PIPELINE SPECTROSCOPY

T. R. Puzak, A. Hartstein, V. Srinivasan, and P. G. Emma IBM - T. J. Watson Research Center Yorktown Heights, NY 10598 trpuzak, amh, viji, pemma @us.ibm.com

ABSTRACT

Pipeline Spectroscopy is a new technique that allows us to measure the cost of each cache miss. The cost of a miss is displayed (graphed) as a histogram, which represents a precise readout showing a detailed visualization of the cost of each cache miss throughout all levels of the memory hierarchy. We call the graphs because they reveal certain signature 'spectrograms' characteristics of the processor's memory hierarchy, the pipeline, and the miss pattern itself. Cache miss spectrograms are produced by analyzing misses according to the miss cluster size, and comparing instruction sequences and execution times that occurred near the miss cluster in a 'finite cache' simulation run to the same set of instructions and execution times in an 'infinite cache' run, then calculating the difference in run times. We show that in a memory hierarchy with N cache levels (L1, L2, ..., L_N , and memory) and a miss cluster of size C, there are $\begin{pmatrix} C + N \\ C \end{pmatrix}$ possible miss penalties. This represent all possible sums from all possible combinations of the miss latencies from each level of the memory hierarchy (L2, L3, ... Memory) for a given cluster size. Additionally, a theory is presented that describes the shape of a spectrogram, and we use this theory to predict the shape of spectrograms for larger miss clusters. Detailed analysis of a spectrograph leads to much greater insight in pipeline dynamics, including effects due to prefetching, and miss queueing delays.

Categories and Subject Descriptors

C.4 [Performance of Systems]: Design studies, Measurement techniques, Modeling techniques, Performance attributes

General Terms

Algorithms, Measurement, Performance, Theory

Keywords

1

Cost of a miss, Cache, Probability Transition Matrix, Convex Combination

1. INTRODUCTION

In order to improve the performance of a processor or an application, designers have increased the amount of parallelism between the levels of the memory hierarchy. This area of research, termed memory-level-parallelism (MLP) has been explicitly studied in [1, 2, 3] while early studies focused on modeling and evaluating performance with ILP processors [4, 5, 6]. Chou, et al. [7, 8] studied several techniques (out-of-order, runahead, value prediction, prefetching, store handling

optimization) for increasing MLP in applications that are dominated by memory delays. They show that substantial amounts of performance gains are possible by increasing the MLP in these applications. Qureshi et al. [9] demonstrates that not all misses have the same cost and measures miss parallelism to improve cache performance by altering the replacement algorithm. In this paper, we build on this work by describing a new technique that quantitatively measures the cost of each cache miss though out all levels of the memory hierarchy. We call this new technique 'pipeline spectroscopy' and the graphs representing the miss cost a 'spectrogram'. The graphs are called spectrograms because they reveal certain signature features of the processor's memory hierarchy, the pipeline, and the miss pattern itself (e.g. amount of overlap between misses in the miss cluster). Using pipeline spectroscopy, we are able to quantify the cost of each cache miss and measure the rate that misses are satisfied from the different levels of the memory hierarchy (L2, L3, ... Memory). This quantification leads to a much greater understanding of the amount of parallelism or overlap (miss cost) that the micro architecture and application allow. Armed with this information designers can individually analyze each miss and improve the performance of the hardware or software.

Several mechanisms that measure pipeline stalls and miss costs are described in the patent literature [10-16]. Each technique describes the difficulty in determining an accurate measure for the cost of the miss and rely on hardware monitors to count events (cycles) that indicate when the decoder or execution unit is delayed (stalled) while waiting for an operand (data) to estimate this cost. However, not all of these events contribute to the loss of performance in a program. Today's processors have superscalar capabilities and parallel execution paths and a delay suffered in one component of a processor can be overlapped with other events to mask any loss due to the miss. For example, consider two events occurring in parallel: a branch miss-prediction and a cache miss. Simply counting the number of cycles an instruction (in the decoder or execution unit) is stalled waiting on a miss is not an accurate measure of the cost of the miss since many of the stall cycles are already overlapped with the delays caused by the branch miss-prediction.

Additional performance tools are described in: Dean et al. [17] a technique for pairwise sampling used to track concurrent events to measure performance, Fields et al. [18] use 'shotgun profiling' to construct dependence graphs and study the performance of an application, while Karkhanis et al. [19] use analytical models to study the performance of concurrent events.

In our work, we apply pipeline spectroscopy to produce a cache miss spectrogram which represents a precise readout showing a detailed histogram (visualization) of the cost of each cache miss through out all levels of the memory hierarchy, with and without overlap. We find this visualization very helped in analyzing cache and pipeline performance (prefetching algorithms and bus queueing). Cache miss spectrograms are produced by comparing instruction sequences and execution times that occurred near a miss in a 'finite cache' simulation run to the same set of instructions and their execution times in an 'infinite cache' run.

Next, we present a theory that describes the underlying shapes visualized in a miss spectrogram. The theory is based on observations and probability distributions drawn from the miss clusters of size 1 and 2 and is able to predict the miss patterns that will occur in larger miss clusters. By understanding the forces that contribute to the size and shape of the spectrogram, we hope to gain insight into pipeline dynamics, and techniques that can be used to improve the hardware and software.

The rest of this paper is organized as follows: Section 2 contains definitions and terminology. Section 3 describes constructing a miss spectrogram. The simulation model is described in Section 4. In Section 5 we measure the cost of a data miss. In Section 6, we make several observations regarding the properties of a spectrogram, which leads to developing a theory that describes the miss patterns, described in Section 7. In Section 8, we show simulation results for the theory. Section 9 describes an instruction miss spectrogram, and Section 10 shows how bus queueing changes the shape of a spectrogram. Summary and conclusions are discussed in Section 11.

2. PERFORMANCE TERMINOLOGY

The overall methodology used to calculate the cost of a miss and the visualization process are explained as a prelude to analyzing a miss spectrogram. First, the definitions and formulas used to calculate the cost of a miss are described, then a description is set forth relative to how misses cluster and affect the standard operation of a high performance processor, followed by a description of the visualization process.

The most commonly used metric for processor performance is, "Cycles Per Instruction" (*CPI*). The overall *CPI* for a processor has two components: an "infinite cache" (*CPI_{INF}*) component and a "finite cache adder" (*CPI_{FCA}*).

$$CPI_{OVERALL} = CPI_{INF} + CPI_{FCA}$$
(1)

 CPI_{INF} represents the performance of the processor in the absence of misses (cache, and TLB). It is the limiting case in which the processor has a first-level cache that is infinitely large and is a measure of the performance of the processor's organization with the memory hierarchy removed. CPI_{FCA} accounts for the delay due to cache misses and is used to measure the effectiveness of the memory hierarchy.

