



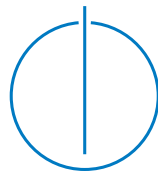
FAKULTÄT FÜR INFORMATIK

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis

Host Compiled Simulation for Timing and Power Estimation

Gaurav Kukreja





FAKULTÄT FÜR INFORMATIK

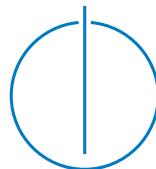
TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis

Host Compiled Simulation for Timing and Power Estimation

Author:	Gaurav Kukreja
Supervisor:	Prof. Michael Gerndt
Advisor:	Dr. Josef Weidendorfer
Advisor:	Mr. Bo Wang

Submission Date: 15th October, 2014



I confirm that this Master's Thesis is my own work and I have documented all sources and materials used.

Munich, 15th October, 2014

Gaurav Kukreja

Abstract

Simulation is a useful technique for Hardware Software Co-development. It is performed at various levels of abstraction to serve different purposes. Instruction Set Simulation is the lowest level of abstraction where the processor pipeline is simulated in detail, to allow hardware developers to test their modifications and evaluate the impact on performance. At higher levels of abstraction, simulation provides developers with a tangible environment for early software development. The focus of this project is on simulation for performance estimation, namely, estimation of time and power consumed in running a benchmark application on a target processor.

While Instruction Set Simulators are known to be highly accurate, they are difficult to develop and slow to execute because of the level of detail they address. Host Compiled Simulation is a technique to accelerate performance estimation with negligible impact on accuracy. The idea is to instrument¹ the source code, by taking into consideration the behaviour of the target processor. The instrumented source code is compiled and run on the Host Machine. The technique relies on the assumption that performance of each basic block² in the binary code can be accurately estimated on a certain processor by emulating the pipeline. Other aspects that affect performance, like resources spent in memory access can be accounted for, and a fairly accurate estimate of the time and power consumed can be estimated.

In this project, a tool to perform Host Compiled Simulation was developed. This thesis discusses the state-of-art in simulation. It explains the approach used to develop this tool. The results showing accuracy of estimations from this approach are presented.

¹Instrumentation is a technique to modify the source code of an application in order to collect statistics at run-time. This may be used to measure performance of the application, or diagnose errors.

²A basic block in a program is a series of instructions which are executed sequentially. The basic block does not contain branch instructions.

Contents

Abstract	iii
1 Introduction	1
1.1 Simulation	1
1.2 Related Work	1
1.2.1 Sampling Based Approach	1
1.2.2 Host Compiled Simulation	2
1.3 Focus	2
1.4 Thesis Outline	2
2 Background	4
2.1 Basics of Computer Architecture	4
2.1.1 Pipeline Execution	4
2.1.2 Cache Hierarchy	5
2.1.3 Branch Prediction	6
2.2 ARM Cortex A5 Processor	7
2.2.1 Processor Pipeline	7
2.2.2 Cache Hierarchy	7
2.2.3 Branch Prediction Unit	8
2.3 Bare Metal Applications on ARM	8
2.4 Cross-Compiling using GCC	9
2.5 Compiler Optimizations	9
2.5.1 Partial or Full Loop Unrolling	9
2.5.2 Conditional Execution	9
2.5.3 Code Rearrangement	10
2.6 Extracting Debug Information from GDB	10
3 Host Compiled Simulation	13
3.1 The Technique	13
3.1.1 Simple Example	14
3.2 The Flow	16
3.3 Source Code to Intermediate Source Code	17
3.4 Mapping Algorithm	17
3.5 Annotation for Execution Cycles and Branch Prediction	20
3.5.1 Branch Prediction	20

3.6	Annotation for Memory Access	21
3.6.1	Cache Simulator	22
3.7	Instruction Cache Simulation	23
3.8	Data Cache Simulation	23
3.8.1	Memory Access Reconstruction	24
3.9	Extracting Debug Information from GDB	25
3.10	Data Cache Simulation	25
3.10.1	Memory Access Reconstruction	27
3.10.2	Annotation for Data Cache Simulation	27
3.10.3	Implementation of Cache Simulator	29
3.11	Instruction Cache Simulation	29
3.12	Annotation for Execution Time in Pipeline	29
3.13	Annotation for Branch Prediction	29
4	Implementation	30
5	Results	31
6	Conclusion	32
	List of Figures	33
	List of Tables	34

1 Introduction

1.1 Simulation

Simulation is the technique to imitate the behaviour of a system. It is generally used when developing or working with the real system is expensive or difficult.

Simulation is widely used in Hardware Software Co-development. The behaviour of a processor is simulated, and benchmarking applications are run. It allows Hardware Developers to analyse and validate the performance impacts of design decisions at early stage of hardware development. This saves the crucial effort and cost involved in the painstaking process of fabricating hardware at each milestone.

The popular approach used for simulation in the industry today is called Cycle Accurate Simulation (CAS). In CAS the processor micro-architecture is simulated in great detail. Each stage of processor pipeline is simulated along with other building blocks of the processor like Cache Memory and Branch Prediction Units. This approach provides cycle accurate estimates of performance. However, CAS is difficult to develop and slow to execute because of the amount of details that are simulated. Long running software benchmarks are used to accurately estimate bottlenecks in the system, and CAS is not suited for this due to the slow execution speed.

In this research, a technique for fast simulation of processors has been explored. The aim is to extract accurate estimate of number of cycles spent and power consumed by a processor in running a software benchmark.

1.2 Related Work

Research in this area has focussed on following approaches.

1.2.1 Sampling Based Approach

Sampling is an approach used in statistical analysis. Small, yet representative samples are chosen from a vast amount of data. These samples are analysed in detail, and the

results are interpolated to gather information about the entire data set.

In this approach, the application is mostly run using Functional Simulation, and some samples are executed using the detailed CAS. The number of cycles spent in execution of the samples is calculated, and the number of cycles spent in executing the entire pipeline is estimated by interpolating.

This approach provides considerable speed up compared to CAS, however accuracy of the estimation is highly dependent on how the samples are chosen. Also, developing this technique is difficult, since CAS is used.

1.2.2 Host Compiled Simulation

Host Compiled Simulation is based on the approach of Source Code Instrumentation (SCI). SCI is the process of modifying source code to collect performance statistics and generate trace information during run-time.

When an application is run on a processor, most of the time is spent in

- Execution of the instructions in the processor pipeline, and
- Fetching data from the memory.

