

MASE: A Novel Infrastructure for Detailed Microarchitectural Modeling

Eric Larson
larsone@eecs.umich.edu

Saugata Chatterjee
saugata@sandcraft.com

Todd Austin
austin@umich.edu

Advanced Computer Architecture Laboratory
Electrical Engineering and Computer Science
University of Michigan

Abstract

MASE (*Micro Architectural Simulation Environment*) is a novel infrastructure that provides a flexible and capable environment to model modern microarchitectures. Many popular simulators, such as SimpleScalar, are predominately trace-based where the performance simulator is driven by a trace of instructions read from a file or generated on-the-fly by a functional simulator. Trace-driven simulators are well-suited for oracle studies and provide a clean division between performance modeling and functional emulation. A major problem with this approach, however, is that it does not accurately model timing dependent computations, an increasing trend in microarchitecture designs such as those found in multiprocessor systems. MASE implements a micro-functional performance model that combines timing and functional components into a single core. In addition, MASE incorporates a trace-driven functional component used to implement oracle studies and check the results of instructions as they commit. The check feature reduces the burden of correctness on the micro-functional core and also serves as a powerful debugging aid. MASE also implements a callback scheduling interface to support resources with non-deterministic latencies such as those found in highly concurrent memory systems. MASE was built on top of the current version of SimpleScalar. Analyses show that the performance statistics are comparable without a significant increase in simulation time.

1. Introduction

Computer system simulation is a vital technology in the computer system design cycle. The flexibility to quickly update software simulation models speeds the evaluation of design changes, permitting architects to explore large portions of the design space. Software modeling infrastructure also decouples hardware and software design efforts so that software development may proceed in parallel with hardware design, thereby reducing time to market for products with hardware and software components.

A microprocessor performance model is a software representation of a hardware design. It tracks the timing of instructions and data through the processor pipeline and memory system. Very detailed performance models may also include I/O device models such as disks and network interfaces. It is important for the performance model to be closely matched to the hardware that is being emulated. An inaccurate model can lead to incorrect or misleading research results [7].

Figure 1 illustrates a trace-driven modeling infrastructure, the most prevalent simulator organization. The performance model is driven by an instruction trace that represents the dynamic stream of instructions executed for a

specific processor architecture and workload. Traces are either read from a file [23], created by instrumented hardware [1], or generated on-the-fly by emulating a program [15]. SimpleScalar [6], SMTSIM [24], and VMW [8] are some examples of trace-driven simulation infrastructures. We use a broad definition of trace-based simulation. We consider any simulation environment that decouples creation of the dynamic instruction stream from the model that computes timing to be trace-based.

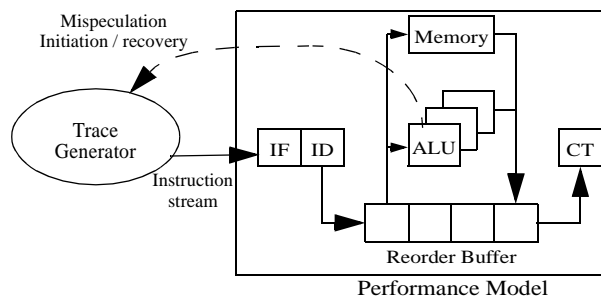


Figure 1: Organization of a trace-driven simulation environment. Instructions are supplied by a trace generator to a performance model that represents a detailed microarchitecture.

Techniques that use external trace files lack a functional simulation component. Since most trace files only contain the non-speculative instruction stream, performance models based on trace readers typically do not model mispeculation, that is, the instructions that are executed in the shadow of mispredicted branches, addresses, and instruction values. This can create large inaccuracies in the performance models; recent studies have shown that mispeculation streams provide instruction and data prefetching that is beneficial [9,21]. Moreover, mispeculation reclamation techniques such as instruction reuse [19] and the misprediction recovery cache [5] cannot be modeled using external trace files. One advantage of external traces created from real hardware is that they can include operating system code and interrupt routines. Embedding this information into a trace allows these segments to be included without building a full-system simulator.

Dynamic trace generation provides the additional flexibility to model mispeculation by providing an interface for the performance model to direct the instruction emulator to compute mispeculation (as shown by the dashed line in Figure 1). The primary advantage of this simulation approach is that it provides a clean division of infrastructure for performance modeling and instruction set emulation (or trace processing). However, there are many drawbacks to this infrastructure that make it increasingly

difficult to accurately model sophisticated, new microarchitectures.

Instruction streams can become inaccurate when they contain timing dependent computation, that is, computation whose inputs (and thus the result) are dependent on *when* the instruction executes. An example of a timing dependent computation is an instruction that is subject to value prediction in the execute stage.¹ We assume that the scheduler of the processor will schedule instructions that have all of their operands first. If there are available functional units after all of the ready instructions are executed, the scheduler will attempt to predict the outcome of an instruction (at the execute stage) based on previous instances of the instruction. This approach would seem to be a good one. Accessing the value prediction table as late as possible would result in higher prediction accuracy. Moreover, only accessing the table when needed would result in less power. The result of the instruction is dependent on when the instruction executes because it may get its value from the value prediction table or from its proper inputs. It is impossible to tell what the precise value is until the instruction is actually executed. Furthermore, if the value prediction is incorrect, the incorrect value is important because it may be an address that subsequently misses in the cache which could further affect the running time of the program. Accurate modeling of the mispredicted path is important since it competes for resources with the non-speculative instruction stream, which may cause additional cache misses or prefetch requests and affect the state of various predictors. It is our belief that the trend in microarchitectures is toward more dependence and speculation optimizations such as value prediction [12], instruction reuse [19], and speculative value coherence [11], making the accurate modeling of timing dependent computation crucial to design-time emulation.