Just as processor performance (for both in-order and out-of-order machines) can be expressed in terms of a *CPI*, the "memory adder" can be expressed as the product of an event rate (specifically, the miss rate), and the average delay per event (cycles lost per miss):

$$CPI_{FCA} = \left(\frac{Misses}{Instruction}\right)\left(\frac{Cycles}{Miss}\right)$$
(2)

Substituting for CPI_{FCA} in (1), the overall performance for a processor can be expressed as:

$$CPI_{OVERAIL} = CPI_{INF} + \left(\frac{Misses}{Instruction}\right) \left(\frac{Cycles}{Miss}\right)$$
(3)
Solving for the average cost of a cache miss, we have:
$$\left(\frac{Cycles}{Miss}\right) = (CPI_{OVERALL} - CPI_{INF}) \left(\frac{Instructions}{Miss}\right)$$
(4)

We use this formula to calculate the amount of time (cycles) a processor loses due to each cache miss. The following example illustrates calculating cycles per miss using equation (4). Consider an application whose entire run length is one million instructions and a processor where each cache miss is satisfied from the L2 that is 20 cycles away. If an infinite cache simulation run takes one million cycles ($CPI_{INF} = 1$), and a finite cache simulation run takes 1.3 million cycles, then cache misses account for 300,000 cycles and the total CPI = 1.3 and $CPI_{FCA}=.3$. If the finite cache simulation run generates 25,000 misses, then $(\frac{Misses}{Instruction}) = \frac{25,000}{1,000,000} = \frac{1}{40}$ and $(\frac{Cycles}{Miss}) = \frac{300,000}{25,000} = 12$. By applying this equation over the entire length of an application,

by applying this equation over the entire length of an application the average cost for all misses can be calculated.

In the example above, we applied Eq. 4 macroscopically to calculate the average cost of a miss over the total run time of an application. However, Eq. 4 can also be used microscopically to calculate the cost of a single miss. We will take a microscopic approach in using Eq 4 to calculate the cost of each miss and produce a miss spectrogram. As presented in Section 5, the information contained in a miss spectrogram represent the cost of all misses throughout all levels of the memory hierarchy, including the amount of overlap (parallelism) achieved between any two misses.

3. SPECTROGRAM CONSTRUCTION

A description of how misses can cluster and affect the performance of a processor is now described. Figure 1 shows the same five instructions executed as an 'infinite cache' sequence of instructions and a 'finite cache' sequence of instructions. In the finite cache sequence, the instruction decode times are shown in bold (without parenthesis) and instruction completion or EndOp times are shown in parenthesis. In the infinite cache run only the instruction decode times (in bold) are shown.



Figure 1. Miss Cluster Patterns for an Application, Miss Cluster of Size 1 and 3 are Shown

Associated with the finite cache run are two miss clusters, where a miss cluster is the maximal continuous interval of time characterized by at least one miss in progress at all times. The size of the miss cluster is the number of misses that started during this interval. In the finite cache run, the first miss cluster has three misses with overlap (size = 3) and the second miss cluster is size = 1 (a miss in isolation).

Next we describe a technique used to calculate the cost of the miss (cluster). The time to process the first miss cluster (in the finite cache run) is bounded by the decode time for instruction I1 and the EndOp time of I3, (I3EndOp - I1Decode)_{Finite Cache} time. Instruction I1 (the decode time) represents the greatest lower bound of the miss cluster, while instruction I3 (EndOp time) is the least upper bound of the cluster. We call these points the infimum and supremum of the miss cluster. (By convention, the infimum of a miss cluster is the greatest instruction (time) that decoded just prior to the beginning of the first miss in the miss cluster and the supremum is the first instruction (time) that completed (EndOp) just after the last miss in the miss cluster finished.)

Similarity, the infimum instruction of the second miss cluster is I4 and the supremum is I5. The time to process the second miss cluster is then (I5EndOp - I4Decode)_{Finite Cache}. To calculate the amount of delay associated with the first miss cluster we must subtract the amount of time to process the same set of instructions in an infinite cache run from the finite cache run. That is, [(I3EndOp - I1Decode)_{Finite Cache} - (I3EndOp - I1Decode)_{Infinite Cache}] equals the number of cycles the pipeline was stalled due the first miss cluster. Similarly, the amount of delay associated with the second miss cluster is [(I5EndOp - I4Decode)_{Finite Cache}].

The reader will note that in the derivation above and in the equations presented in Section 2, no mention was made as to whether the processor is in-order or out-of-order. That is because out-of-order processing will not change the analysis. However, it may affect the manner in which the infinite cache running times for the sequence of instructions that surround a miss need to be determined. For example, if instruction processing is from an out-of-order processor, it may be necessary to save the sequence of instructions between the infimum and supremum of the miss (from the finite cache run) and use this same sequence of instructions (and their order) while determining the infinite cache run time.

By applying the above technique repeatedly, we can calculate the cost of a miss or miss cluster for both in-order or out-of-order processors, one cluster at a time. In the example above, I1, I2, and I3 can even be from three different threads running on a multithreaded processor (or three out-of-order instructions), but as long as the same three instructions (and their order) are used to determine the infinite cache run time, the cost of the miss cluster can be determined.

3

Confidential - Do Not Distribute

There are certain boundary conditions that must be considered when determining the infimum and supremum of a miss cluster. For example, the infimum of a miss cluster can only be established after the supremum of the previous miss cluster has been determined. This ensures that one miss cluster is terminated before another starts. If the upper and lower bounds of a miss cluster cannot be uniquely established, the two adjoining miss clusters are combined into a large miss cluster.

Also, when determining the infinite cache running time for an instruction sequence that occurred during a miss cluster, it may be necessary to prime the processor's pipeline with some of the instructions that occurred prior to the infimum instruction. This ensures that the correct execution and EndOp times of the infimum instruction are preserved as it passes through the processor's pipeline. By grouping misses according to their cluster size and calculating the delay associated with a miss cluster (number of stall cycles) using the method described above, the amount of time a processor loses due to cache misses is produced.

