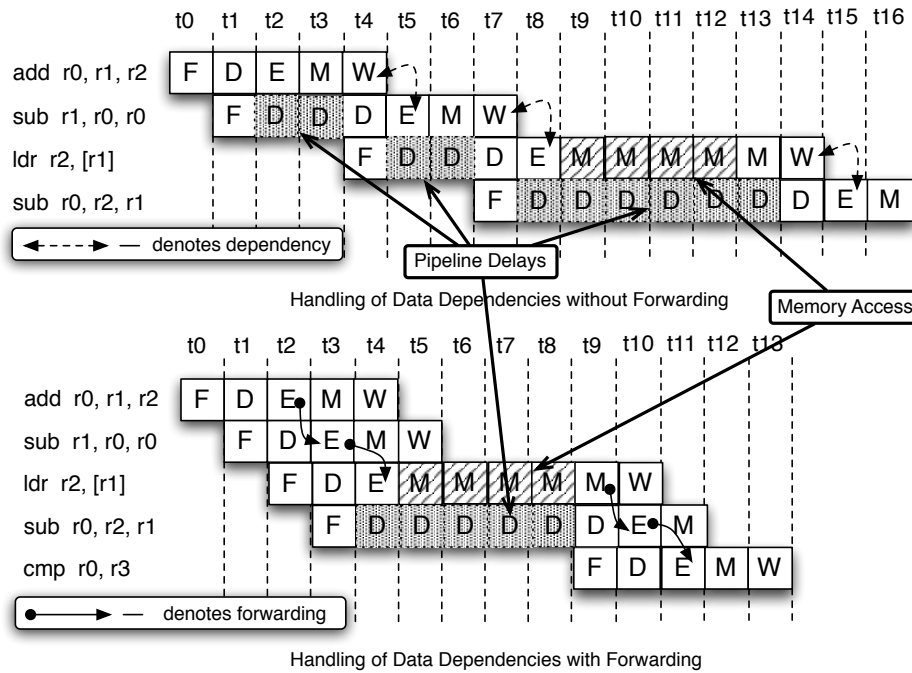


Figure 2.4: Handling of data dependencies in single threaded pipelines

can be shown. Typical memory access latencies could range from 10 to 1000 cycles. We can thus see that forwarding can address the data-dependencies caused by pipelining – the read-after-write of register computations. However, they cannot address the data-dependencies caused by other long latency operations such as memory operations, so pipeline stalls are still needed.

Multithreaded architectures were introduced to improve instruction throughput over instruction latency. The architecture optimizes thread-level parallelism over instruction-level parallelism to improve performance. Multiple hardware threads are introduced into the pipeline to fully utilize thread-level parallelism. When one hardware thread is stalled, another hardware thread can be fetched into the pipeline for execution to avoid stalling the whole pipeline. To lower the context switching overhead, the pipeline contains physically separate copies of hardware thread states, such as registers files and program counters etc, for each hardware thread. Figure 2.5 shows a architectural level view of a simple multithreaded pipeline. It contains 5 hardware threads, so it has 5 copies of the Program Counter (PC) and Register files. Once a hardware thread is executing in the pipeline, its corresponding thread state can be selected by signaling the correct selection bits to the multiplexers. The rest of the pipeline remains similar to a traditional 5 stage pipeline as introduced in Hennessy and Pattern(ToDo: citation). The extra copies of the thread state and the multiplexers used to select them thus contribute to most of the hardware additions needed to implement hardware multithreading.

Ungerer et al. [22] surveyed different multithreaded architectures and categorized them based upon the (ToDo: thread selection?) policy and the execution width of the pipeline. The thread selection policy is the context switching scheme used to determine which threads are executing, and how often a context switch occurs. Coarse-grain policies manage hardware threads similar to the way operation systems manage software threads. A hardware thread gain access to