If the number of cycles spent in each of these phases can be accurately estimated, the total number of cycles spent in running the application can be calculated. In this approach, the source code is instrumented to do this. The instrumented source code is compiled for and run on the Host Machine (the machine where the simulation is run) and hence the name, Host Compiled Simulation.

1.3 Focus

The focus in this research is to develop a tool to perform Host Compiled Simulation of a processor. The tool should be able to automatically instrument a given source code. The instrumented source code will be compiled and run, and accurate estimates of the performance will be reported.

1.4 Thesis Outline

Chapter 2 covers some basic background knowledge that will be helpful in understanding the project. This chapter can be skipped by the knowledgeable readers. References

to appropriate sections of the chapter is provided in the next chapters. The readers may choose to revisit these concepts if needed.

Chapter 3 presents an overview of the technique. A flow chart explaining the approach is described and the challenges in each stage are briefly illustrated using examples.

Chapter 4 deeply explains each stage of the technique in further detail. The challenges, and how they are tackled is discussed. It also sheds some light into the limitations of the tool and how it can be enhanced.

Chapter 5 presents a deep analysis of the results. The test setup is explained, and the accuracy of results from this technique is demonstrated.

2 Background

This chapter covers background knowledge that may be required to understand the project. The topics briefly covered in this chapter are as follows.

- Basics of Computer Architecture: Pipeline, Caching, Branch Prediction.
- Details of ARM Cortex A5 Processor
- Execution of Bare Metal Applications on ARM
- Cross-Compiling using GCC
- Compiler Optimization
- Extracting Debug Information using GDB

Basic understanding of these topics is needed to understand this thesis. Readers may choose to skip this chapter. Further in the thesis, appropriate references are made to sections in this chapter. Reader may choose to revisit these sections if required.

2.1 Basics of Computer Architecture

2.1.1 Pipeline Execution

Most modern processors use a multi-stage pipeline architecture for the execution unit. Execution of an instruction can be divided into multiple stages as follows:

- Instruction Fetch
- Instruction Decode
- Data Fetch
- Execute
- Write Back

This is only a general representation of the break down, and different processors use a slightly modified break down. The task performed in each of the stage is self-explanatory and usually, one clock cycle is spent in executing each stage. The execution unit of the processor is divided into separate units for each stage of the instruction

execution.

This allows for parallel execution of instructions. At the beginning of execution, the pipeline is empty and the first instruction is fetched from the memory. While this instruction is being decoded the subsequent instruction is being fetched and so on. This helps in achieving a performance of nearly 1 Cycle Per Instruction (CPI), when no branching instructions are encountered and all data is available in registers.

This performance may not always be achievable due to interlocking between instructions. Interlocking occurs when subsequent instructions have control or data dependency. If an instruction depends on the result of the previous instruction, the pipeline must be stalled unless the result is made available.

This performance is also hindered by conditional branching instructions. A simple *if-then-else* branch is implemented as a compare instruction followed by a conditional branch instruction in the binary code. Depending on the result of the compare instruction, the branch instruction is either executed or ignored. The processor uses Branch Prediction Algorithms to predict the next instruction that will be executed and already feeds the instruction to the pipeline. Incorrect prediction leads to Pipeline Flush, which results in a wastage of a couple of cycles.

2.1.2 Cache Hierarchy

Memory access is a critical bottleneck in modern computers. Main memories used today take over hundreds of cycles to fetch data, while the processor halts execution waiting for the data to be delivered. This leads to considerable waste of resources.

To optimize this, processors use a hierarchy of low-latency cache memories. The level 1 or L1 Cache is closely coupled with the processor. It has very low-latency for read and write operations, but is much smaller in size compared to the main memory. L2 Cache has higher latency than L1 Cache, but is much faster compared to main memory. Further levels of caching may be used, but it has been observed that the best performance is achieved by using 2 or 3 levels of caching.

Different techniques can be used to store and lookup data in a cache. The most popular technique is to use an N-way Set Associative Cache which has been explained graphically in figure [TODO]. The cache memory is divided into *n-ways* of equal size. Each *way* has a number of fixed size cache lines which can be indexed. Each entry in the cache contains a *tag*, the *data* and *flags* which tell the validity of the data.

When data on an address is requested by the processor, the address is divided into *tag* and *index*. The data can only reside in one of the cache ways in the cache line indexed by *index*. The cache controller looks for valid data in the cache lines, which has the same *tag* as the *tag* derived from the address. If such a data is found, a cache hit occurs

and the data is made available to the processor. If the data is not found, a cache miss occurs and lookup is carried out in further levels of cache. If none of the caches contain the data, it must be fetched from the memory. At this point, the data along with the tag is stored in the cache. If the same data is requested again, it is now available in cache and can be delivered with little latency.

For a write operation, different policies can be used. On receiving a write operation request from the processor, a cache with Write-Through policy will immediately write the data back to the main memory. The processor must remain in an idle state until this write operation is completed. In Write-Back policy however, the data will only be written to the cache, and it will be flagged as *dirty* data. The data will be written to the main memory later, and the processor does not need to wait, thereby improving performance. Most processors today use Write-Back Policy.

Separate caches are usually used to store Instruction and Data, as the access pattern for the two is generally quite different.

Obviously, there is a limit to the amount of data that can be stored in the cache, and hence the cache controllers use a cache replacement policy to replace data which is less likely to be used again. Popular replacement policies are Least Recently Used and Pseudo Random Replacement.

The memory management unit also tries to analyse the memory access pattern, and prefetches the data that the processor could use in future. This reduces the number of cache misses, and optimizes performance.

This is just a basic understanding of how caches work, in a computer with a single processor. In multi-processor systems each processor may have a private L1 Cache. Different processors may try to read and write to the same address simultaneously leading to race conditions. The cache controllers must ensure data coherency. This topic is not discussed here, because the focus for this research is on single-processor systems.

2.1.3 Branch Prediction

The issue of pipeline flushing in presence of branching instructions was discussed in section 2.1.1. Processors use Branch Prediction Algorithms to improve performance by trying to predict the result of the conditional branch instruction. On fetching a branch instruction, the Branch Prediction Unit tells the Instruction Fetch unit in the pipeline the address of the next instruction to be fetched according to the prediction. If the prediction holds true, a few clock cycles are saved on execution. If the prediction was incorrect, the pipeline is flushed.

The general policy for Branch Prediction utilizes a track record of all branch instructions

encountered previously. It uses a state machine (figure [TODO]) for each branch instruction, and maintains the state of the machine in the table using 2 bits as follows.