The alternative to trace-based techniques is to employ a *micro-functional* performance model that not only times the activities of the program but also executes the program at that time, thereby reproducing the execution and timing of the program in a fashion identical to the simulated hardware. This approach can lead to more accurate models, an observation that is well recognized by architects and simulator developers [7]. However, most high-level simulation infrastructures remain trace-based due to two important drawbacks in micro-functional performance models. The first drawback comes from the coupling of timing and correctness. If the micro-functional performance model is incorrect in any way, the simulation can fail. While this exposes errors more readily in the functional model, this is not always the ideal design scenario. Inaccuracies may be the result of ongoing simulation development where details are left out because they were determined to be secondary. For instance, the forwarding of partial store values (*e.g.* a byte store forwarded to a word load) is a case that is complicated to get completely correct but happens so infrequently the overall

impact on performance is negligible. Trace-driven models are more tolerant of infrequent inaccuracies because they don't fail, and if the inaccuracies are infrequent, the overall simulation remains accurate. The second drawback with micro-functional performance models is that the simulation approach does not lend itself to oracle studies. Instruction inputs, results, and next PC values are not known until the instruction executes. For oracle studies, such as perfect branch or address prediction, not having this information earlier in the pipeline prevents performance bounds studies.

In this paper, we present MASE (Micro Architectural Simulation Environment), a novel performance modeling infrastructure that is built on top of the popular SimpleScalar toolset [6]. We address many of the drawbacks present in SimpleScalar that make it difficult for researchers to accurately model features present in complex high-performance microarchitectures. The goal of MASE is to provide a flexible infrastructure to researchers that they can use to create accurate models of the hardware they are studying. We provide four enhancements to the current implementation of SimpleScalar: (i) an oracle and checker that allows performance bound studies and removes the burden of correctness from the simulator core, (ii) a micro-functional core that increases modeling accuracy, (iii) fine-grain state management facilities that simplify the implementation of control and data speculative optimizations, and (iv) an interface to support resources with non-deterministic latencies.

MASE possesses all the benefits of trace-based modeling, by decoupling model accuracy from simulator correctness and providing extensive support for oracle studies. We implement our new simulation strategy by combining a micro-functional performance simulation infrastructure with a trace-driven oracle execution unit. The oracle execution unit executes instructions at the front end of the simulated processor pipeline, producing instruction information suitable for directing perfect speculation and other oracle studies. The burden of correctness is lifted from the performance model through the use of a checker execution component at the retirement stage of the performance model. Instruction results are checked as they retire into the architected state of the machine, and if they do not match those computed by the oracle, the performance model is flushed and restarted with correct simulation state. The checker provides a powerful model validation mechanism as well as a backup source of reference semantics that incomplete or inaccurate performance models may rely on to correctly complete any instruction. At the core of MASE is a micro-functional performance model; instructions are not only timed but executed in the core, accurately modeling timing-dependent computation.

We added to SimpleScalar a flexible speculative state management facility that permits restarting from any instruction. The current version of SimpleScalar only allows a single branch instruction to mispeculate and force a restart. The ability to restart from any instruction allows optimizations such as load address speculation and value prediction to be implemented. In these optimizations, instructions other than branches could be mispeculated, making it necessary to restart at the offending instruction. This approach also simplifies external inter-

1. This example is based on work done by Calder *et al.* [4]. The main difference is that the value prediction table is accessed in the fetch stage in their implementation while it is accessed in the execute stage in our example. We tried to find a real example of timing dependent value prediction from the literature but could not. The lack of any example is likely the result of deficiencies in existing simulation infrastructures.

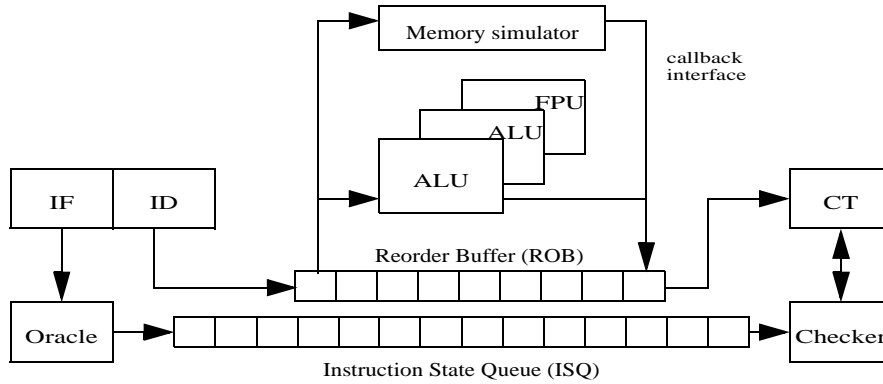


Figure 2: Block diagram of the new performance model infrastructure. The oracle executes instructions in advance and the results are stored in the ISQ so they can be checked by the checker when the instructions commit. The memory interface supports a callback interface to handle non-deterministic memory latencies. The functional units execute the instruction instead of just returning a latency.

rupt handling since any instruction could follow an interrupt request, forcing a rollback. The checker also uses this mechanism to recover from any errors that are detected since any instruction could potentially cause an error.

In addition, MASE addresses another deficiency in SimpleScalar by providing an interface to model components with non-deterministic latency. Modern DRAM systems can reorder requests to reduce the overall access time, allowing later requests to affect the access time of the current request. In MASE, a callback interface is used that allows the memory system (or any resource) to invoke a callback function once the memory system has determined an operation’s true latency. The callback interface provides for a more flexible and accurate method for determining the latency of non-deterministic resources. Several simulation infrastructures, such as RSIM [16] and SMTSIM [24], already provide this capability.

The remainder of the paper is organized as follows. The design and implementation of MASE is described in Section 2. Section 3 gives results from detailed analyses of our implementation. This was done by comparing the MASE model against the SimpleScalar baseline performance model (sim-outorder), comparing its accuracy and simulation speed. We also present a case study that explores the implementation of a blind dependence speculation technique, a speculation technique that cannot easily be implemented without these new facilities. Section 4 discusses related work. Section 5 concludes and suggests additional work.