4. SIMULATION METHODOLOGY

To date, pipeline spectroscopy has been implemented in three proprietary processor simulators. Each timer has produced results similar to those shown in Sections 5 and 9 below. Each timer is cycle accurate and has been thoroughly validated against existing hardware. The processor model used in this paper is shown in Figure 2 and described in [20, 21], is a 4 issue superscalar, with address generation and cache access an independent out-of-order process. We use trace tapes produced for the IBM zSeries processor family. Instructions that typically produce addresses (LA, BXLE, BALR, BAL, ...) are pre-executed after the decode stage of the pipeline to avoid future pipeline stalls due to address interlocks. Loads are executed as soon as the datum fetched returns from the cache and the results are forwarded to all dependent instructions. The instruction window was set at 32 entries. Separate L1 instruction and data caches were modeled at 64 KB, the L2 size varies from 256K to 1 MB, and the L3 (when modeled) was varied from 1 MB to 4 MB. This processor model was chosen to illustrate the technique used to construct a spectrogram, and does not represent any existing or planned processor design. In our initial studies, Endop and those instructions not pre-executed after the decode stage are completed and executed in-order. Future work is planned to measure the benefits of prefetching, SMT and SMT out-of-order execution.



Fig. 2 Pipeline modeled for study. Stages include: Decode, Rename, Agen Q, Agen, Cache Access, Execute Q, Execute, Completion and Retire.

In order to stress different levels of the memory hierarchy (L1, L2 or L3), we use applications with large instruction and data footprints capable of stressing caches up to 4 Megabytes. Typically, commercial database applications have these characteristics [7]. In our study, we use eight workloads drawn from database workloads, SPEC 2000, and a C++ application. We use three proprietary commercial database applications running on zSeries servers, (oltp, and oltp2, and oltp3 described in [20, 21]); mcf, gcc, and perl from SPEC 2000, SPECjbb 2000, and perf1 a large-processor simulator written in C++. Typically, trace tape lengths are 5 to 100 million instructions. The simulation environment can handle all of the SPEC suite; the application subset used for this work was chosen for its ability to stress L2 and even L3 cache usage.

5. SPECTROGRAM FOR DATA MISSES

In order to examine the miss spectrogram for data misses alone, we model an infinite or perfect instruction L1 cache, and set the data L1 cache to 64KB. The L2 is set to 256KB with a 15 cycle latency, and set L3 latency to 100 cycles. All L2 misses are resolved in the L3. The line size and bus width are set at 128 bytes. No data prefetching was modeled. In fact, data prefetching is very difficult for many of these applications. However, prefetching will be explored when we examine the instruction miss spectrogram. Using the techniques described above, Figure 3 shows the miss spectrogram for the oltp workload. The overall hit ratio of the L2 was approximate 50%.

The miss spectrograms for cluster sizes = 1, 2, 3, and 4 are shown. The X axis represents the cost of the miss for the specified cluster size. The Y axis shows the percent of misses that had that delay.

The cluster =1 plot (in Figure 3) shows two peaks. The first peak is centered near 15 cycles (the L2 hit latency), and the second peak is near 100 cycles (the L3 hit latency). The area under each peak is approximately the percentage of L1 misses resolved in the corresponding level of the memory hierarchy (i.e., the hit rates for the L2 and L3, or 50% in each).



The cluster size = 2 plot shows peaks at 15 and 30, 100 and 115, and 200 cycles. Notice that each peak represents the cost of a miss cluster (two L1 misses) and identifies one of all possible hit/miss combinations and all possible overlap/no-overlap patterns that the L1 misses had in the L2 or L3. Additionally, each peak identifies the degree of overlap/no-overlap (parallelism) between the two misses. Calculating the area under each peak, we see that the five miss costs (15, 30, 100, 115, and 200 cycles) have a probability of .138, .168, .288, .191, and .215, respectively. The peaks at 15 and 30 represent two L1 misses that both hit in the L2 but highlight two distinctively different outcomes. In the first case (peak at 15), both misses had a high degree of overlap (parallelism) and the overall cost was approximately the L2 hit latency while in the second case there was little overlap and the cost of the miss cluster was the sum of two L2 hits.

The peak at 100, again identifies two misses that were overlapped (had a high degree of MLP). Whether it was two misses that hit in the L3, one miss that hit in the L2 and one that hit the L3, the overall cost of the misses in the cluster was just the L3 hit latency.

The peak at 115, identifies two misses that had little or no overlap. Here, one miss hit in the L2 and one miss hit in the L3 but the cost of the miss cluster was the sum of the individual miss latencies. Finally, the peak at 200 identifies two misses that were resolved in the L3 and there was little overlap.

The peaks in the cluster = 3 graph again represent all possible hit/miss combinations (with and without overlap) of length 3 using the two miss latencies (15, 100) for the L2, and L3. For example, the peaks at 15, 30, and 45 present three L2 hits where two misses were overlapped, one miss was overlapped or no miss was overlapped with the other misses in the cluster. However, the peak at 300 represents three L3 hits with little overlap. Obviously, three dependent misses that are resolved in the L3 can cause this miss penalty. Finally, the peaks in the cluster = 4 graph show all of the hit/miss, overlap/no-overlap, combinations of length 4 using the miss latencies 15 and 100.

Each peak represents the amount of time the group of cache misses (cluster) stalled the pipeline. By summing the 'stall cycles' calculated for each miss cluster, we can reconstruct the finite-cache-adder for the entire run, one cluster at a time. In many cases this involves summing the delay associated with 10s of thousands to over 100,000 miss clusters. Using this technique, we have always been able to calculate the total finite-cache-adder to within 5% (one cluster at a time), and in many cases the error is less than 2%. This shows how accurately we can identify miss clusters and evaluate their costs.

Prefetching and bus delays can change the shape of the peaks in a spectrogram. Prefetching can broaden the left shoulder of any peak and show the degree that a prefetch is being issued in advance of the nominal miss penalty. Queueing and bus delays can increase the right shoulder of a peak, adding miss latency. The effects of prefetching, bus queueing, and changes in the space of a miss spectrogram for both instruction and data misses are analyzed and discussed Sections 9 and 10.

Figures 4 and 5 plot the number of misses and cycles per miss versus cluster size, respectively. Even though the maximum miss cluster for the run was well over 1000 misses, typically the average miss cluster size is much smaller. For example, over 80% of the misses occur to miss clusters of size 6 or less and over 40% of the misses are a miss in isolation. This was observed for most



applications used in this study.

In Figure 5, notice how the average miss penalty decreases as the cluster size grows. The slope of the line indicates the degree that miss parallelism or miss overlap is occurring. Obviously, the greater the amount of miss overlap the greater the absolute value of the slope of the line. In this example the cost of a miss at a cluster size = 10 is approximately two thirds the cost of an isolated miss. This is far less than the potential for complete overlap.