- 00 - Strongly Not Taken (SNT)
- 01 - Weakly Not Taken (WNT)
- 10 - Weakly Taken (WT)
- 11 - Strongly Taken (ST)

When in SNT and WNT state, the Branch Prediction Unit predicts that the branch will not be taken. Similarly, in WT and ST states it predicts that the branch will be taken. When a branch instruction is seen for the first time, the state is set to SNT. The state changes depending on whether the prediction was correct or incorrect.

2.2 ARM Cortex A5 Processor

ARM Cortex A5 has been used as the target processor to analyse the results in this research. In this section, features of the processor have been briefly described. Only the details relevant for understanding the thesis have been covered.

2.2.1 Processor Pipeline

The Cortex A5 has an 8 stage processor pipeline. Each stage in the pipeline takes 1 cycle to execute. Interlocking may occur between instructions, which may lead to pipeline stalls.

2.2.2 Cache Hierarchy

The Cortex A5 uses separate L1 cache for data and instructions. The size of the L1 Cache is configurable from 2 KB to 64 KB. On the target processor used for testing, the 32 KB of cache memory was used for Data and Instruction Caches. Each cache line has a fixed size of 32 Bytes.

The Data Cache is 4-way Set Associative Cache, with Write-Back Policy. The Instruction cache is a 2-way Set Associative Cache. The replacement policy used is Pseudo Random Replacement.

The L2 Cache is a unified Cache for data and instruction. It is 256 KB in size and is a 16-way Set Associative Cache with 32 Byte cache line length. It also uses Write-Back Policy, and Pseudo Random Replacement strategy.

The Cache controller ensures that each data is only cached in one of the caches. When a data is fetched from main memory, it is stored in the L1 Cache. When a data needs to be evicted from L1 Cache, it is stored in L2 Cache. In case of a L2 Cache Hit, the data is moved from L2 cache to L1 cache.

2.2.3 Branch Prediction Unit

The branch prediction unit uses the algorithm described in section 2.1.3. It maintains a 125 entry history table for the branches.

2.3 Bare Metal Applications on ARM

Bare Metal Applications are programs that are run on the processor without an underlying layer of Operating System. Running a bare metal application is quite different from running an application on top of an Operating System such as Linux.

The job of an Operating System is to manage the resources of the computer system, and provide a mechanism to load and execute programs at the users request. Operating Systems take control of the main memory of the system, and try to make the most efficient use of it by appropriately sharing it among different applications. When an application needs to be executed, the Operating Systems creates space in the main memory, to load the code and the data for the application.

The application code may be loaded at any address in the main memory, at the behest of the Operating System. Hence the code must be compiled in a way such that it is relocatable. The addresses used within the code are relative to the actual (physical) address where the application code will be loaded. During execution, each access to the relative address is translated by the Memory Management Unit, with support from Operating System, to the physical address.

In some sense, it is much more easier to understand how a bare metal application is executed. When a program is compiled to be run as bare-metal application, a specific address needs to be provided at which the program will be loaded. All memory addresses used in the binary code, will be actual (physical) addresses on the target processor. To run the application, the binary code needs to be loaded at the predefined start address. The processor starts executing the instruction the instruction at this address in a sequential manner, until a branch instruction is found. The processor then jumps to the address provided in the branch instructions and continues executing the code until the end of the program is reached.

It is important to note, that for each instruction and data in the bare metal application, the exact address where it will be located in the memory of the target processor can be

known. This information will be used later in the project.

Apart from this, there are some limitations that arise with Bare Metal Applications. Since there is no Operating System, interaction with input and output devices is very difficult. Reading and writing of files from the disk is not possible, because the logic to understand the file system is not available. Using heap memory for dynamic memory allocation, while possible, is complicated. Using Shared libraries in bare-metal applications is not possible, because the Operating System is needed to perform dynamic linking at run-time. All applications must be statically linked.

2.4 Cross-Compiling using GCC

2.5 Compiler Optimizations

The compiler uses a number of techniques to optimize the code, without changing the functionality. Often, this results in a substantially modified control flow in the binary code, compared to the flow in the source code.

Some of these optimization strategies have been discussed here. This will be useful in understanding some of the challenges that the project aims to solve.

2.5.1 Partial or Full Loop Unrolling

To reduce the number of branching instructions that need to be executed, the compiler checks if the number of iterations of a loop can be reduced by merging a few iterations together.

In cases where the number of iterations of the loop are known at compile time, the compiler may try to eliminate the loop by replicating the instructions inside the loop. This may be useful when the number of iterations of the loop is less, and the body of the loop contains only a few instructions. It eliminates the need of branching instructions.

2.5.2 Conditional Execution

Conditional Execution or Predication is a feature that is provided by the Instruction Set for the target processor. Instructions can be predicated with a condition. The instruction is fed to the processor pipeline, but the result of the instruction is only written back (committed) if the condition holds true.

Conditional Execution is used to eliminate conditional branching instructions. A simple

if-then-else construct is implemented in ARM binary code as a compare instruction, followed by a conditional branch instruction. The condition checks for the result of the compare instruction, and appropriately tells the processor to branch or continue executing the consecutive instruction. Since the result of the compare instruction will not be available until the last stage of the pipeline execution, the next instruction to be executed can not be known with certainty.

The Branch Prediction Unit, predicts whether the branch will be taken or not and accordingly tells the Instruction Fetch Unit to feed the appropriate instruction into the pipeline. If the prediction turns out to be incorrect, the instructions in the pipeline must be flushed and the correct instructions must be reloaded. On the ARM Cortex A5 with an 8-stage pipeline, this will lead to a loss of crucial 5-6 clock cycles.

Conditional Execution can be used to eliminate the use of branching instructions. This will lead to an improvement in performance, especially when the code in *then* and *else* body contained only 1 or 2 instructions. Each instruction in the *then* and *else* body is predicated with an appropriate condition. After the compare instruction, the predicated instructions are fed into the pipeline. The predicated instructions progress in the pipeline, meanwhile the result of the compare instruction is available. The result of the predicated instruction is only written back (committed) if the condition evaluates to true.

If the *then* and *else* bodies contain one instruction each, only one cycles is wasted by using Conditional Execution.

The example in listing [TODO] shows Conditional Execution. Figure [TODO] shows how the control flow of the binary code varies significantly from the control flow of the source code.

2.5.3 Code Rearrangement

The compiler tries to remove data and control dependency among consecutive instructions by rearranging instructions. This leads to reduced pipeline stalls and overall improvement in performance.