2. MASE modeling architecture

This section describes the architecture of MASE’s performance model. A high-level view of the new architecture is shown in Figure 2. The following sections describe each of our key additions in more detail: the oracle and dynamic checker, the modernization of the performance model, the micro-functional performance model, the ability to restart from an arbitrary point, and a callback interface that supports non-deterministic resource latencies.

2.1 Oracle and checker execution component

The oracle sits in the fetch stage of the pipeline and exe-

cutes instructions in program order. The oracle is a functional emulator as it does not model any timing. The oracle contains its own register file and memory. To minimize the overhead, the oracle does not have a separate copy of memory but uses a hash table that contains all new values that have not been committed to architectural memory (all store values currently in the pipeline). Oracle loads are implemented by initially looking at the memory table. If there is a hit in the table, the value from the table is used. If there is a miss, architectural memory is accessed to obtain the value. Stores are executed by simply adding an entry to the oracle memory table. When an instruction is executed, the oracle state is updated and the result of the instruction is stored in the instruction state queue (ISQ) as shown in Figure 2. This queue serves several purposes. It holds data computed by the oracle that is used by the checker to verify that the instruction executed correctly. The data can also be used to facilitate performance bounds studies where correct data is needed. The queue serves as a record of executed instructions in case a recovery is necessary due to a branch misprediction or some other event. The oracle must always be synchronized with the fetch stage of the microarchitectural model, as a result, if instructions are flushed from the microarchitectural model, the oracle must also flush state. Since any instructions could potentially cause an error within the checker, the oracle and microarchitecture must have the ability to rollback from any instruction as described in Section 2.4.

Figure 3 shows an example of how the instruction state queue and oracle state is updated. The head of the queue contains the oldest instruction in the machine (the next instruction to commit) and the tail of the queue points to the next available entry. For each address in the oracle memory hash table, there is a linked list of values with the most recent store at the head of the list. In the example, there are two stores to address 0x500. The most recent store to this address wrote the value of 5; so it appears first. It is not sufficient to overwrite the previous value of 8 because the branch instruction between the two stores may mispredict causing a removal of any entries that occur after the branch (the store of 5 to address 0x500 in this case). When the store to address 0x514 is committed, architectural memory is updated and the corresponding

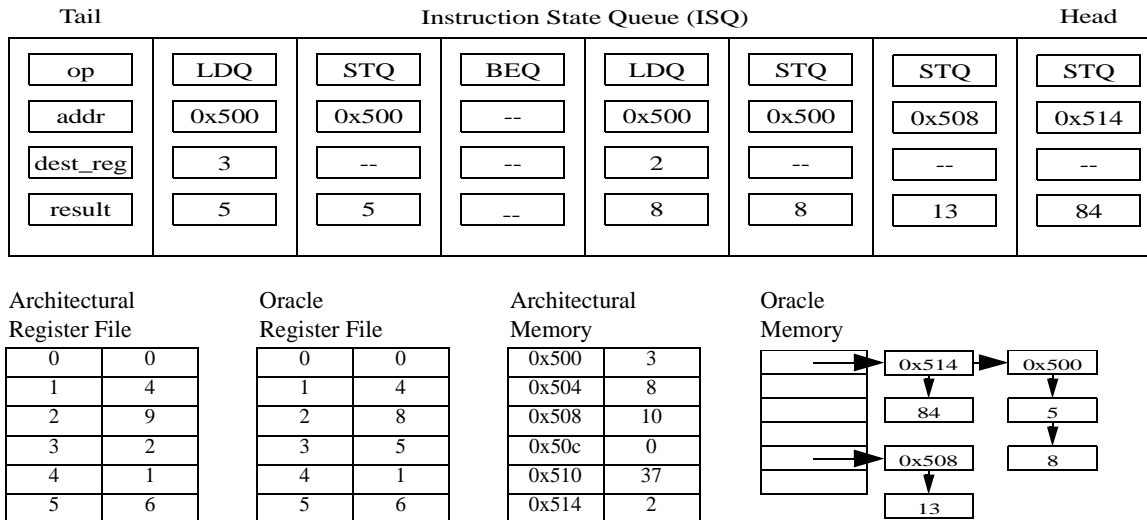


Figure 3: Example of oracle state and instruction state queue. The instruction state queue holds the instructions that currently reside in the pipeline and is used to check the results when the instruction is committed. The oracle memory is stored as a hash table and only contains values for stores that are in the pipeline.

entry is removed from the hash table. If a load attempts to access an address not in the hash table (address 0x504 for example), it will initially look in the hash table. After it discovers it is not there, architectural memory is accessed.

The oracle register file is stored in an array. It contains the latest write to all registers. In the example in Figure 3, register 3 has a value of 5. Unlike memory, it contains values for all of the registers, so it is only necessary to access the oracle register file. Mispredictions are also handled differently. If a branch is mispredicted, the architectural register file is copied into the oracle register file and the instruction state queue is scanned, from the head to the mispredicted branch, looking for instructions that write registers. If an instruction writes a register, it will update the register in the oracle register file with the appropriate value. This scanning process can recreate correct register state for any point in the dynamic instruction stream.

Oracles are commonly used to provide “perfect” behavior to do studies that measure the maximum benefit of an optimization. A common case of this is perfect branch prediction where all branch mispredictions are eliminated. In order to provide this capability, the oracle resides in the fetch stage so it knows the correct next PC to fetch. When executing an instruction, the oracle also serves as a decoder in the sense of extracting various information about the instruction such as the type of instruction or branch target. To minimize the running time of the simulation, the oracle saves the decoding information and passes the information along to the dispatch/decode stage.