Next, we repeat the above experiment but use the oltp2 workload and add an L3. We use the following memory hierarchy: data L1=64KB, L2=256KB 15 cycle latency, L3=1MB 75 cycle latency, and 300 cycle memory latency. Figure 6 shows the data miss spectrogram for cluster sizes = 1, 2, and 3.

Again, each peak in the miss spectrogram identifies a signature feature of the memory hierarchy: either a hit/miss combination of the miss cluster throughout all levels of the memory hierarchy or a



overlap/no-overlap condition among the misses (amount of parallelism between the misses in the cluster).

The cluster = 1 plot shows three peaks. The first peak is centered near 15 cycles (the L2 hit latency), the second peak is near 75 cycles (the L3 hit latency) and the third peak at 300 cycles (the memory latency). Calculating the area under each peak is approximately the hit ratio (regarding L1 misses) for that level of the memory hierarchy (i.e., the hit ratios for the L2, L3, and memory). Examining the plots for cluster sizes equal 2 and 3, we see that they show all of the hit/miss, overlap/no-overlap combinations for a miss cluster of size 2 and 3 using the 3 miss latencies: 15, 75 and 300. Obviously the peak at 15 for the cluster size = 3 indicates a great deal of overlap among the three misses. However, the peak at 390 (for cluster size =3) indicates very little overlap among the three misses. Here, the three miss cluster had one miss hit in the L2, one hit in the L3, and one went all the way to the memory but the total miss penalty for the whole cluster is the sum of the individual miss latencies (15, 75, and 300). Similar hit/miss patterns, and overlap/no-overlap conclusions can be drawn for examining any peak in the miss spectrogram.

We will refer to spectrograms like the one shown in Figures 3 and 6 as the 'canonical' representation for the cost of a miss in a multilevel memory hierarchy. It is a canonical form because it represents the most general form (combinations) of the miss patterns in a memory hierarchy. Obviously, prefetching and bus queueing can alter the miss patterns, and costs. By including the possibility of a peak at zero, a miss spectrogram can have all possible sums from all possible combinations of the miss latencies from each level of the memory hierarchy for a given cluster size. We show in Appendix A that for a memory hierarchy with N cache levels (L1, L2, L3,..., L_N, memory) and a miss cluster of size $\begin{pmatrix} C+N\\ C \end{pmatrix}$ C, there are (5)

possible penalties (peaks) that characterize the canonical form hit/miss and overlap patterns. A peak at zero has the physical meaning that a miss or cluster of misses has zero delay. Prefetches, if issued far enough in advance of their use, speculative misses that do not interfere with any other cache accesses, or unused prefetches have the possibility of causing zero delay. Using (5), and the memory hierarchy described in Figure 6, we see that plotting miss clusters of size 4, 5, and 6 could have 35, 56, and 84 peaks, respectively.

6. OBSERVATIONS

Before we develop a theory that describes the miss patterns of a spectrogram it is necessary to understand the ways hits and misses form clusters. We start by defining a notation, then examine the individual probabilities for each hit/miss combination.

Let random variable X_i denote the event that the i_{th} L1 miss in a miss cluster of size N is a hit or miss in the L2. Then, the sequence HM in a miss cluster of size 2 (a hit followed by a miss) can be expressed as the pair X_1X_2 where $X_1 = H$ and $X_2 = M$ and has the conditional probability $\Pr[X_2 = M \mid X_1 = H]$.

Let *T* be the total number of L1 misses and let T = M + H where *M* and *H* are total number of hits and misses in the L2. Let *p* represent the probability of a miss in the L2, then $\Pr[M] = p = \frac{M}{H+M} = \frac{M}{T}$ and the probability of a hit is $\Pr[H] = 1 - p$. Let the maximum miss cluster size be *N* and let M_i , H_i , and T_i equal the total number of misses, hits, and total accesses to clusters size *i*. Then the probability of a miss in cluster *i* is $p_i = \frac{M_i}{M_i+H_i} = \frac{M_i}{T_i}$.

The L2 miss ratio p is related to the miss ratios of the individual clusters in the following manner, $p = \frac{M}{T} =$

$$\frac{\sum M_i}{T} = \frac{M_1}{T_1} \left(\frac{T_1}{T}\right) + \frac{M_2}{T_2} \left(\frac{T_2}{T}\right) + \dots \frac{M_N}{T_N} \left(\frac{T_N}{T}\right)$$
$$= p_1\left(\frac{T_1}{T}\right) + p_2\left(\frac{T_2}{T}\right) + p_3\left(\frac{T_3}{T}\right) + \dots p_N\left(\frac{T_N}{T}\right)$$

That is, the miss ratio is the weighted sum of the miss ratios for each miss cluster size, where the weights are the probabilities of a hit or miss to a given cluster size. We refer to P as the global miss ratio for the L2 and p_i as the local miss ratio relative to all of the misses that make up cluster size *i*. Typically, each local miss ratio is within +3 percent of the global miss ratio.

We now examine the individual probabilities for each hit/miss combination that occurs within a miss cluster. We use the spectrogram shown in Figure 3. Recall, in this experiment we set the L1 and L2 cache sizes to 64KB and 256KB to produce a L2 hit/miss ratio as close to 50% as possible. The actual hit/miss ratios were .493 and .507, respectively.

The cluster = 1 spectrogram comes from a miss in isolation that either hit or missed in the L2. Each of these events has its own (unique) spectrogram. We show these in Figure 7. We see that the hits have values around 15 and the misses have values around 100 cycles. The individual hit/miss probabilities are $P_1 = \Pr[X_1 = M] = .519$ or $\Pr[X_1 = H] = .481$.



There are four different hit/miss combinations that make up the cluster = 2 spectrogram shown in Figure 3. The random variable pair X_1X_2 can either be a HH, MM, HM, or MH. Each of these hit/miss combinations have their own spectrograms. These are shown in Figure 8 along with their overlap/no-overlap probabilities. The overlap, no-overlap probabilities were obtained by calculating (integrating) the area under each of the two peaks in the graph. Each spectrogram describes a different probability distribution for the cost of the miss cluster. For example, the only spectrogram that has penalties near 15 and 30 cycles is the HH spectrogram. The peak at 15 identifies two L1 misses (L2 hits) that were almost entirely overlapped while the peak at 30 indicates two dependent misses with nearly no overlap. Calculating the area under each peak determined the probability of overlap and no overlap. The HH spectrogram had a miss overlap probability of .45 and .55 no overlap probability.