2.6 Extracting Debug Information from GDB

GDB is an open source debug tool which can be used to debug applications. It can be used to extract a lot of information. In this project, GDB will be used to extract information about the variables used in the code. This section briefly discusses the scope of information that can be collected from GDB.

For debugging ARM binaries, GDB specially compiled for the target must be used. The following command starts GDB and uses as input *example.elf* which is a binary of a Bare-Metal application compiled to run on an ARM processor. This gives the GDB prompt. GDB can be used for remote debugging on a target hardware, and it also provides a functional simulator for the ARM processor. The following commands are given to tell GDB to use a simulator as a target, and load the binary.

```
$ arm-none-eabi-gdb example.elf
(gdb) target sim
(gdb) load
(gdb) _
```

Often, many variables used in the source code are eliminated (optimized out) during compiler optimization. GDB can be used to extract information about the Global Variables being used in the binary code. It can provide the name, address and sizes of the global variables. To get this information, the following commands can be used.

```
(gdb) info variables
From example.c
unsigned int array[100];
unsigend int result;

(gdb) print &result
$1 = 0x7c8

(gdb) _
```

This shows that global variables *array* and *result* are declared in file *example.c*. It shows the address of the variable *result* using the second command.

Getting similar information about local variables is also possible. The local variables are only visible inside the scope of the function in which they are defined. We need to configure a break point in the beginning of the function, then run the simulation and issue following commands to extract information about local variables defined in the function.

```
(gdb) b main
Breakpoint 1 at 0x368: file example.c, line 10.

(gdb) info locals
i = 0
j = 0

(gdb) print &i
$1 = 0x17b40
```

This shows that local variables i and j are declared inside the function `main`. Address of i is printed, and it must be noted that this address resides in the stack, as i is a local variable.

3 Host Compiled Simulation

Host Compiled Simulation (HCS) is an important area of research in the space of fast simulation techniques for performance estimation. It is popular because it is simple to understand and develop. It is expected to be faster in execution compared to Cycle Accurate Simulators (CAS) and more accurate than other approaches used for Fast Simulation like the Sampling Based Approach.

In this chapter, the concept of HCS is illustrated using a simplified example. The steps involved in performing HCS are outlined. An overview of each step and the challenges is provided.

Detailed implementation of the technique and handling of challenges has been discussed in the next chapter.

3.1 The Technique

HCS is based on the approach of Source Code Instrumentation (SCI). Instrumentation is the technique to modify the source code of an application, so as to extract useful information during the run-time of the application. In HCS, the source code is instrumented to gain information of the time spent in executing the code on a particular Target Processor. The term Target Processor refers to the processor being simulated, and Host Machine is the computer on which the simulation will be run.

When an application is run on a processor, most of the time is spent in following phases.

- Execution of instructions.
- Fetching Data from memory, while the execution is stalled.

The technique is based on the assumption, that number of cycles spent in each phase of the execution can be accurately predicted using instrumentation. The predicted cycles can be accumulated during the run-time of the application to calculate the total cycles spent in execution of the program. The generated data can further be used to estimate the total power spent in the execution.

3.1.1 Simple Example

This simple example will be able to illustrate the concept of Host Compiled Simulation.

Consider the source code in Listing 3.1. The function `sum` calculates the sum of elements in an array and returns the result. The source code is cross-compiled to run on ARM Bare-Metal. The load address for bare-metal binaries is decided upon compile time. Listing 3.2 shows the object dump of the binary code. The address of each instruction is the physical address in the memory of the target processor, where the instruction will reside. Bare-Metal execution was briefly discussed in Section [TODO].

The instrumented code is shown in Listing 3.3, where the annotations are highlighted.

```

1  int sum(int array[20])
2  {
3      int i;
4      int sum = 0;
5
6      for (i=0; i<20; i++)
7          sum += array[i];
8
9      return sum;
10 }
```

Listing 3.1: Simple C Code

```

1  00008068 <sum>:
2  8068:    mov     r3, #0
3  806c:    mov     r2, r3
4  8070:    ldr     r1, [r0, r3]
5  8074:    add     r2, r2, r1
6  8078:    add     r3, r3, #4
7  807c:    cmp     r3, #80 ; 0x50
8  8080:    bne     8070 <sum+0x8>
9  8084:    mov     r0, r2
10 8088:    bx      lr
```

Listing 3.2: Objdump Code

```

1  unsigned int execCycles;
2  unsigned int memAccessCycles;
3
4  int sum(int array[20])
5  {
6      int i;
7      int sum = 0;
8      execCycles += 2;
9      memAccessCycles += simICache(0x8068, 8);
10
11     for (i=0; i<20; i++)
12     {
13         sum += array[i];
14         memAccessCycles += simDCache(&array + i, READ);
15         execCycles += 5;
16         memAccessCycles += simICache(0x8070, 40);
17     }
18
19     execCycles += 2;
20     memAccessCycles += simICache(0x8084, 8);
21     return sum;
22 }
```

Listing 3.3: Instrumented Code

In the instrumented code, two global variables **execCycles** and **memAccessCycles** have been declared on lines 1 and 2 respectively. **execCycles** will store the number of cycles spent in actual execution of instructions, when the processor is in active state. **memAccessCycles** will store the number of cycles spent in performing read/write operations to the memory.

From the binary code, 3 basic blocks can be identified. These blocks can be matched to corresponding blocks in the source code. The mapping is shown in the Table 3.1.

Basic Block in Binary		Matching block in Source	
BlockID	Lines	BlockID	Lines
1	1-2	1	3-4
2	4-8	2	7
3	9-10	3	9

Table 3.1: Mapping of Basic Blocks

For each basic block, cycles spent in executing the instructions needs to be annotated. For simplicity, let us assume that the processor has a single stage pipeline, and each instruction takes 1 cycle to execute when input data is available in the registers. The number of cycles spent in execution of each basic block is estimated, and the result is added to the global variable **execCycles** on lines 8, 15 and 19.

The object code contains a load instruction on line 4, which corresponds to loading of the elements of the array. To estimate the time spent in fetching this data, the cache hierarchy of the target system must be simulated. The cache simulator offers an API **simDCache** to simulate data cache access. It takes as parameters the **address** of the data and a **flag** to tell whether it is a read or write operation. The cache simulator returns the number of cycles spent in performing the memory access. This value is accumulated in the global variable **memAccessCycles**, as seen on line 14.

Also, the instruction cache access must be simulated. This is done at the basic block granularity. The cache simulator offers API **simICache** which takes as parameters **address** of the first instructions in the basic block, and **size** of the basic block in bytes. The cache simulator returns the number of cycles spent in fetching the instructions. The value is also accumulated in global variable **memAccessCycles**, as seen on lines 9, 16, and 20.