The checker monitors all instructions that are committed, correcting any incorrect results due to model bugs or inaccuracies. The results of an instruction are passed along with the instruction until it is committed. The checker will compare these results with the results obtained by the oracle in the front-end of the machine. If the results match, the result will be committed to architectural state and the simulation will progress as normal. If the results do not

match, the oracle result will be committed to architectural state and a recovery will be initiated. The remaining instructions in the pipeline are flushed and the front-end is redirected to the next instruction. The instruction with the bad result is allowed to commit (with its result corrected) in order to ensure forward progress. This policy is used because, depending on the nature of the error, the bad instruction may repeatedly get the same error if it is re-executed, causing livelock in the simulation. Some instructions are non-speculative and have non-deterministic values, i.e. shared memory loads that are subject to race conditions. The results of these instructions can differ with the checker. These instructions are marked as such and if their result differs from the checker, the value from the micro-functional core will be regarded as correct and used to synchronize the checker.

The checker may be used in several different ways. The first way is to verify that any changes or enhancements to the simulator code are indeed correct. Since not all errors directly cause an error in the output, it provides extra security that a model enhancement did not violate any microarchitectural dependencies or program semantics. Another advantage is to allow the checker to handle tricky infrequent corner cases to save programming time. For example, it is difficult to program all of the cases involving partial store forwards where the base addresses are different. Instead of adding code to handle this situation, the checker can detect the error and recover from that point. If the event is rare, the effects on the overall performance will be negligible. A third use is to have the checker be part of the microarchitecture itself. Dynamic on-chip verification is being investigated to reduce the burden of correctness in modern microarchitectures [2]. For these cases, the number of errors detected by the checker provides a measure of quality for the microarchitecture and its implementation in software. Fewer errors indicates a more complete microarchitectural core or a more well-behaved simulator. This can be useful when comparing aggressive microarchitectures.

2.2 Modernization of the performance model

MASE also includes several enhancements to the baseline performance model that capture several trends in microarchitecture. The register update unit (RUU) has been removed and replaced with reservation stations and a reorder buffer (ROB). Reservation station entries and reorder buffer entries are allocated when the instruction is dispatched. Dispatch is stalled if there are no reservation station or reorder buffer entry available. The reservation station entry is reclaimed when the instruction issues and the reorder buffer entry is reclaimed when the instruction commits.

The scheduler now includes a queue that instructions must wait in before entering a functional unit. The number of cycles an instruction must wait is a user-specified parameter. The scheduler will predict the latency of any variable latency instruction (such as a load), thereby speculating on when functional units will become available. If the latency prediction is too small, instructions are removed from the scheduler queue and are replayed when it is safe to do so. Options exist to do this by either flushing the entire scheduler pipeline or by only squashing dependent instructions.

Similarly, a delay has also been added to the fetch queue. Instructions cannot dispatch until they have been held in the fetch queue for a set number of cycles. This combined with the scheduler queue improves the accuracy of misprediction modeling in MASE. In SimpleScalar, there is only one delay element that was attached to the beginning of the fetch stage and all mispredictions were subject to the same delay regardless of the location of the misprediction. This approach is not accurate since mispredictions that are signalled earlier in the pipeline will have a smaller delay than mispredictions that are signalled later. For example, misfetches are signalled in the dispatch stage and are only subjected to the fetch queue delay. Branch mispredictions are signalled at writeback and are subject to both the fetch queue delay and scheduler queue delay. Front-end queue delays can be increased to simulate additional stages in the front-end pipeline.

2.3 Micro-functional performance model

The current version of SimpleScalar has no infrastructure for micro-functional simulation. Instructions are executed using an oracle in the dispatch stage - input values are read directly from architected state and results are written back right away. Values do not propagate through the pipeline and there is no notion of architectural storage, unlike true microarchitectures. The models of the various microarchitectural features in the current version of SimpleScalar concentrate exclusively on timing and performance aspects, and 'true' execution in the micro-architecture is not modeled at all.

MASE models execution as it would be in a real microarchitecture. At dispatch, new instructions get reservation stations and reorder buffer entries allocated. Register renaming takes place and input operand values are either obtained from architectural storage or from a 'completed' creator's reorder buffer entry and written into the instruction's reservation station. If the creator of an operand has not completed execution, the instruction is attached to the

dependence chain of the creator, and it gets its operand value when the creator completes execution. When all input operands of the instruction are ready, it is put into the ready queue. The issue stage attempts to issue instructions from the ready queue depending on the availability of functional resources. It is at this point that the actual execution of the instruction takes place. Input values for the instructions come from reservation station entries and the result is written to the reorder buffer entry. The instruction is entered into an event queue from which it will emerge as a completed instruction when the required amount of latency for executing the instruction has passed. During the commit stage, completed instructions are committed in-order and the results of instructions become visible at architectural storage.

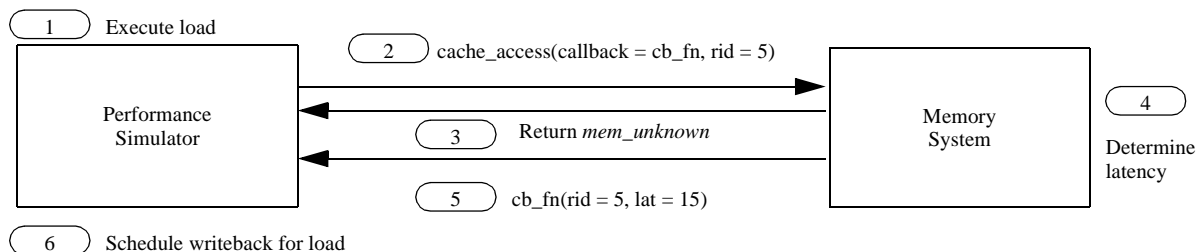
The micro-functional performance model allows for forwarding of values from completed stores to waiting loads using a load-store queue (LSQ). We have not implemented partial store forwarding, the case where a store produces only part of the value used by a later load. There is indeed a possibility that a load ends up with a stale value because we do not have complicated overlapping address checks in the LSQ. Since such situations occur rarely in practice, we decided that should such an event arise, we will address the deficiency by letting the checker capture and correct the error instead of implementing the complicated but rarely used circuitry needed for partial store forwards. We confirmed this result for the SPEC benchmarks in Section 3.3.