Figure 8. Cluster Size=2 Individual Miss Spectrograms, L1=64KB, L2=256KB, 15 Cycle Latency, L3 = 100 Cycles Latency, OLTP

The HM and MH spectrograms each have miss probabilities centered near 100 and 115. Again, the peak near 100 indicates two misses with overlap (one was a L2 hit and one was a L2 miss), while the peak at 115 indicates two misses with little overlap (the miss penalty is the sum of a L2 hit and L2 miss). The MM spectrogram has two peaks: one near 100, one near 200. The miss penalty near 100 indicates both misses were overlapped, while the peak near 200 (only in the MM spectrogram) indicates very little overlap. Combining each of four hit/miss spectrogram produces the cluster = 2 spectrogram shown in Figure 3.

The local miss probabilities was $P_2 = .512$, and the four hit/miss combinations had the following probabilities:

$$MM \sim \Pr[X_1, X_2 = MM] = .331, HM \sim \Pr[X_1, X_2 = HM] = .188$$

 $MH \sim \Pr[X_1, X_2 = MH] = .175, HH \sim \Pr[X_1, X_2 = HH] = .306$

Examining the probabilities more closely, we see that the probability of a hit following a hit (HH) and a miss following a miss (MM) are much higher than independent events would predict. If the probability of a hit or miss were independent, each of the four hit/miss combinations would have a probability close

to .25. For example, if they were independent we could use the binomial distribution to describe the probability of having k L2 misses out of N L1 misses as $\binom{N}{k}p^{k}(1-p)^{N-k}$ and then $\Pr[X_1 = M \cap X_2 = M] = \Pr[X_1 = M] \Pr[X_2 = M] = .5 \times .5 = .25.$ However, this is not the case. We observe that the MM and HH events have a much higher probability than the MH and HM probabilities. Thus, in the random variable pair X_1X_2 , X_2 is strongly correlated with X_1 .

This observation fits nicely with the spatial locality properties of a cache and Denning's [22] working set model for program behavior. If X_1 is an L2 hit, then X_2 is likely to be a hit (indicating the L2 contains the working set of the application). However, if X_1 is an L2 miss, then X_2 is likely to miss (indicating the application is changing working sets).

It is possible to construct the cluster = 2 spectrogram (from Figure 3) from the four hit/miss probabilities and their respective overlap, no-overlap probabilities. Figure 9 illustrates this process as a tree. The figure shows events X_1X_2 with its four hit/miss combinations and overlap/no-overlap probabilities. The values are produced by multiplying the corresponding hit/miss probability with the appropriate overlap/no-overlap probability at the appropriate level. Note, the four hit/miss overlap/no-overlap probabilities were draw from the cluster size = 2 hit/misscombinations shown in Figure 8, and use the four, cluster size =2hit/miss probabilities shown above. The top of the tree shows X_1 . The right fork represents a miss and the left fork a hit. Level two shows X_2 and its hit/miss possibilities. The four hit/miss combinations along with their probabilities are shown on level The overlap, no-overlap probabilities are shown as three. branches off of each hit/miss possibility. The cost of each miss cluster is shown as a leaf node of the tree. For example, there are three events that have a possible 100 cycle miss penalty: the HM with overlap, the MH with overlap, and the MM with overlap. Summing all events with the same cost produces the probability





Confidential - Do Not Distribute

distribution representing the cluster = 2 spectrogram shown in Figure 3.

Returning to Figure 3, the cluster = 3 has eight possible hit/miss patterns (MMM, MMH, MHM, MHH, HMM, HMH, HHM, HHH). Thus a miss cluster of size N has 2^N hit/miss combinations. The local miss probability was $P_3 = .502$ and the eight hit/miss combinations had probabilities:

 $MMM \sim \Pr[X_1, X_2 X_3 = MMM] = .219, MMH \sim \Pr[X_1, X_2 X_3 = MMH] = .083$ $MHM \sim \Pr[X_1, X_2 X_3 = MHM] = .094, MHH \sim \Pr[X_1, X_2 X_3 = MHH] = .103$ $HMM \sim \Pr[X_1, X_2 X_3 = HMM] = .098, HMH \sim \Pr[X_1, X_2 X_3 = HMH] = .090$ $HHM \sim Pr[X_1, X_2 X_3 = HHM] = .105, HHH \sim Pr[X_1, X_2 X_3 = HHH] = .208$

Again we see that the events $X_1X_2X_3$ (in a cluster of size 3) are not independent. The probabilities for the HHH and MMM events have over twice the probability of the other hit/miss combinations. Each of the eight hit/miss combinations has their own individual spectrogram. Because of space limitations, we only present the spectrograms for the MMM, HHH, MHM, and HMH combinations along with their overlap/no-overlap probabilities in Figure 10. Each graph has three or four peaks, indicating the overlap/no-overlap probabilities of the hit/miss pattern. The HHH graph has peaks near 15, 30 and 45 indicating the degree of overlap between the 3 L2 hits (12, 48, and 40 percent). The sum of all eight individual spectrograms produces the nine peak (cluster = 3) spectrogram shown in Figure 3.



Summarizing the above experiments, we observe that in a cluster of N misses, there are 2^N possible hit/miss combinations. Also, hit/miss patterns are not independent. In the miss sequence $X_1X_2X_3,...,X_N$ the probability that X_{i+1} has the same outcome as X_i is much higher than independent miss probabilities. Finally, a miss spectrogram can be reconstructed by knowing the individual probabilities of each of the 2^N hit/miss combinations and their associated overlap/no-overlap probabilities. In the next

section we create a model for predicting the miss costs and associated probabilities found in a spectrogram.

7. THEORY

In this section we develop a theory that allows us to predict the size and shape of a spectrogram. The theory draws on six parameters: the overall L2 miss ratio, a correlation parameter describing the probability of a hit following a hit or miss following a miss in the L2, and the four overlap/no-overlap probabilities calculated from the four hit/miss combinations found in a cluster size =2 miss pattern. Knowing only these parameters (nostly drawn from cluster size = 2 misses) we can predict the patterns found in spectrograms for miss clusters of size 3, 4, and 5.

In order to predict the miss costs and associated probabilities found in a spectrogram, we describe a two step process. The first step involves predicting the probabilities for the individual 2^N hit/miss combinations that make up a cluster of size N. Recall, it was shown that these probabilities are dependent. The second step involves determining the overlap/no-overlap probabilities for adjacent misses in a cluster size = 2 miss pattern. Recall, there are four of them. A miss spectrogram in produced from the product of these two probabilities, then combining alike terms (cost).