The values of **execCycles** and **memAccessCycles** are returned at the end of the simulation.

3.2 The Flow

Figure 3.1 shows a flow-chart depicting the stages involved in performing automatic instrumentation. Each stage is explained in further sections.

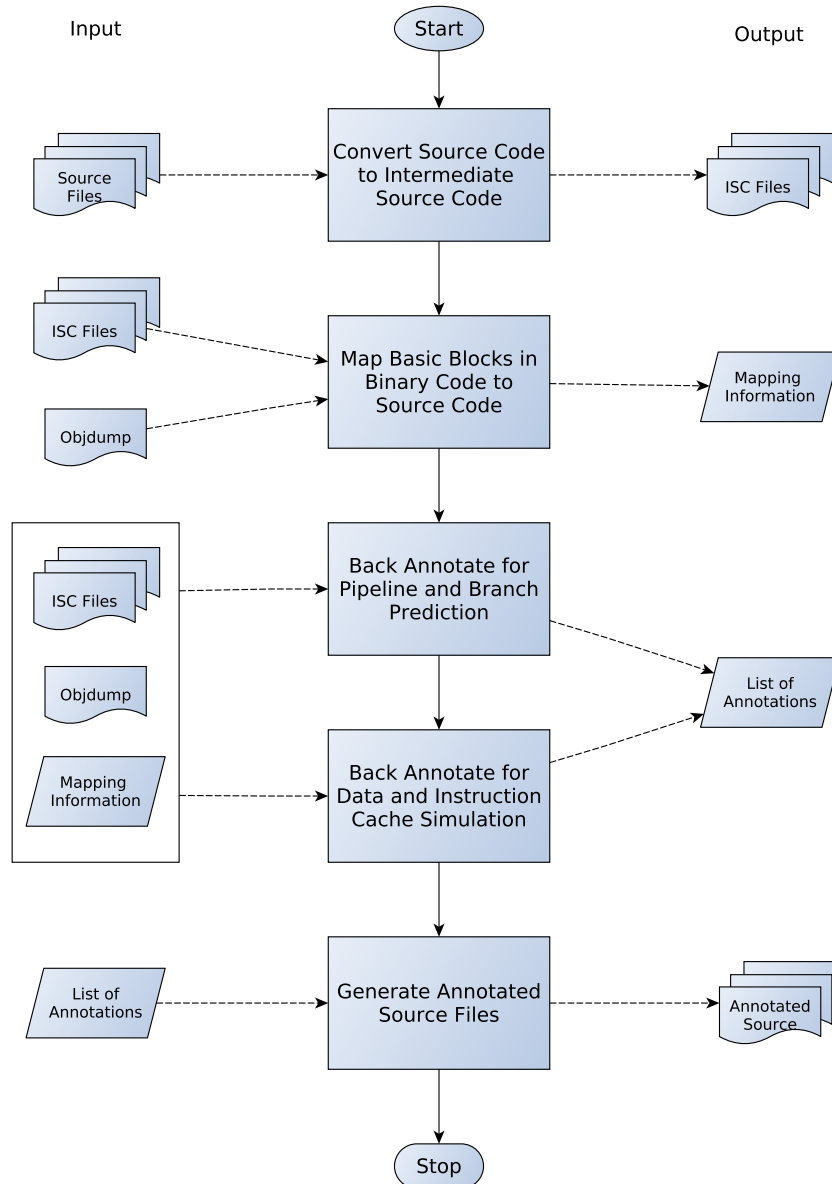


Figure 3.1: Flow Chart

In this section, a brief overview of the project architecture is provided. Each part of the project is independent, and has been treated as a separate problem. The flow chart in figure [TODO] shows each step of the project. The approach to solve these techniques

is presented in the subsequent sections in this chapter. Details of the implementation have been presented in the next chapter.

3.3 Source Code to Intermediate Source Code

From the example, it is clear that mapping between source code and binary code is very important for accurate instrumentation. Unfortunately, this mapping is destroyed during the optimization phases of the compiler. The compiler optimizes code by moving instructions around, eliminating branch instructions and modifying loops. Some of these optimization strategies have been discussed in Section [TODO].

The compiler performs these optimizations in two stages. The front-end of the compiler, translates the source code from a high-level language like C, to an Intermediate Representation (IR) called GIMPLE. The processor independent optimization strategies are applied to the IR code. In the back-end of the compiler, the optimized IR code is translated into Machine Language for the target processor. The processor dependent optimization strategies are applied in this phase.

The Intermediate Code is expected to have a control flow similar to the control flow of the Binary Code, since most optimization strategies have already been applied. It should be much easier to perform mapping between the IR Code and the Binary Code. However, instrumentation of IR Code is difficult.

In this project, the source code of the benchmark application is cross-compiled, and the optimized IR Code is translated back into a high-level language, C. The generated code is called Intermediate Source Code (ISC). The generated code is optimized version of the source code, with exactly similar functionality. It is easier to parse and instrument.

This approach is inspired from the research published in [TODO]. The python tool to translate GIMPLE code to C Code has been reused with minor modifications.

3.4 Mapping Algorithm

For accurate instrumentation, mapping at the Basic Block granularity is needed. GDB provides mapping information, but this has been observed to be highly inaccurate. This is a challenging problem and considerable research has been performed to develop algorithms for extracting accurate mapping.

The compilers employ a number of optimization strategies in which the control flow of the source code is modified. One such strategy is briefly discussed below, to illustrate the complexity of the mapping problem. Listing 3.4 shows a simple **if-then-else**

construct written in C. The code checks if **a** is greater than **b**. If true, the value of **a** is assigned to **max**, else the value of **b** is assigned to **max**.

Listing 3.5 is representative of how the assembly code may look like without any optimization.

```
1  int max(int a, int b)
2  {
3      int ret;
4      if (a > b)
5          max = a;
6      else
7          max = b;
8      return ret;
9  }
```

Listing 3.4: Example C Code

```
1  00008068 <max>:
2      8068:      cmp     r1, r0
3      806c:      ble     8078
4      8070:      mov     r0, r1
5      8074:      b       807c
6      8078:      mov     r0, r0
7      807c:      bx      lr
```

Listing 3.5: Object Code

Each instruction is loaded in the processor pipeline. While the compare instruction on line 4 is being executed, the branch instruction on the next line is loaded into the pipeline. Depending on the result of the compare instruction, the branch will be taken or not taken. Since the result of the compare instruction is not yet available, the move instruction on line 6 is loaded into pipeline by assuming that the branch will not be taken. However, if the content of register **r1** is lesser or equal to that of **r2**, the pipeline will be flushed and move instruction on line 8 will be executed. The pipeline flush leads to a loss of a few cycles.