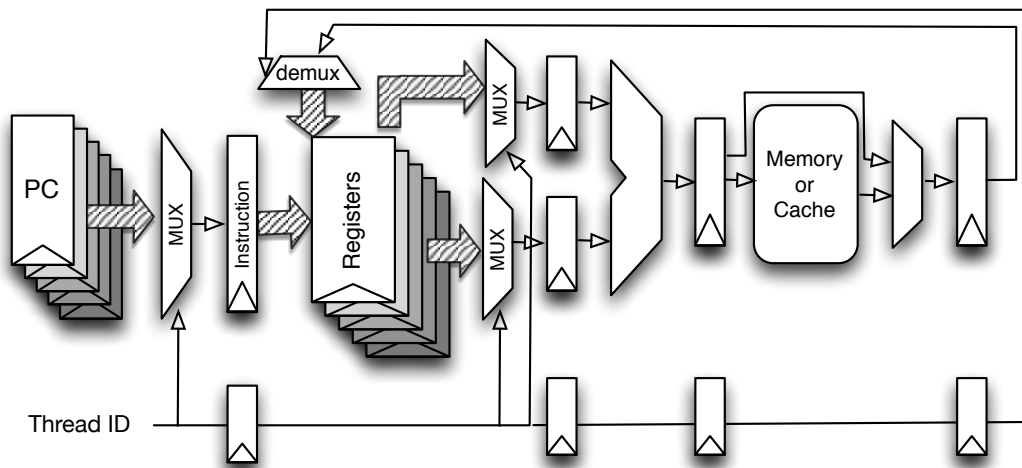


Figure 2.5: Simple Multithreaded Pipeline

the pipeline and continues to execute until a context switch is triggered. Context switches occur less frequently via this policy, so less hardware threads are required to fully utilize the processor. Different coarse-grain policies trigger context switches with different events. Some trigger on dynamic events, such as cache miss or interrupts, and some trigger on static events, such as specialized instructions. Fine-grain policies switch context much more frequently – usually every processor cycle. Both coarse-grain and fine-grain policies can also have different hardware thread scheduling algorithms that are implemented in a hardware thread scheduling controller to determine which hardware thread is switched into execution. The width of the pipeline refers to the number of instructions that can be fetched into execution in one cycle. For example, superscalar architectures have redundant functional units, such as multipliers and ALUs, and can dispatch multiple instructions into execution in a single cycle. Multithreaded architectures with pipeline widths of more than one, such as Simultaneous Multithreaded (SMT) architectures, can fetch and execute instructions from several hardware threads in the same cycle.

Multithreaded architectures typically bring additional challenges to execution time analysis of software running on them. Any timing analysis for code running on a particular hardware thread needs to take into account not only the code itself, but also the thread selection policy of the architecture and sometimes even the execution context of code running on other hardware threads. For example, if dynamic coarse-grain multithreading is used, then a context switch could occur at any point when a hardware thread is executing in the pipeline. This not only has an effect on the control flow of execution, but also the state of any hardware that is shared, such as caches or branch predictors. Thus, it becomes nearly impossible to estimate execution time without knowing the exact execution state of other hardware threads and the state of the thread scheduling controller. However, it is possible for multithreaded architectures to fully utilize thread-level parallelism while still maintaining timing predictability. Thread-interleaved pipelines use a fine-grain thread switching policy with round robin thread scheduling to achieve high instruction throughput while still allowing precise timing analysis for code running on its hardware threads. Below, its architecture and trade-offs are described and discussed in detail along with examples and explanation of