We start by developing a model that increases (boosts) the probability of a miss following a miss (or hit following a hit). Let $\Pr[M] = p$ and $\Pr[H] = (1-p)$, we seek a function that increases the probability of a MM being greater than the product of $p \times p$. Consider the random pair $X_i X_{i+1}$, we need to define transition probabilities that describe the four hit/miss combinations in a cluster of size 2. To increase $\Pr[X_2 = M \mid X_1 = M]$, we define *a* (alpha), $0 \le a \le 1$ and form a convex combination between the values of *p* and 1 such that $\Pr[X_2 = M \mid X_1 = M] = ap + (1-a) = (1-a+ap)$.

Note that the function (1 - a + ap) has the properties we desire. For $0 \le a \le 1$ and $p \le 1$, $(1 - a + ap) \ge p$. In Appendix B, we show $a = 1 - \rho$ where ρ (rho) is the correlation coefficient of $X_i X_{i+1}$. As $\rho \to 1$, (the correlation between $X_i X_{i+1}$ becomes stronger) $a \to 0$, and the $\Pr[X_2 = M | X_1 = M] \rightarrow 1$ signifying that the probability of a miss following a miss is 1 (perfect correlation). Similarly, as there becomes less correlation between $X_i X_{i+1}$ then $\rho \to 0$ and $a \to 1$ then $\Pr[X_2 = M \mid X_1 = M] \to p$ signifying that the probability of a miss following a miss is p (as if they were independent). We will determine the value for a by fitting the individual hit/miss probabilities for the cluster = 2 data (spectrogram) and use this value of a to predict the hit/miss patterns for larger miss clusters.

The same technique is used to increase the probability of a hit following a hit. Again we use a (the same a used in boosting the probability of a miss) and define a convex combination between that $\Pr[X_2 = H | X_1 = H] =$ (1-p)and 1 such (1-p)a + (1-a) = (1-ap). This also has the properties we desire, for $a \ge 0$ and $p \le 1$, $(1-ap) \ge (1-p)$. Additionally, as $\rho \to 1$, $\Pr[X_2 = H | X_1 = H] \to 1$ (perfect correlation) and as $\rho \to 0$, $\Pr[X_2 = H | X_1 = H] \to p$ (as if they were independent). The remaining two hit/miss probabilities are defined using the HH and MM probabilities. That is, MH has $\Pr[X_2 = H | X_1 = M] = 1 - \Pr[X_2 = M | X_1 = M] = a - ap$ and HM has $\Pr[X_2 = M | X_1 = H] = 1 - \Pr[X_2 = H | X_1 = H] = ap$.

Figure 11 shows each of the four hit/miss combinations as a state transition matrix. Each entry shows the transition probability for one of the four hit/miss combination in the sequence $X_i X_{i+1}$. Using Figure 11, we see that a MM sequence has probability p(1-a+ap), the MH probability p(a-ap), HM probability (1-p)(ap),and HHprobability (1-p)(1-ap). From the simple assumption used to define a single parameter (a) to increase the HH or MM probabilities, we are able to generate the probabilities of larger hit/miss sequences (using Figure 11). For example, three hit/miss sequences HMM, MHHM, and HMMMH (representing the hit/miss patterns for L1 misses that hit/miss in the L2) have the following probabilities: (1-p)(ap)(1-a+ap),p(a-ap)(1-ap)(ap), and, $(1-p)(ap)(1-a+ap)^2(a-ap)$. Hit/miss probabilities of any length can be calculated using the generating function $(1-p)(1-a+ap)^w(a-ap)^x(1-ap)^y(ap)^z$ (if the hit/miss pattern begins with а hit) or $p(1-a+ap)^{w}(a-ap)^{x}(1-ap)^{y}(ap)^{z}$ (if the sequence begins with a miss). Exponents w, x, y, and z then have the property that w+x+y+z=C-1, where C is the length of the cluster.



8. SIMULATION RESULTS

To test the robustness of the theory, each workload was simulated twice: first with a 256KB L2, and then with a 1MB L2. (Additionally, mcf, oltp, oltp2, and perf1 were simulated a third

time using a 2MB L2 for 20 simulation runs in all.) The L1 was 64KB in all cases. The spectrograms, along with the probabilities for each hit/miss combination were produced for clusters 1 through 5. The value for a was determined by least square fit using the four hit/miss equations from the cluster size = 2 data. The global miss probability p was used in all cases along with the hit/miss state transition defined in Figure 11. Using the oltp workload (shown in Figure 3), Table 1 shows the results of the fit. The best fit value for a = .73. As can be seen, the least squares fit is very good. Next, a and the global miss ratio p are

Cluster = 2 HH		HM	MH	MM	
Real	.306	.175	.188	.331	
Predicted	.308	.182	.182	.328	

Table 1. Hit/Miss Probabilities for Cluster = 2. Data is For the OLTP workload. L1=64KB, L2=256KB

used to predict the hit/miss probabilities for larger cluster sizes. Table 2 shows the eight hit/miss probabilities for cluster size = 3. As can be seen the agreement between theory and experiment is quite good. Next we simulated all eight workloads varying the L2 cache size. We determined a using the cluster = 2 data and then predicted the 2^N hit/miss probabilities for clusters 3, 4 and 5. We

CL 3	ннн	HHM	HMH	HMM	MHH	MHM	MMH	MMM
Real	.208	.105	.090	.098	.103	.094	.083	.219
Pred	.193	.115	.065	.117	.115	.068	.117	.210

Table 2. Hit/Miss Probabilities for Cluster = 3. Data is For the OLTP, L1=64KB, L2=256KB

measured the error between our prediction and the true hit/miss probabilities using the Kolmogorov-Smirnov (K-S) [23] test, where the K-S test measures the maximum difference between two cumulative distributions $T_N(x)$ (the true distribution) and $P_N(x)$ (the predicted distribution) max | $T_N(x) - P_N(x)$ | over all x. To represent each of the 2^N hit/miss probabilities as a distribution we ordered the



Figure 12. Kolmogorov-Smirnov Test For Predicted Hit/Miss Probabilities For Clusters = 2, 3, 4, and 5.

9

probabilities from 0 to $2^N - 1$, based on a hit = 0 and a miss = 1. Thus, in a cluster size = 3, HHH=0 and MMM=7.