2.4 Recovery from any instruction

In the current version of SimpleScalar, branches are the only instructions that can cause misprediction. A branch misprediction causes a speculation mode bit to be set. When this bit is set, certain data structures are not updated, eliminating the need to provide a rollback mechanism. This implementation has two major deficiencies: (i) it prevents speculative optimizations, such as blind load speculation, where it is impossible to determine if a recovery will occur until the later stages of the pipeline and (ii) it increases the inaccuracies in the misprediction modeling since data structures are not updated during a misprediction.

MASE addresses these deficiencies by adding rollback mechanisms for several data structures throughout the machine. This includes updates to the ROB, rename table, and branch predictor. The oracle and checker also need to be synchronized so the instruction state queue, oracle memory, and oracle register file also have rollback mechanisms. All of these structures will continue to be updated even if the instructions are executed speculatively, increasing the accuracy of mispredicted instructions. When a misprediction occurs, a recovery occurs by restarting the fetch engine with the proper PC and rolling back the data structures to the appropriate point.

The rollback mechanism is illustrated by an example describing our implementation of the rename table. The rename table is used to break false register dependencies by keeping track of the latest creator (instruction that last created/wrote a value) for each logical register. In our implementation, each entry includes a linked list of all instructions that currently reside in the pipeline and create



1. Load is executed.
2. Call to memory system with unique request id (rid) 5.
3. Unknown latency, so *mem_unknown* status is returned.
4. Memory system determines latency (15 cycles).
5. Callback function is invoked with the latency value.
6. The load is scheduled to writeback.

Figure 4: Processing a load with non-deterministic latency. This diagram shows the steps of the execution for a load that requires use of the callback mechanism.

a value for the logical register. The head of the list points to the most recent creator in the pipeline. When an instruction that created a register value is committed, the corresponding entry is removed from the rename table. During a recovery event, the linked lists are scanned. Entries that refer to instructions that were squashed during the recovery are removed. Since the linked lists are ordered by position in the pipeline, once an entry is found that corresponds to an instruction that survived the recovery, the search can be terminated for that logical register.

2.5 Support for non-deterministic latencies

MASE provides a callback interface for resources that have non-deterministic latency. The callback mechanism was used to update the memory interface in SimpleScalar. Currently, any call to the memory interface returns a latency value immediately for that particular operation. This is not strictly an accurate reflection of modern day memory systems, where the actual latency for a particular access might not be immediately known. For example, DRAM systems typically have a page cache that saves the last page accessed in the memory, eliminating the row access time in subsequent requests to the same page. Assume that requests access pages in the following order: A, B, A. Since an access to page A was first, page A will reside in the page cache. If the second request to page A is seen before the access to page B is initiated, the memory system can reorder the memory requests so the second access to page A comes before the access to page B. This significantly speeds up the access time for page A since page A was already in the page cache, but it slows down the access to page B since it has to wait for the additional memory access. The memory interface in SimpleScalar requires a latency to be returned immediately after the memory request is sent. This is not sufficient in the scenario described above because it is not known if there will be another memory access to page A before page B gets to access memory, thus the true latency is not known immediately. The actual latency is only known when the DRAM memory component is accessed.

In MASE, calls to the memory interface return two values - status and latency. The status can be of three types - *mem_known*, *mem_unknown* or *mem_invalid*. If the status returned is *mem_known*, then the memory system was able to determine the latency of the operation (most likely due to a L1/L2 hit) and returns that latency immediately.

If the status returned is *mem_unknown*, the memory system is indicating that it will handle the request but the latency cannot be determined at this point in time. One of the parameters passed to the memory call is a callback function, which will be invoked by the memory system when the latency has been determined. There are different types of callback functions for different types of memory accesses. For instruction cache misses, the fetch stage is blocked indefinitely. The callback function notifies the fetch stage how much longer the operation will take so the fetch can be unblocked at the appropriate time. For store data cache misses, the latency is ignored because the store data is either in the cache or in the writeback buffers; either can be read if required. For load data cache misses, there is a table that keeps track of callback memory requests. Each entry in the table contains a unique index and a reservation station entry. When the callback function is invoked, it scans the table for a matching index. If it finds the entry and it is still valid (the entry could be invalidated if it was squashed), it schedules a writeback event for the load with the given latency (see Figure 4). Finally, if the status returned is *mem_invalid*, it means that the memory system cannot immediately handle this request for some reason (for example, a full buffer). In instruction cache accesses, the fetch stage is blocked and we keep retrying until it gets a different status. We do the same for store data cache accesses - block the commit stage and keep retrying. Loads that cannot access the data cache are not scheduled to execute. The simulator will attempt to schedule the load each cycle until it gets a status of *mem_known* or *mem_unknown*.

3. Early analyses

This section describes some analyses we have done with MASE. We compare our infrastructure to the current version of SimpleScalar to compare the performance and running time of the two simulators. In addition, we analyze the benefits of the dynamic checker and look at blind load speculation, a study that would be difficult to implement with the current release of SimpleScalar.

3.1 Simulation methodology

We used the SPEC95 integer benchmarks, compiled for the Alpha ISA using the Compaq C (version 5.9) and Fortran (version 5.3) compilers using full compiler optimization (-O4). The train input set was used for all

experiments. The benchmarks were simulated to completion or for a maximum of 250 million instructions. The baseline SimpleScalar performance model is sim-outorder from SimpleScalar/Alpha version 3.0 (the most recent release as of this writing).