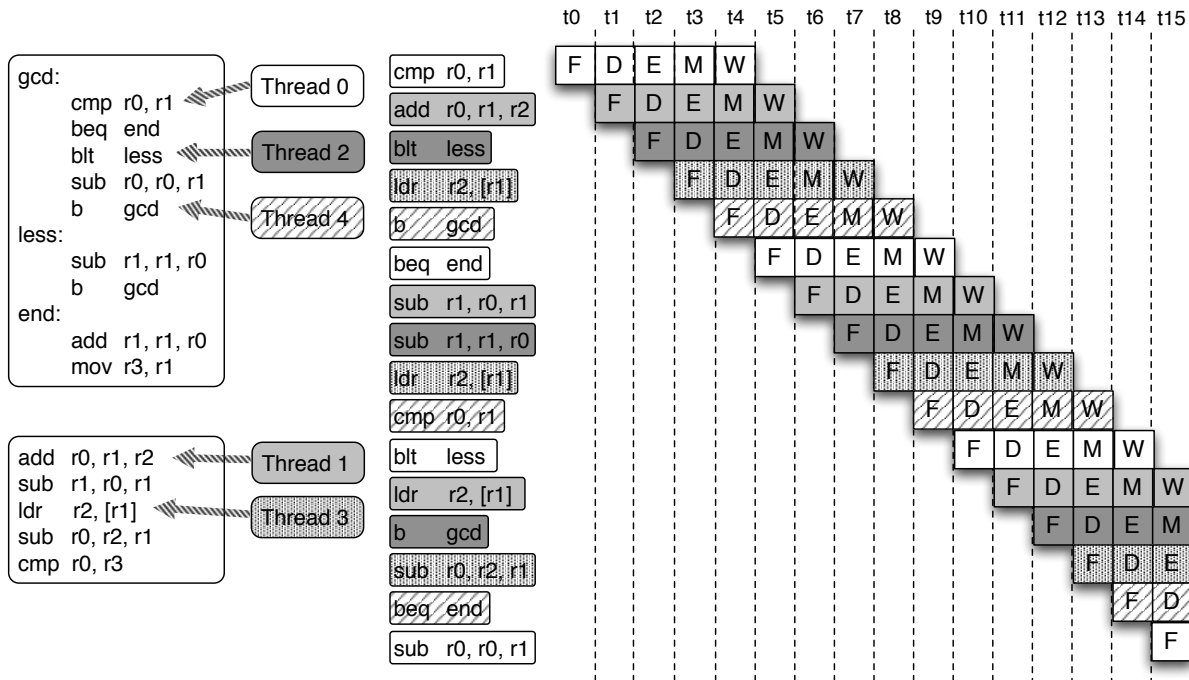


Figure 2.6: Sample execution sequence of a thread-interleaved pipeline with 5 threads and 5 pipeline stages

how timing predictability is maintained. Through the remainder of this chapter, we will use the term “thread” to refer to explicit hardware threads that have physically separate register files, program counters, and other thread states. This is not to be confused the common notion of “threads”, which is assumed to be software threads that is managed by operating systems with thread states stored in memory.

## 2.1 Thread-Interleaved Pipelines

The thread-interleaved pipeline was introduced to improve the response time of handling multiple I/O devices (Todo: citation). I/O operations often stall from the communication with the I/O devices. Thus, interacting with multiple I/O devices leads to wasted processor cycles that are idle waiting for the I/O device to respond. By employing multiple hardware thread contexts, a hardware thread stalled from the I/O operations does not stall the whole pipeline, as other hardware threads can be fetched and executed. Thread-interleaved pipelines use fine-grain multithreading; every cycle a context switch occurs and a different hardware thread is fetched into execution. The threads are scheduled in a deterministic round robin fashion. This also reduces the context switch overhead down to nearly zero, as no time is needed to determine which thread to fetch next, and barely any hardware is required to implement round robin thread scheduling; a simple  $\log(n)$  bit up counter (for  $n$  threads) would suffice. Figure 2.6 shows an example execution sequence from a 5 stage thread-interleaved pipeline with 5 threads. The thread-interleaved pipelines shown and presented in this thesis are all of single width. The same code segments from figure 2.2 and figure 2.3 are being executed in this pipeline. Threads 0, 2 and 4 execute GCD (figure 2.2) and threads 1



and 3 execute the data dependent code segment (figure 2.3). Each hardware thread executes as an independent context and their progress is shown in figure 2.6 with thick arrows pointing to the execution location of each thread at  $t_0$ . We can observe from the figure that each time step an instruction from a different hardware thread is fetched into execution and the hardware threads are fetched in a round robin order. At time step 4 we begin to visually see that each time step, each pipeline stage is occupied by a different hardware thread. The fine-grained thread interleaving and the round robin scheduling combine to form this unique property of thread-interleaved pipelines, which provides the basis for a timing predictable architecture design.