Figure 12 shows the distribution of the maximum error between 20 simulation runs. The cluster = 2 results are included to indicate the 'goodness' of the fit in determining a. As can be seen, the predictions are very good. Cluster = 3 results show nearly 75% of the simulation runs had less than 7% error.

As the cluster size increases, the amount of error grows. However, even in the cluster = 5 predictions, over two thirds have less than a 10% maximum error.

Once the 2^{N} hit/miss probabilities are determined (predicted), we only need to determine the four (MM, HM, MH, and HH) overlap/no-overlap probabilities for the cluster = 2 case (as shown in Figure 8), then a spectrogram for a larger cluster can be constructed. Intuitively, we believe that the same forces (program dependencies, micro architecture features) that determine the overlap probabilities for the cluster = 2 case will also be present and determine the overlap probabilities in larger clusters.

Again, we simulated all workloads and determined the overlap/no-overlap probabilities from the cluster = 2 case. Next, we produced a overlap/no-overlap tree for each of the 2^N hit/miss combinations for that cluster.

Figure 13 shows the cluster = 3 trees for the HHH and MMH combinations using values from the simulation run described in Figures 3, and Tables 1 and 2. The top tree shows the transitions for the HHH miss pattern. At the root node, the overlap/no-overlap probabilities are shown for X_1X_2 . A left branch indicates that the two misses were overlapped, a right branch indicates no overlap. The second level of the tree shows the overlap probabilities for X_2X_3 , and the leaf nodes show the



Figure 13. Overlap/No-Overlap Probability Tree For HHH and MMH. Data For OLTP

calculations for determining the probabilities for miss penalties 15, 30, and 45 cycles. Note, each node uses the predicted probability for the HHH pattern (.193), shown in Table 2.

The bottom tree shows transition probabilities for the MMH pattern. Note, that miss pattern X_1X_2 uses the MM overlap probabilities, while sequence X_2X_3 uses the MH overlap probabilities. Thus, the overlap probabilities for miss sequence X_iX_{i+1} are drawn from its unique hit/miss pattern (either HH, HM, MH, or MM). The leaf nodes show the calculations to produce miss penalties of 100, 115, 200 and 215.

Using this technique on all of the 2^N hit/miss combinations with their predicted probabilities along with the four hit/miss overlap/non-overlap probabilities from a cluster size = 2, a spectrogram for a larger cluster is produced. Figure 14 shows our

15, 30 45, 100 ...). In future work we plan a theoretical treatment of the position for each peak.

As can be seen, the agreement between theory and experiment is very good. The difference between our prediction and the true probability of a miss pattern (the error) agrees closely with the errors produced in estimating the 2^N hit/miss probabilities shown in Figure 12.

9. INSTRUCTION SPECTROGRAM

In the next set of experiments, we set the instruction cache to 64KB, L2=256KB with a 15 cycle latency, L3=1MB 75 cycles away, and a 300 cycle memory latency (the reverse of the experiment shown in Figure 6). Instruction fetching is guided by a 32K-entry branch target buffer (BTB) that runs well ahead of instruction fetching and the decoder, predicting branches and aiding prefetching. The benefits of using a BTB to guide



prediction for the cluster size = 3 and 4 spectrograms for the oltp workload shown in Figure 3. The red series with symbol marker is the measured spectrogram, while the blue series is our predicted spectrogram. Since our theory only predicts the area of a peak, we plot our predictions as a pulse under the true spectrogram for that cluster. The height of the peak is scaled to match the height of the true spectrogram according to the amount of error between our prediction and true measurement. For example, if there is a 10% error between our prediction and the true area under a peak, then there is a 10% difference between the heights of the peaks. Also, we have not attempted to calculate the position of the peak, but instead use the miss latencies to generate their position, (i.e.. instruction fetching and prefetch instructions has been well documented [24-28].

Figure 15 shows the instruction spectrogram for cluster sizes = 1, 2, 3, and 4 for oltp3. These spectrograms are different from their data counterparts. The spectrograms for cluster sizes 1, 2, and 3 have dominant peaks at 15, 75, and 300, with few misses (clusters) costing more than 300 cycles. The amount of overlap, in the cluster = 2 and 3 spectrograms, is substantial.

The cluster size = 2 spectrogram is still dominated by peaks at 15, 75 and 300, even though there are two misses. Obviously, one of the misses is overlapped with the first. In the cluster size = 3 spectrogram, there are four dominate peaks: 15, 30, 75, and 300.



Here we have three misses but only the peak at 30 indicates that two out of the three misses were not overlapped. In the other cases (peaks at 15, 75, and 300), two out of the three misses were overlapped with the first. Recall that the data spectrograms for these cluster sizes had miss penalties extending to 600 and 900 cycles, respectively. Even the cluster size = 4 spectrogram has peaks at 15, 75, and 300 cycles, indicating three out of the four misses were overlapped. Analyzing the peak at 150 indicates two of the four misses were resolved in the L3, while the other two misses were overlapped. Similarly, the peak at 375 indicates one miss when to memory and the other was resolved in the L3 (the other two misses were overlapped).

Notice the large left shoulders on each of the peaks in the cluster = 1 spectrograms. There is also a large peak at 0 indicating that nearly 25% of the instruction misses had zero delay. Even the spectrograms for cluster sizes = 2, 3 and 4 show peaks at zero. For example, the cluster = 4 spectrograms shows approximately 6% of the miss clusters had zero cycles penalties. Here we have four misses that occur in parallel but yet the BTB was able to prefetch them far enough in advance of their use to cause zero cycles difference between the finite cache simulation run and an infinite cache simulation run.

In Figure 16, we change the range of the X and Y axis of the cluster=1 spectrogram to emphasize the left shoulder of the 15 and 75 peaks (zoom in). We see that the left shoulder of the peak



at 75 extends to nearly 20 (indicating that prefetches are issued nearly 55 cycles ahead of their use).

10. SPECTROGRAM SHAPE AND BUS OUEUEING

In our last experiment we highlight the effects bus queueing has on the shape of a spectrogram. We repeat the experiment shown in Figure 3 but use the oltp3 workload and change the memory latency to 50 cycles and model a 16 byte bus that transfers a packet of information every other cycle, as opposed to the optimistic 128B bus transfer used previously. Now, it takes 16 cycles to transfer a line between the caches. Figure 17 shows the spectrograms for clusters = 1, 2 and 3. Notice the effects that bus queueing has on the shape of the right shoulder.