The simulation configuration models a modern out-of-order processor microarchitecture. The processor model executes user-level instructions, including execution down any speculative path until the detection of a fault, TLB miss, or branch misprediction. The processor has a large window of execution; it can fetch and issue up to 4 instructions per cycle. It has a 128 entry register update unit with a 64 entry load/store buffer. To implement a comparable model in MASE, there are 128 reservation stations and 128 reorder buffer entries. Loads can only execute when all prior store addresses are known (except for the blind load speculation experiment). In addition, all stores are issued in program order with respect to prior stores. A 4k entry gshare branch predictor was used. The processor has 4 integer ALU units, 2 load/store units, 2 FP adders, 1 integer MULT/DIV unit, and 1 FP MULT/DIV unit. The latencies are: integer ALU 1 cycle, integer MULT 7 cycles, integer DIV 12 cycles, FP adder 4 cycles, FP MULT 4 cycles, and FP DIV 12 cycles. All functional units, except the dividers, are fully pipelined allowing a new instruction to start each cycle.

The processor we simulated has 64k 2-way set-associative instruction and data caches. Both caches have block sizes of 32 bytes. The data cache is write-back, write-allocate, and is non-blocking with two ports. The data cache access latency is two cycles (for a total load latency of three cycles). There is a unified second-level 512k 4-way set-associative cache with 32 byte blocks, with a 12 cycle cache hit latency. If there is a second-level cache miss it takes a total of 80 cycles to make the round trip access to main memory. There is a 16 entry 4-way associative instruction TLB and a 32 entry 4-way associative data TLB, each with a 30 cycle miss penalty.

3.2 Comparison to original model

This experiment compares the simulated performance of each benchmark to the baseline infrastructure. The purpose of this experiment was to ensure that the models implemented by the two simulators are comparable. This strategy tests the new recovery mechanism, micro-functional core, and the dynamic checker. The results of the experiment, shown in Table 1, confirm that the benchmark throughput is comparable with negligible differences.

To test the performance of the memory interface requires the use of a memory system that takes advantage of the new callback mechanism. We adopted the existing SimpleScalar memory system, randomly selecting accesses to use the callback mechanism.

Our next experiment was to see the effect of the added functionality on the running time of the MASE simulator. We compared the running time of our simulator to the baseline and the results are shown in Table 1. The runtimes (listed in seconds) indicates our simulator is twice as slow on average. This is not surprising given that our simulator is more accurate and implements additional facili-

ties such as the checker. The code has not been profiled yet to see where the bottlenecks are, so we hope to improve upon these results before the code is released. The size of the source code increased by about 50%.

Table 1: Performance validation statistics. This table compares MASE’s performance model to the baseline model. The performance is similar but MASE runs twice as slow.

Benchmark	Performance (IPC)			Run time (seconds)	
	MASE	Baseline	Diff.	MASE	Baseline
cc1	1.8581	1.8599	0.10%	6,170	2,964
compress	2.3540	2.3530	0.04%	1,208	607
go	1.6786	1.6784	0.01%	6,853	3,329
jpeg	2.8211	2.8207	0.01%	4,869	2,358
li	2.3211	2.3204	0.03%	5,784	2,856
m88ksim	2.1919	2.1969	0.22%	3,098	1,601
perl	2.2399	2.2404	0.02%	1,047	524
vortex	2.1086	2.0925	0.76%	5,870	2,895

3.3 Use of the dynamic checker

One benefit of the dynamic checker is to reduce the burden of correctness on the performance model. When implementing an optimization, most of the time is spent on corner cases that happen very rarely. The programmer can simply not implement infrequent corner cases and rely on the dynamic checker to correct any instructions that have the incorrect value. In our implementation, we decided not to implement partial store forwarding from the load-store queue when the base addresses do not match. For instance, if there is a two-byte store to address 0x102h that is followed by a four-byte load from address 0x100h, it would not be detected by the load forwarding mechanism. Instead, we rely on the checker to fix this problem. Table 2 shows the number of checker errors that were obtained. In five of the eight benchmarks, this scenario was not encountered so there were no checker errors. In the other three benchmarks, the number of errors was very small. The benchmark *vortex* was highest with 195 errors, a very small number considering that millions of instructions were executed.

Table 2: Checker errors. Number of checker errors due to not handling partial store forwarding when base addresses do not match.

Benchmark	Errors	Benchmark	Errors
cc1	82	li	0
compress	0	m88ksim	0
go	0	perl	0
jpeg	18	vortex	195

An additional benefit of the checker is that it serves as a debugging aid. This is hard to quantify in terms of numbers so we describe our debugging experience instead. The oracle and checker were among the first additions to our infrastructure so it could be used when implementing some of our other ideas. In most simulators, it is difficult to determine precisely where an error occurred when there is a difference in the output. Furthermore, some modeling errors don’t surface as errors in the output. When the

checker is used, each instruction that produces a result (writes to a register or memory) is checked. If an error is encountered, the number of errors is incremented. Optionally the error can be printed out to indicate what instruction has failed, allowing a programmer to precisely identify the point where the error occurred. The checker came in handy when implementing the micro-functional component. The first thing we realized during debugging is that most of the failing instructions referred to Alpha register \$31 (the zero register). Almost immediately, we were able to determine that the processing of this special register was incorrect. Once that problem was flushed out, we noticed that most of the problems dealt with conditional move instructions and how the output was incorrectly zero most of the time. We concentrated our debugging efforts at the conditional move and quickly identified that when the conditional move was not executed, it was not handled properly. The checker was also useful in implementing our blind speculation case study (described in the next section). As one might expect, loads were the only instruction that failed so the error message provided by the checker did not provide as much insight as in the previous cases. Instead, we focused on the first error that was signalled. We used *gdb* to debug the simulator and set a breakpoint on the failing instruction. Once we found the failing instruction, we analyzed the state of the machine at that time and were able to isolate the problem relatively quickly.