To reduce this loss, the Instruction Sets of most modern processors support a mechanism known as Predication or Conditional Execution. Each instruction like **mov** or arithmetic instructions can be predicated with a condition. The instruction is loaded into the pipeline, but the result is only written back (or committed) if the condition evaluates to true. This feature can be used to eliminate conditional branch instructions. The optimized binary code is as shown in Listing 3.6.

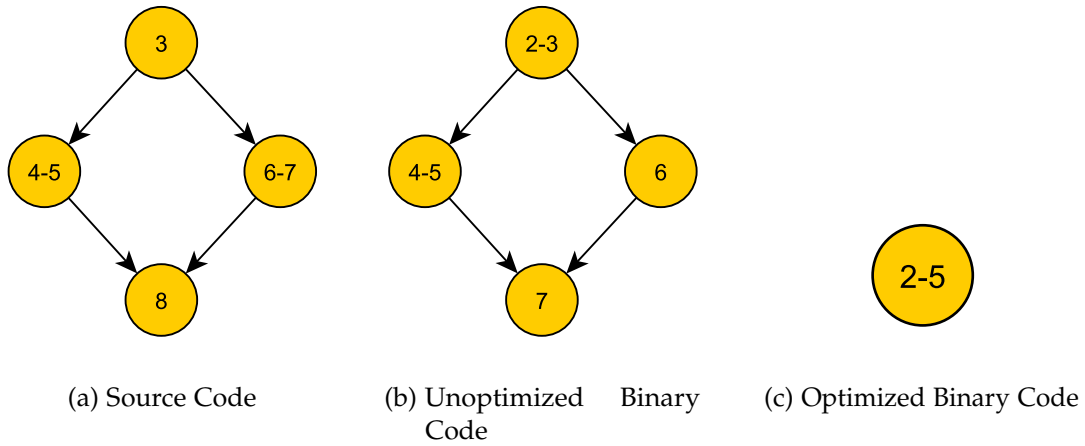


Figure 3.2: Control Flow Graphs

1	00008068	<max>:
2	8068:	cmp r1, r0
3	806c:	movge r0, r1
4	8070:	movlt r0, r0
5	8074:	bx lr

Listing 3.6: Optimized Object Code

Since the branching instructions are eliminated, the code in Listing 3.6 will be treated as a single basic block. The control flow graphs for the given source code, and the unoptimized binary code have been represented in figures 3.2a and 3.2b. They are similar, and hence easy to map. Figure 3.2c shows the Control Flow Graph of the optimized binary code. As expected, it is a single block.

Similarly, other optimization strategies may modify the Control Flow. To map the Control Flow Graphs, a standard graph matching algorithm using Depth First Traversal has been implemented, and special handling has been done for each optimization strategy used by the compiler. Additionally, mapping information from GDB is used in corner cases.

During compile time, developers can tell the compiler which optimizations to perform using the option **-O** followed by a number between 0 and 3 for the level of optimizations, 3 being the highest level of optimization. The tool currently provides promising results at optimization level 1, and can be extended to handle other optimization techniques.

The Control Flow Graphs generated from ISC and the Optimized Binary Code, along with the mapping information generated by the algorithm is stored in a data structure. The data is used in the next steps.

3.5 Annotation for Execution Cycles and Branch Prediction

The annotation for number of cycles spent in execution of the instructions is done at a basic block granularity. Cycles are estimated for each basic block in the binary code, and back annotated to the mapped block in the source code. A global variable **execCycles** is declared, and the annotated cycles for each basic block are accumulated in this variable.

Most processors use multi-stage pipelined execution units to improve the throughput (number of instructions executed in unit time). The execution cycle for instructions is divided into stages. Each stage takes one cycle to execute. The first instruction is fed into the first stage of the pipeline, and while it advances to the next stage, the next instruction is fed into the first stage of the pipeline. Pipeline Architecture has been explained in Section [TODO].

Some instructions have data and control dependency on the previous instruction. For instance, an instruction may need as input, the result generated by the previous instruction. This is called as Interlocking, and it leads to pipeline stall because the result of the previous instruction must be made available, before proceeding with the next instruction. The number of cycles for which the pipeline is stalled may vary, depending on the type of interlocking.

To accurately estimate the number of cycles spent in execution of a basic block, instructions are parsed sequentially and simulated on the pipeline. For each basic block, it is assumed that the pipeline is initially empty, which may not be the case and will be accounted for later. It is also assumed that all the data needed by the instructions is available with no latency. The cycles spent in fetching the data from memory will also be accounted later. Interlocking between instructions is identified and appropriate penalty is added.

The total number of cycles consumed in execution of a basic block in the binary code, is annotated into the matching basic block in the ISC. The annotated cycles for each basic block are accumulated in a **execCycles**.

3.5.1 Branch Prediction

To calculate the cycles spent in execution, it was assumed that the pipeline is empty at the start of each basic block. This may not always be the case. This will be taken into account here.

Conditional branching is implemented in binary code as a compare instruction followed by a conditional branch instruction. The branch is taken, depending on the result of the compare instruction. Since the result of the compare instruction may not be available

in time, the instruction to be fetched after the branch instruction cannot be known with certainty. The processor predicts an outcome of the compare instruction, and accordingly whether the branch will be taken or not. The appropriate instructions are loaded into the pipeline. If the processor predicted the result of the branch instruction incorrectly, the subsequent instructions may need to be flushed. This is known as a Pipeline Flush, and leads to the waste of a few cycles.

To minimize the occurrence of pipeline flush, most processors implement a sophisticated Branch Prediction Unit (BPU), that uses heuristics to predict the outcome of a branch instruction. To accurately account for the cycles spent in switching between basic blocks, the Branch Prediction Algorithm needs to be simulated.

The Branch Prediction Algorithm simulated in this project uses a table to maintain the history of branch instructions seen in the recent past. The address of the branch instruction is stored along with the outcome. When a branch instruction is loaded into the pipeline, the BPU checks if the information for the branch is present in the history table. If an entry is found, the prediction is made based upon the outcome stored in the table. If an entry is not found, the processors predicts that the branch will not be taken. Eventually, when the outcome of the instruction is known the history tables are updated. The pipeline will be flushed if the prediction was incorrect, but wastage of cycles will be prevented if the prediction was correct.