For thread-interleaved pipelines, if there are, at a minimum, the same number of threads as there are pipeline stages, then at each time step no dependency exists between the pipeline stages since they are each executing on a different thread. As a result, data and control pipeline hazards, the results of dependencies between stages within the pipelines, no longer exist in the thread-interleaved pipeline. We've already shown from figure 2.1 that when executing the GCD code segment on a single-threaded pipeline, control hazards stem from branch instructions because of the address calculation for the instruction after the branch. However, in a thread-interleaved pipeline, the instruction after the branch from the same thread is not fetched into the pipeline until the branch instruction is committed. Before that time, instructions from other threads are fetched so the pipeline is not stalled, but simply executing other thread contexts. This can be seen in figure 2.6 for thread 0, which is represented with instructions with white backgrounds. The *cmp* instructions, which determines whether next conditional branch *beq* is taken or not, completes before the *beq* is fetched at time step 5. The *blt* instruction from thread 0, fetched at time step 10, also causes no hazard because the *beq* is completed before *blt* is fetched. The code in figure 2.3 is executed on thread 1 of the thread interleave pipeline in figure 2.6. The pipeline stalls inserted from top of figure 2.4 are no longer needed even without a forwarding circuitry because the data-dependent instructions are fetched after the completion of its previous instruction. In fact, no instruction in the pipeline is dependent on another because each pipeline stage is executing on a separate hardware thread context. Therefore, the pipeline does not need to include any extra logic or hardware for handling data and control hazards in the pipeline. This gives thread-interleaved pipelines the advantage of a simpler pipeline design that requires less hardware logic, which in turns allows the pipeline clock speed to increase. Thread-interleaved pipelines can be clocked at higher speeds since each pipeline stage contains significantly less logic needed to handle hazards. The registers and processor states use much more compact memory cells compared to the logic and muxes used to select and handle hazards, so the size footprint of thread-interleaved pipelines are also typically smaller.

For memory operations that have longer latencies, thread-interleaved pipelines hides the latency with its execution of other threads. Thread 3 in figure 2.6 shows the execution of a *ld* instruction that takes the same 5 cycles as shown in figure 2.4. We again assume that this *ld* instruction accesses data from the main memory. While the *ld* instruction is waiting for memory access to complete, the thread-interleaved pipeline executes instructions from other threads. The next instruction from thread 3 that is fetched into the pipeline is again the same *ld* instruction. As memory completes its execution during the execution of instructions from other threads, we replay the same instruction to pick up the results from memory and write it into registers to complete the execution of the *ld* instruction. It is possible to directly write the results back into the register file when the memory operation completes. This would however require non-trivial hardware additions to support and manage multiple write-back paths in the pipeline so contention can be avoided with

the existing executing threads. Thus, in our design we simply replay the instruction for write-backs to simplify design and piggy back on the existing write-back datapath.

In the real world, main memory of systems are typically made of DRAMs, which have access latencies that range from 10 to 100 cycles; one round of the round robin thread schedule would be insufficient time to hide the memory latency access. Often times multithreaded architectures mark these threads inactive and the thread is not rescheduled until the data is ready. This is done to maximize the throughput of the pipeline, since threads waiting for data from other hardware components cannot make any progress until the data is returned. However, this leads to unpredictable timing behaviors in the threads. When threads are scheduled and unscheduled dynamically, the other threads in the pipeline would dynamically execute more or less frequently depending on which threads are active and inactive. This greatly complicates any timing analysis on the software running on each thread as the execution frequency of the threads would depend on the execution of other threads. Thus, our thread-interleaved pipeline does not mark threads inactive, but simply replays the instruction from the thread. The effects of latency hiding is still present, as other threads continue to progress while one thread is replaying its multi-cycle instruction.

In addition, the shared main memory is a structural hazard, a hazard that occurs when a hardware unit is needed by two or more instructions at the same time. While a thread is accessing main memory, other threads could just as well fetch a *ld* or *str* instruction, and cause contention to access main memory. Care must also be taken when adding datapaths that take multiple cycles, or else the interference introduced could easily disrupt the timing analysis of threads. If the added datapath isn't able to support pipelined or simultaneous operations, then it will introduce contention amongst the threads. As one thread executes a floating-point division instruction, any other thread that also executes a floating-point division must now wait until the first instruction finishes. If other threads also executes the same instruction, then queuing mechanisms must be introduced, for threads that are contending for the floating-point divider. This would greatly complicate the timing analysis, as the execution time of floating-point division instructions now depend on the execution context of other threads. Pipelining the floating-point divider would increase the throughput at the cost of area and latency. However, by pipelining the floating-point divider unit, each thread that executes a floating-point division can now access it without contention. The replay mechanism also hides the long latency of the instruction, and benefits from the improved latency. Because there is no contention, the timing analysis of floating-point operations are now trivial and predictable.

In summary, blah blah blah...