3.4 Case study: Blind speculation

This section describes a case study we performed with our new performance simulator. We implemented blind speculation, a technique that allows loads to execute before all previous stores have executed [14]. Without blind speculation, loads are required to wait in the LSQ (load-store queue) until all addresses for previous stores have been resolved. With blind speculation, the scheduler assumes that any unknown store address will not match the address of the load, as such a load can be executed as soon as its address is known and there isn't an earlier matching store address in the LSQ. The benefit of blind speculation is that it allows loads to execute earlier, but with a potential mispeculation penalty that is incurred when an earlier unknown store matches a speculated load address. In this situation, instructions after the load are flushed and the fetch engine is directed to restart at the instruction after the mispeculated load.

Blind speculation is difficult to implement in the current version of SimpleScalar because SimpleScalar requires that mispeculations be detected at fetch. Blind speculation misses can only be detected in the core since they are dependent on the contents of the LSQ when a load becomes ready.

Adding blind speculation to our infrastructure was relatively straightforward. Each cycle, the LSQ is scanned to determine if any loads can be scheduled to execute. A load can be scheduled if the address is known and there are either no store instructions that are known to write to the same address as the load or the store data must be known for the most recent store to this address. Once a load is scheduled, the LSQ is scanned, checking all unknown store instructions that are older than the load but

younger than the last known store to the load address. If the load address matches the store address (obtained from the instruction state queue), the store is marked. At write-back, the store will transfer data into the load. If the load is still valid and completed execution, a recovery will occur that squashes all instructions after the load. If the instruction hasn't completed, instructions dependent on the load haven't executed so it is safe to continue using the forwarded value.

Table 3: Blind speculation results. Results are given as speedups over the baseline. Speedups range from no speedup to 5.4%. As the LSQ size increases, the benefit flattens or lessens as in the case of *m88ksim*.

LSQ size	4	16	64	256	1,024
cc1	0.3%	1.5%	2.3%	2.5%	2.4%
compress	0.1%	0.5%	0.2%	0.7%	0.7%
go	0.2%	1.8%	2.3%	2.3%	2.3%
jpeg	0.2%	2.2%	2.5%	2.6%	2.6%
li	-0.2%	1.6%	3.2%	3.6%	3.6%
m88ksim	1.4%	5.4%	4.1%	0.7%	1.7%
perl	0.2%	0.0%	-0.1%	-0.1%	-0.1%
vortex	0.3%	0.8%	1.6%	1.8%	1.7%

After implementing blind speculation, we performed a small performance study. The baseline model is MASE with blind speculation turned off. The simulation includes different LSQ sizes and the ROB is twice the size of the LSQ for each run. The results, shown in Table 3 ranged from no impact (*compress* and *perl*) to a speedup of 5.4% for *m88ksim*. Negligible speedups are obtained with a LSQ size of 4 due to a lack of opportunity to apply speculation. The speedups increase when the LSQ size increases up until the point where the LSQ is not being utilized due to frequent mispeculation. In fact, the benefit of *m88ksim* decreased with larger LSQs since there is a higher probability of a blind speculation miss.

4. Related Work

There are a number of performance modeling infrastructures available to researchers today that implement various forms of these technologies. The Pentium Pro simulator [17], Dinero [10], and Cheetah [22] are examples of simulators that read external traces of instructions. Turandot [15], SMTSIM [24] and VMW [8], are simulators, like SimpleScalar, that generate instructions traces through the use of emulation. RSIM [16] is an example of a micro-functional simulator; instructions are emulated in the execution stage of the performance model. Unlike MASE, it does not have a trace-driven component in the front-end. This prevents oracle studies such as perfect branch prediction. The idea of dynamic verification at retirement was inspired by Breach's Multiscalar processor simulator [3]. Other simulation environments include SimOS [18] and SimICS [13] which focus on system-level instruction-set simulation. MINT [26] and ATOM [20] concentrate on fast instruction execution.

5. Summary and future work

There are two primary components in current modeling infrastructures: a functional simulator and a performance model which does timing analyses of instructions and data as they flow through the processor pipeline. While this

division has its advantages, there are drawbacks to this approach such as the inability to model timing dependent computation. The solution is to integrate functional execution with performance modeling, such that programs are executed in the performance model. This approach, however, has two disadvantages. First, inaccuracies in the performance model, even if they are infrequent, will cause programs to fail. Second, a micro-functional performance model does not lend itself to oracle-type studies.

We proposed MASE, a new performance infrastructure that addresses problems with the trace-driven approach found in SimpleScalar. Our simulation infrastructure has an oracle at the front end that is directed by the fetch stage and executes instructions using its own state. This component provides a setting for all types of oracle studies including perfect branch and value prediction. MASE contains a detailed micro-functional performance model that not only does timing analyses on instructions and data, but also executes instructions similarly to real microprocessors. The performance model, however, can afford to be inaccurate for infrequent cases (*e.g.*, partial store forwards). This is because MASE has a checker that provides online validation for our simulator. The checker will detect and correct infrequent corner cases by comparing the results of instructions from the performance model and the oracle. The checker also serves as a valuable debugging aid by identifying precisely the point at which a model failed. MASE also supports arbitrary rollback to any point in the instruction window, which enhances our ability to implement aggressive speculation techniques such as blind speculation. Lastly, our memory interface is more realistic than SimpleScalar's as it includes non-deterministic memory latency modeling for memory accesses that cannot immediately determine their latency.

We compared our new performance model to the current version of SimpleScalar. The performance of programs were comparable on both versions. The running time of our simulator was approximately twice as slow as compared to SimpleScalar due to increased accuracy and micro-functional execution in our model. As a case study, we modeled and studied blind load speculation using MASE, something which would have been very difficult to do with the baseline SimpleScalar performance model.