To simulate this, the source code is annotated to call the Branch Prediction API at the beginning of each basic block. The Branch Prediction Simulation keeps a track of the block which was being executed previously. By simulating the algorithm, it checks whether the branch from the previous block to the current block would have been predicted correctly by the BPU, and returns the result. If the branch would have been correctly predicted, the pipeline would not be empty at the beginning of the basic block, and hence a few cycles are subtracted from **execCycles**. Otherwise, a pipeline flush must have occurred and the assumption holds true.

3.6 Annotation for Memory Access

The cycle spent in fetching data and instructions from the memory must be accounted. To do this, a Cache Simulator has been implemented. Annotation is made to the source code to call the API offered by the Cache Simulator, in order to simulate the memory access at run-time.

3.6.1 Cache Simulator

For accurate estimation, the memory system of the target processor needs to be simulated. The cache simulator keeps track of the data stored in the caches. Whenever a memory access is made, the simulator checks whether the data would have been available in either of the caches. If the data is not found in any cache, the data must be fetched from the memory. The simulator returns the number of cycles that would have been spent on the target processor to perform the memory access.

Processors can use a variety of Cache Organizations, and features. The Cache Simulator has been developed keeping this in mind.

Most processors use separate cache for storing Data and Instructions. The cache simulator offers the following API to simulate the different types of memory access.

```
/**
 * @brief Function to simulate Data Cache Access.
 *
 * @param Address of the data to be fetched
 * @param True, if read access
 *        False, if write access
 *
 * @return Number of cycles spent in performing access.
 */
unsigned long long simDCache(unsigned long address,
                             unsigned int isReadAccess);

/**
 * @brief Function to simulate Instruction Cache Access.
 *
 * @param Start Address of the basic block.
 * @param Size of the basic block in Bytes.
 *
 * @return Number of cycles spent in performing access.
 */
unsigned long long simICache(unsigned long address,
                             unsigned long size);
```

For every data access in the application, annotation is made in the source code to call **simDCache** and trigger Data Cache simulation. The address of the data being fetched is provided as parameter, along with a flag to signify whether it is a read or write access. Simulation of Instruction Cache is done at basic block granularity. **simICache** is called, with the start address of the basic block and the size of the basic block in bytes as parameters.

Caches may have varying parameters and features such as,

- Sizes
- Approach for Associativity of Data (Direct Mapped, N-way Set Associative)
- Hierarchy (multiple levels of caching)
- Replacement Policies

The implementation of the cache simulation, is independent of the API and can be easily plugged using a portable interface mechanism. This allows easy maintenance and modification of the Cache Algorithms, to analyse the impact on performance.

3.7 Instruction Cache Simulation

As previously stated, for estimating cycles spent in fetching the instructions from memory, instrumentation for instruction cache simulation is performed at the basic block granularity.

The basic blocks in the cross-compiled binary code are analysed. The address of the first instruction in each basic block is extracted from the cross-compiled binary. Note that since the binary is compiled to be run on Bare-Metal, the load address of the binary was decided at compile time. This means, that the address at which each instruction will be located in the memory of the target processor can be estimated. The size of the basic block is also extracted by analysing the basic block.

The corresponding basic block in the source code is known from the mapping information. Instrumentation is performed at the beginning of the basic block, and the instruction access is simulated using the API provided by the cache simulator.

3.8 Data Cache Simulation

The example shown in Section [TODO] was simplified for illustration, but had a flaw which would result in inaccuracy of estimation. To simulate load operation to fetch elements of `array`, the address of `array` from the Host Machine was used. For accurate estimation, data access must be simulated using target addresses. This is especially important since the Host and Target Memory Systems may vary significantly. For instance, the Host and Target system may use different sizes for basic data types like integer. This will lead to severe inaccuracy when accessing a sequence of elements from an array of integers.

However, resolving the addresses of load/store instructions cannot be done statically. The process to do this is inspired from research published in [TODO] and is called Memory Access Reconstruction.

3.8.1 Memory Access Reconstruction

The first step in the process of memory access reconstruction is to find the physical addresses where each variable will be stored in the memory of the target system. An application may use different types of memory. Separate approach is needed to resolve addresses of each type of variable.

Global Variables

Global Variables are accessible by all functions, and are stored in the Data Section of the application memory. The physical addresses of the global variables are decided at compile time. These addresses can be extracted by static analysis of the cross-compiled binary using GDB.

For each Global Variable, another variable is declared to store its physical address from the target system. This address is later used for simulating access to the global variable. Refer to the example [TODO].

Local Variables

Local Variables are defined inside a function definition, and are only accessible inside the function. Memory for Local Variables is only allocated when the function is called, and is located in the stack frame of the function. The actual physical address of the local variables can not be resolved by static analysis, as the stack grows and compacts at run-time. However, the address of each local variable relative to the value of the stack pointer, can be known.

If the value of the stack pointer is known, the relative address can be added to it to resolve the physical address of the local variable. Whenever a function is called, the stack frame of the function is pushed to the stack. The stack frame is popped when the function returns. To keep track of the stack frame at run-time, a global variable **CSIM_SP** is declared and initialized with the initial value of the stack pointer. The stack frame size of each function is extracted by static analysis of the binary. **CSIM_SP** is incremented at the beginning of each function by the stack frame size of the function.

Example [TODO] illustrates the annotations needed to resolve the physical address of local variables.

Dynamically Allocated Memory

Dynamically Allocated Memory is stored in the Heap Section. To allocate and free heap memory, the application uses API provided by system libraries. The memory allocation algorithm of the target system needs to be emulated at run-time of simulation, to resolve physical addresses of dynamically allocated memory. This approach is discussed in [TODO]. However, this is complicated to achieve, and has been ignored for this project. Only benchmark applications that do not use Dynamically Allocated Memory can be used. The project can later be extended to include this functionality.

Pointers as Function Parameters

The address of each Global and Local Variables has been resolved, however when pointer to a variable is passed as a parameter to a function, the address which is being accessed can not be deduced accurately. To handle this the signature of each function which takes pointers as parameters is modified. Another parameter is added, to pass the address of the variable on the target system.

3.9 Extracting Debug Information from GDB

As discussed briefly in the example, each load/store instruction in the binary code must be matched to a variable in the source code to simulate Data Cache Access. To be able to do this, information about each variable used in the code must be extracted. GDB has been used to extract this information.

Section ?? describes how GDB can be used to extract names, addresses and sizes of each Global and Local Variable. GDB also offers a mechanism to use scripts, which operate a series of commands and generate an output, which can be later parsed to gather the information required. This mechanism of GDB has been used.