## 2.2 Memory System

## Chapter 3

# Programming Models

Intro text here

### 3.1 PRET Programming model Section Header

Fig. 3.1 shows an image

### 3.2 Pret Programming model Section Header 2

Here is another header

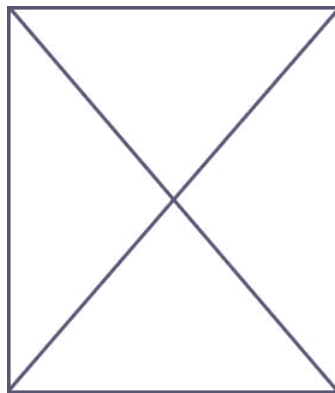


Figure 3.1: Image Placeholder

## Chapter 4

# Related Work

### 4.1 Architectural Modifications

Fig. 4.1 shows an image

### 4.2 Related Section Header 2

Here is another header

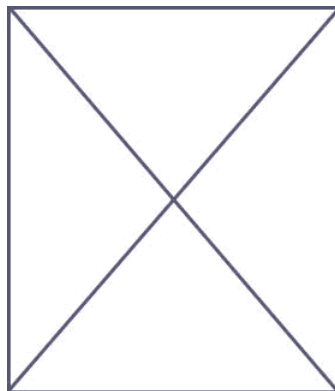


Figure 4.1: Image Placeholder

## **Chapter 5**

# **Implementation of PTARM**

This section documents the implementation of PTARM

### **5.1 Implementation**

#### **5.1.1 PTARM Simulator**

#### **5.1.2 PTARM VHDL Softcore**

#### **5.1.3 Worst Case Execution Time Analysis**

## Chapter 6

# Applications

### 6.1 Eliminating Side-Channel-Attacks

Encryption algorithms are based on strong mathematical properties to prevent attackers from deciphering the encrypted content. However, their implementations in software naturally introduce varying run times because of data-dependent control flow paths. Timing attacks [14] exploit this variability in cryptosystems and extract additional information from executions of the cipher. These can lead to deciphering the secret key. Kocher describes a timing attack as a basic signal detection problem [14]. The “signal” is the timing variation caused by the key’s bits when running the cipher, while “noise” is the measurement inaccuracy and timing variations from other factors such as architecture unpredictability and multitasking. This signal to noise ratio determines the number of samples required for the attack – the greater the “noise,” the more difficult the attack. It was generally conceived that this “noise” effectively masked the “signal,” thereby shielding encryption systems from timing attacks. However, practical implementations of the attack have since been presented [6, 9, 25] that clearly indicate the “noise” by itself is insufficient protection. In fact, the architectural unpredictability that was initially believed to prevent timing attacks was discovered to enable even more attacks. For example, computer architects use caches, branch predictors and complex pipelines to improve the average-case performance while keeping these optimizations invisible to the programmer. These enhancements, however, result in unpredictable and uncontrollable timing behaviors, which were all shown to be vulnerabilities that led to side-channel attacks [3, 20, 2, 8].

In order to not be confused with Kocher’s [14] terminology of *timing attacks* on algorithmic timing differences, we classify all above attacks that exploit the timing variability of software implementation *or* hardware architectures as *time-exploiting attacks*. In our case, a *timing attack* is only one possible *time-exploiting attack*. Other time-exploiting attacks include branch predictor, and cache attacks. Examples of other side-channel attacks are power attacks [17, 13], fault injection attacks [4, 10], and many others [25].

In recent years, we have seen a tremendous effort to discover and counteract side-channel attacks on encryption systems [4, 8, 15, 11, 1, 12, 7, 24, 23]. However, it is difficult to be fully assured that all possible vulnerabilities have been discovered. The plethora of research on side-channel exploits [8, 4, 15, 11, 1, 12, 7, 24, 23] indicates that we do not have the complete set of solutions as more and more vulnerabilities are still being discovered and exploited. Just recently, Coppens et al. [8] discovered two previously unknown time-exploiting attacks on modern x86 pro-

cessors caused by the out-of-order execution and the variable latency instructions. This suggests that while current prevention methods are effective at *defending* against their particular attacks, they do not *prevent* other attacks from occurring. This, we believe, is because they do not address the root cause of time-exploiting attacks, which is that run time variability *cannot be controlled* by the programmer.