In the future, we want to incorporate microcode enhancements in our simulator. Instructions can be divided into a sequence of micro-ops and the functionality of the micro-ops can be defined in the simulator code. Instructions would be executed by executing the required sequence of micro-ops. We also plan to optimize the performance of the MASE code.

Acknowledgments

We would like to thank the anonymous reviewers for their valuable comments. This work was supported under a National Science Foundation Graduate Fellowship and by the NSF CADRE program, grant no. EIA-9975286. Equipment support was provided by Intel.

References

- [1] A. Agarwal, R. Sites, and M. Horowitz. ATUM: A new technique for capturing address traces using microcode. *Proc. of the 13th Annual Int. Symposium on Computer Architecture*, June 1986.
- [2] T. Austin. DIVA: A Dynamic Approach to Microprocessor Verification. *Journal of Instruction-Level Parallelism Vol. 2*, Jun. 2000.
- [3] S. Breach. Design and Evaluation of a Multiscalar Processor. *Ph.D. thesis, University of Wisconsin-Madison*, 1999.
- [4] B. Calder, G. Reinman, D. Tullsen. Selective Value Prediction. *Proc. of the 26th Annual Int. Symposium on Computer Architecture*, May 1999.
- [5] J. Bondi, A. Nanda, and S. Dutta. Integrating a misprediction recovery cache (MRC) into a superscalar pipeline. *Proc. of the 29th Annual Int. Symposium on Microarchitecture*, Dec. 1996.
- [6] D. Burger and T. Austin. The SimpleScalar Tool Set, Version 2.0. *University of Wisconsin Computer Sciences Technical Report #1342*, June 1997.
- [7] R. Desikan, D. Burger, and S. Keckler. Measuring Experimental Error in Microprocessor Simulation. *Proc. of the 28th Annual Int. Symposium on Computer Architecture*, July 2001.
- [8] T. Diep. VMW: A Visualization-based Microarchitecture Workbench. *Ph.D. thesis, Carnegie Mellon University*, June 1995.
- [9] J. Dundas and T. Mudge. Improving Data Cache Performance by Pre-executing Instructions Under a Cache Miss. *Proc. of the 1997 Int. Conference on Supercomputing*, July 1997.
- [10] J. Edler and M. Hill. Dinero IV Trace-Driven Uniprocessor Cache Simulator. <http://www.neci.nj.nec.com/homepages/edler/d4>.
- [11] A. Lai and B. Falsafi. Memory sharing predictor: The key to a speculative coherent DSM. *Proc. of the 26th Annual Int. Symposium on Computer Architecture*, May 1999.
- [12] M. Lipasti, C. Wilkerson, J. Shen. Value locality and load value prediction. *Proc. of the 7th Int. Conf. on Architectural Support for Programming Languages and Operating Systems*, Oct. 1996.
- [13] P. Magnusson, F. Dahlgren, H. Grahm, M. Karlsson, F. Larsson, F. Lundholm, A. Moestedt, J. Nilsson, P. Stenström, and B. Werner. SimICS/sun4m: A Virtual Workstation. *Usenix Annual Technical Conference*, June 1998.
- [14] A. Moshovos and G. Sohi. Memory Dependence Speculation Tradeoffs in Centralized, Continuous-Window Superscalar Processors. *The 6th Annual Int. Symposium on High Performance Computer Architecture*, Jan 2000.
- [15] M. Moudgill, J. Wellman, J. Moreno. Environment for PowerPC Microarchitecture Exploration. *IEEE Micro*, May/June 1999.
- [16] V. Pai, P. Ranganathan, and S. Adve. RSIM Reference Manual. Version 1.0. *Technical Report 9705, Department of Electrical and Computer Engineering, Rice University*, July 1997.
- [17] D. Papworth. Tuning the Pentium Pro Microarchitecture. *IEEE Micro*, April 1996.
- [18] M. Rosenblum, S. Herrod, E. Witchel, and A. Gupta. Complete computer system simulation: the SIMOS approach. *IEEE Parallel & Distributed Technology: Systems & Applications*, Winter 1995.
- [19] A. Sodani and G. Sohi. Dynamic Instruction Reuse. *Proc. of the 24th Int. Symposium on Computer Architecture*, June 1997.
- [20] A. Srivastava and A. Eustace. ATOM: A system for building customized program analysis tools. *Proc. of the 1994 Symposium on Programming Language Design and Implementation*, June 1994.
- [21] K. Sundaramoorthy, Z. Purser, and E. Rotenberg. Slipstream Processors: Improving both Performance and Fault Tolerance. *Proc. of the 9th Int. Conference on Architectural Support for Programming Languages and Operating Systems*, Nov. 2000.
- [22] R. Sugumar and S. Abraham. cheetah - Single-pass simulator for direct-mapped, set-associative and fully associative caches. *Unix Manual Page*, 1993.
- [23] TraceBase. *Parallel Architecture Research Laboratory, New Mexico State University*. <http://tracebase.nmsu.edu/tracebase.html>.
- [24] D. Tullsen, S. Eggers, and H. Levy. Simultaneous Multithreading: Maximizing On-Chip Parallelism. *Proc. of the 22nd Annual Int. Symposium on Computer Architecture*, June 1995.
- [25] R. Uhlig, R. Fishtein, O. Gershon, I. Hirsh, and H. Wang. SoftSDV: A Pre-silicon Software Development for the IA-64 Architecture. *Intel Technology Journal*, 4th quarter 1999.
- [26] J. Veenstra and R. Fowler. MINT: a front end for efficient simulation of shared-memory multiprocessors. *Proc. of the 2nd Int. Workshop on Modeling, Analysis and Simulation of Computer and Telecommunications Systems*, Jan. 1994.