The tool automatically generates GDB scripts and runs them, to extract information about the variables used in the program. For each variable the following information is collected.

- name of the variable
- scope of the variable. Empty for Global Variables, name of the containing function for local variables.
- address of the variable. For Global Variables it is the physical address. For local variables, it is the address relative to the Stack Pointer.

- type of the variable.
- size of the variable.

3.10 Data Cache Simulation

To estimate total time spent in fetching data from the memory, the tool performs detailed Cache Simulation. Each load/store instruction in the binary code is matched with to an instruction in the source code. The address from which the data is being fetched is identified, and memory access for the address is simulated.

Generally speaking, the host and target processors may use a different memory layout. For example, in some processors memory needs to be allocated in an aligned fashion for better performance, however this may not be the case for another processor. Also, the sizes for the basic data types may differ among different architectures. Size of an integer in one architecture may be 4 Bytes and in another may be 2 Bytes. For accurate cache simulation, the address at which the variable resides in the Host Machine, can not be used. Instead, each memory access as it would occur on the target processor must be simulated.

The challenge lies in extracting the address from which the data is being fetched. This address can not be easily extracted from the cross-compiled binary by static analysis. To extract this address, the tool implements a mechanism to reconstruct each memory access as it occurs on the target processor. The approach is based on the research published in [TODO].

An application uses memory to read and write input and output data. For simplicity, let us consider how an application written in C Programming language uses memory. The memory can be used by the application in 3 ways.

- **Global Memory.** Data stored in Global Memory is accessible by all functions in the program. The memory is stored in a fixed size Data Section. Size of each Global Variable must be known at compile time, so it is called statically allocated memory. In bare metal applications, the physical address of each global variable is known after compilation, and can be extracted from the binary.
- **Local Memory.** Each function can define variables which can only be used inside the function. The content of local variables is stored in the stack frame of the function. The size of the variable must be known at compile time. The addresses used by the local variables at run-time can not be known by static analysis, however the address relative to the Stack Pointer can be extracted.
- **Heap Memory.** Applications can also allocate memory at run-time. The content of this dynamically allocated memory is stored in a special section known as Heap. The heap can grow and contract at run-time. The address and sizes of this

type of memory can not be known statically.

The current implementation of the tool, only focuses on simulating access from the Local and Global Memory. Cache Simulation for Heap Memory is complicated, because the memory allocation algorithm used in the target processors have to be emulated to know the exact address where the memory will be allocated. This means, that the tool can only be used with Benchmark applications that do not use Dynamic Memory. This does not limit the importance of this tool since the goal is to simulate for performance analysis, and not functional verification.

The following section explains how each load/store instruction in binary code is mapped to a variable in the source code. Further, the approach to reconstruct the memory access is explained. Implementation of the cache simulator is explained in brief.

3.10.1 Memory Access Reconstruction

The tool parses binary code and emulates each instruction. It maintains the state of registers and updates it, as per the instructions. Branch instructions are ignored. The load/store instructions use register addressing modes to access data. From the content of the register the address of memory being accessed can be known. The memory being accessed may be an array, and this can not be clearly identified at this point.

The addresses of the Global Variables are known. Also, the addresses of Local Variables relative to the Stack Pointer are known. Using this information, the address of the load/store instruction found by emulation can be mapped to one of the variables.

Variables accessed in each basic block of the binary code are recorded, and this information will be used to map the load/store instruction to a specific line in the source code. For each basic block in the binary code, the basic block in the source code is parsed to find the line that causes the memory access. This is needed because for accessing an array, the index to be added to the base pointer can only be extracted from the source code.

Once this information is known, the memory address for each access can be reconstructed.

3.10.2 Annotation for Data Cache Simulation

Global Variables

For each Global Variable, say "*var*" of any data type, another global variable "*var_addr*" of type "*unsigned long*" is declared in the instrumented code to store the address where "*var*" will be held in the target memory. This address was extracted using the method described in Section 3.9.

The line in the source code where the global variable is being accessed is identified using the technique presented in Section [TODO]. Instrumentation to simulate cache is added after this line. The Cache Simulator is implemented in C language. It offers the following API to simulate Data Cache Access.

```
/**
 * @brief API Function to simulate Data Cache Access
 *
 * @param address The address being accessed.
 * @param isReadAccess Flag to indicate type of access.
 *                True, if Read Access.
 *                False, if Write Access.
 *
 * @return number of cycles spent in performing the memory access.
 */
unsigned long long simDCache (unsigned long address,
                             unsigned int isReadAccess)
```

The source code is appropriately instrumented using the above API. For accesses to elements in an array, the index multiplied by the size of the data type is added to the base address. The return value is the number of cycles spent in performing the memory fetch. This value is accumulated in a global variable, in a similar way as shown in the simple example above.

```
1  int result;
2  unsigned long result_addr = 0x88ac;           // <--
3  int input_array[20];
4  unsigned long input_array_addr = 0x88b0;      // <--
5
6  void foo()
7  {
8
9
10 }
```

Local Variables

The approach for simulating access of Local Variables is quite similar to the approach used for global variables. A new local variable is declared for each local variable used in the function, with a suffix "_addr" and data type "unsigned long". This variable contains the address of the local variable, relative to the current stack pointer. To accurately estimate the physical address where the memory resides, the value of the Stack Pointer is needed.

The stack grows and contracts during the run-time of an application. Whenever a function is called, a stack frame is created and the stack pointer is incremented. The stack frame contains the values of the function parameters, the local variables used in the function and the return address for the function. The size of the stack frame for each function can be extracted from the binary, and the start address of the stack is fixed at compile time.

To maintain the value of the stack pointer, a global variable "CSIM_SP" is added to the source code, which is initialized to the start address of the stack. At the beginning of each function, the value of "CSIM_SP" is incremented by the size of the stack frame for the function. The address of the local variable, relative to the stack pointer is added to "CSIM_SP", to calculate the physical address of the local variable, as would occur on the target processor.

The memory access is simulated using the same API as above.

Function Parameters

A func

Register Spilling

3.10.3 Implementation of Cache Simulator

3.11 Instruction Cache Simulation

3.12 Annotation for Execution Time in Pipeline

3.13 Annotation for Branch Prediction

4 Implementation

5 Results

6 Conclusion

List of Figures

3.1	Flow Chart	16
3.2	Control Flow Graphs	19

List of Tables

3.1 Mapping of Basic Blocks	15
---------------------------------------	----