It is important to understand that the main reason for time-exploiting attacks is *not* that the program runs in a varying amount of time, but that this variability *cannot be controlled* by the programmer. The subtle difference is that if timing variability is introduced in a controlled manner, then it is still possible to control the timing information that is leaked during execution, which can be effective against time-exploiting attacks. However, because of the programmer’s *lack of control* over these timing information leaks in modern architectures, noise injection techniques are widely adopted in attempt to make the attack infeasible. These include adding random delays [14] or blinding signatures [14, 7]. Other techniques such as branch equalization [18, 25] use software techniques to rewrite algorithms such that they take equal time to execute during each conditional branch. We take a different approach, and directly address the crux of the problem, which is the *lack of control* over timing behaviors in software. We propose the use of an embedded computer architecture that is designed to allow predictable and controllable timing behaviors.

At first it may seem that a predictable architecture makes the attacker’s task simpler, because it reduces the amount of “noise” emitted from the underlying architecture. However, we contend that in order for timing behaviors to be controllable, the underlying architecture *must* be predictable. This is because it is meaningless to specify any timing semantics in software if the underlying architecture is unable to honor them. And in order to guarantee the execution of the timing specifications, the architecture must be predictable. Our approach does not attempt to increase the difficulty in performing time-exploiting attacks, but to eliminate them completely.

In this paper, we present the PREcision Timed (PRET) architecture [16] in the context of embedded cryptosystems, and show that an architecture designed for predictability and controllability effectively eliminates all time-exploiting attacks. Originally proposed by Lickly et al [16], PRET provides instruction-set architecture (ISA) extensions that allow programmers to control an algorithm’s temporal properties at the software level. To guarantee that the timing specifications are honored, PRET provides a predictable architecture that replaces complex pipelines and speculation units with multithread-interleaved pipelines, and replaces caches with software-managed fast access memories. This allows PRET to maintain predictability without sacrificing performance. We target embedded applications such as smartcard readers [15], key-card gates [5], set-top boxes [15], and thumbpods [21], which are a good fit for PRET’s embedded nature. We demonstrate the effectiveness of our approach by running both the RSA and DSA [19] encryption algorithms on PRET, and show its immunity against time-exploiting attacks. This work shows that a disciplined defense against time-exploiting attacks requires a combination of software and hardware techniques that ensure controllability and predictability.

## 6.2 Real Time 1D Computational Fluid Dynamics Simulator

Here is another header

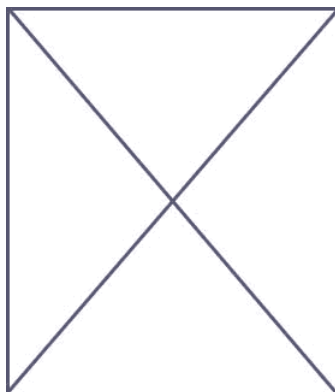


Figure 6.1: Image Placeholder



## **Chapter 7**

# **Conclusion and Future work**

### **7.1 Summary of Results**

This is my summary

### **7.2 Future Work**

Here is what you can keep doing

# Bibliography

- [1] O. Aciicmez, Çetin Kaya Koç, and J.-P. Seifert. On the Power of Simple Branch Prediction Analysis. In *ASIACCS '07: Proceedings of the 2nd ACM symposium on Information, computer and communications security*, pages 312–320, New York, NY, USA, 2007. ACM.
- [2] O. Aclicmez, J. pierre Seifert, and C. K. Koc. Predicting secret keys via branch prediction. In *in Cryptology CT-RSA 2007, The Cryptographers Track at the RSA Conference 2007*, pages 225–242. Springer-Verlag, 2007.
- [3] D. J. Bernstein. Cache-timing Attacks on AES, 2004.
- [4] E. Biham and A. Shamir. Differential Fault Analysis of Secret Key Cryptosystems. *Lecture Notes in Computer Science*, 1294:513–525, 1997.
- [5] S. C. Bono, M. Green, A. Stubblefield, A. Juels, A. D. Rubin, and M. Szydlo. Security analysis of a cryptographically-enabled rfid device. In *SSYM'05: Proceedings of the 14th conference on USENIX Security Symposium*, pages 1–1, Berkeley, CA, USA, 2005. USENIX Association.
- [6] D. Brumley and D. Boneh. Remote timing attacks are practical. In *SSYM'03: Proceedings of the 12th conference on USENIX Security Symposium*, pages 1–1, Berkeley, CA, USA, 2003. USENIX Association.
- [7] D. Chaum. Blind Signatures for Untraceable Payments. In *Advances in Cryptology: Proceedings of Crypto 82*, pages 199–203. Plenu Press, 1983.
- [8] B. Coppens, I. Verbauwhede, K. De Bosschere, and B. De Sutter. Practical Mitigations for Timing-Based Side-Channel Attacks on Modern x86 Processors, 2009.
- [9] J.-F. Dhem, F. Koeune, P.-A. Leroux, P. Mestre, J.-J. Quisquater, and J.-L. Willems. A Practical Implementation of the Timing Attack. In J.-J. Quisquater and B. Schneier, editors, *Proceedings of the Third Working Conference on Smart Card Research and Advanced Applications (CARDIS 1998)*. Springer-Verlag, 1998.
- [10] M. Feng, B. B. Zhu, M. Xu, S. Li, B. B. Zhu, M. Feng, B. B. Zhu, M. Xu, and S. Li. Efficient Comb Elliptic Curve Multiplication Methods Resistant to Power Analysis, 2005.
- [11] R. Karri, K. Wu, P. Mishra, and Y. Kim. Fault-Based Side-Channel Cryptanalysis Tolerant Rijndael Symmetric Block Cipher Architecture. In *DFT '01: Proceedings of the IEEE International Symposium on Defect and Fault Tolerance in VLSI Systems*, page 427, Washington, DC, USA, 2001. IEEE Computer Society.

- [12] J. Kelsey, B. Schneier, D. Wagner, and C. Hall. Side Channel Cryptanalysis of Product Ciphers. In *Journal of Computer Security*, pages 97–110. Springer-Verlag, 1998.
- [13] P. Kocher, J. J. E. and B. Jun. Differential Power Analysis. In *Lecture Notes in Computer Science*, pages 388–397. Springer-Verlag, 1999.
- [14] P. C. Kocher. Timing attacks on implementations of Diffie-Hellman, RSA, DSS, and other systems. In *Lecture Notes in Computer Science*, pages 104–113. Springer-Verlag, 1996.
- [15] O. Kömmerling and M. G. Kuhn. Design Principles for Tamper-Resistant Smartcard Processors. In *USENIX Workshop on Smartcard Technology proceedings*, pages 9–20, 1999.
- [16] B. Lickly, I. Liu, S. Kim, H. D. Patel, S. A. Edwards, and E. A. Lee. Predictable Programming on a Precision Timed Architecture. In *CASES '08: Proceedings of the 2008 international conference on Compilers, architectures and synthesis for embedded systems*, pages 137–146, New York, NY, USA, 2008. ACM.
- [17] T. S. Messerges, E. A. Dabbish, and R. H. Sloan. Investigations of Power Analysis Attacks on Smartcards. In *In USENIX Workshop on Smartcard Technology*, pages 151–162, 1999.
- [18] D. Molnar, M. Piotrowski, D. Schultz, and D. Wagner. The Program Counter Security Model: Automatic Detection and Removal of Control-Flow Side Channel Attacks. In *In Cryptology ePrint Archive, Report 2005/368*, 2005.
- [19] National Institute of Standards and Technology. "Digital Signature Standard". Federal Information Processing Standards Publication 186, 1994.
- [20] C. Percival. Cache missing for fun and profit. In *Proc. of BSDCan 2005*, page 05, 2005.
- [21] P. Schaumont, K. Sakiyama, Y. Fan, D. Hwang, S. Yang, A. Hodjat, B. Lai, and I. Verbauwhede. Testing ThumbPod: Softcore bugs are hard to find. In *Eighth IEEE International High-Level Design Validation and Test Workshop, 2003*, pages 77–82, 2003.
- [22] T. Ungerer, B. Robič, and J. Šilc. A survey of processors with explicit multithreading. *ACM Comput. Surv.*, 35:29–63, March 2003.
- [23] Z. Wang and R. B. Lee. Covert and Side Channels Due to Processor Architecture. In *ACSAC '06: Proceedings of the 22nd Annual Computer Security Applications Conference*, pages 473–482, Washington, DC, USA, 2006. IEEE Computer Society.
- [24] Z. Wang and R. B. Lee. New cache designs for thwarting software cache-based side channel attacks. In *Proceedings of the 34th annual international symposium on Computer architecture*, pages 494 – 505, San Diego, CA, June 2007 2007.
- [25] Yongbin. Side-channel attacks: Ten years after its publication and the impacts on cryptographic module security testing.