In the cluster = 1 spectrogram, we clearly have peaks centered at 15 and 50 cycles, but the peaks are not as sharp as when the line was transferred in one cycle. However, in the cluster = 2 spectrogram, the original five peaks shown in Figure 3, broaden and merge. Now, there are three large peaks, (not five as seen earlier) with one peak between 30 and 60, one between 60 and 90, and one between 90 and 120 with a large right shoulder. Obviously, there are four hit/miss combinations for this cluster but the cost of each miss cannot be identified as clearly as when the bus transfer interval was 1 cycle. Finally, in the cluster = 3 spectrogram, all three peaks start to merge between a range of 30 to 180. Obviously queueing is increasing the cost of each miss or miss cluster well beyond the nominal miss latency.



11. SUMMARY

A new technique has been presented for calculating the cost of a miss and displaying images that represent their cost. We call this technique pipeline spectroscopy. The underlying principles of this technique are very simple: the cost of a miss can be determined by knowing the finite cache and infinite cache execution times for the same sequence of instructions. The difference between these two times is the cost of the miss (cluster).

We use this principle to produce a miss spectrogram, which represents a precise readout of the cost of every miss. A miss spectrogram has enormous value in analyzing the performance of an application or microarchitecture. Detailed analysis of a spectrogram leads to insights in pipeline dynamics, including effects due to prefetching, bus queueing, and underlying architectural features that allow or inhibit memory level parallelism.

We also presented a theory that describes the observed properties that make up a spectrogram. The theory uses 6 parameters to predict the shape of a spectrogram: the global miss probability for the L2 (p), a correlation parameter (a), and four hit/miss overlap/no-overlap probabilities drawn from hit/miss combinations found in the cluster size = 2 misses (from an application). We applied this theory and were able to predict the miss patterns for larger miss clusters.

In this study, we demonstrated that pipeline spectroscopy can be used to explore the amount of memory level parallelism an application can achieve. Future work is needed to study in-order versus out-out-order effects on MLP, as well as SMT processor organizations, There also appears to be much more information in the shape of the spectrogram than just the cost of a miss. Additional work is needed to understand all of the peaks and sub-peaks that are visible in a spectrogram (left and right shoulders), as well as understand the amount of cycles a peak shifts.

12. REFERENCES

- [1] A. Glew, "MLP yes! ILP no!," in ASPLOS Wild and Crazy Idea Session, October 1998.
- [2] V. Pai and S. Adve, "Code Transformations to Improve Memory Parallelism," in 32nd International Symposium on Microarchitecture, November 1999.
- [3] H. Zhou and T. Conte, "Enhancing Memory Level Parallelism via Recovery-Free Value Prediction," in International Conference on Supercomputing, June 2003.
- [4] D. Sorin et al, "Analytic Evaluation of Shared-Memory Systems with ILP Processors," in 25th International Symposium on Computer Architecture, 1998.
- [5] V. Pai, P. Ranganathan and S. Adve, "The Impact of Instruction- Level Parallelism on Multiprocessor Performance and Simulation Methodology," in HPCA February 1997.
- [6] P. Ranganathan, K. Gharachorloo, S. Adve and L. Barroso, "Performance of Database Workloads on Shared-Memory Systems with Out-of-Order Processors," in ASPLOS-VIII, 1998.
- [7] Y Chou, B. Fahs, and S Abraham, "Microarchitecture Optimizations for Exploiting Memory-Level Parallelism Exploiting Memory-Level Parallelism" in 31st International Symposium on Computer Architecture, 2004.
- [8]Yuan Chou, Lawrence Spracklen, Santosh G. Abraham. "Store Memory-Level Parallelism Optimizations for Commercial Applications," pp. 183-196, 38th MICRO 2005.
- [9] M. Qureshi, D. Lynch, O. Mutlu, Yale Patt, "A Case for MLP-Aware Cache Replacement" in 33rd ISCA June 2006
- [10] A. Zahir, V. Hummel, M. Kling, T Yeh, US. Patent 6,353,802, "Apparatus and Method for Cycle Accounting in Microprocessors"

- [11] B. Gaither, R. Smith, US Patent 6,892,173 B1, "Analyzing Effectiveness of a Computer Cache By Estimating a Hit Rate Based on Applying a Subset of Real-time Addresses to a Model of the Cache"
- [12] H. Ravichandran, US Patent 6,341,357 B1, "Apparatus and Method for Processor Performance Monitoring",
- [13] R. Trauben, US Patent 5,594,864, "Method and apparatus for unobtrusively monitoring Processor States and Characterizing Bottlenecks in a Pipeline Processor Executing Grouped Instructions"
- [14] G. Brooks, US Patent 5,845,310 "System and Methods For Performing Cache Latency Diagnostics in Scalable Parallel Processing Architectures Including Calculating CPU Idle Time and Counting Number of Cache Misses.
- [15] W. Flynn, US Patent 6,256,775 B1, "Facilities For Detailed Software Performance Analysis in a Multithreaded Processor"
- [16] F. Levine, B. McCredie, W. Starke, E. Welbon, US Patent 5,862,371, "Method and System for Instruction Trace Reconstruction Utilizing Performance monitor outputs and bus Monitoring"
- [17] J. Dean, J. E. Hicks, C. A. Waldspurger, W. E. Weihl, and G. Z. Chrysos. *ProfileMe : Hardware support for instruction-level profiling on out-of-order processors*. In MICRO'97: pages 292--302, 1997.
- [18] Brian A. Fields, Rastislav Bodik, Mark D. Hill, Chris J. Newburn., Interaction cost and shotgun profiling. ACM Tracsactions on Architecture and Code Optimization, Vol 1, No. 3. Sept 2004.
- [19] Tejas Karkhanis, James E. Smith, A First-Order Superscalar Processor Model. Proceedings of the 31st ISCA. pages 338–349, June 2004.
- [20] A. Hartstein and T. Puzak. The optimum pipeline depth for a microprocessor, 29th ISCA, pages 7-13 May 2002.
- [21] A. Hartstein and T. Puzak. Optimum power/performance pipeline depth. 36th Annual IEEE/ACM International Symposium on Microarchitecture In *MICR*O, Dec. 2003.
- [22] P.J. Denning, "The Working Set Model for Program Behavior", CACM 19(5) pp.285-294 (1976).
- [23] R. Bartoszynski, M Niewiadomska-Bugaj, Probability and Statistical Inference, (Wiley series in probability and statistics) 1996
- [24] J.S. Liptay, "Design of the IBM Enterprise System/9000 High-End Processor," IBM Journal of Research and Development, Vol. 36, No. 4, pp. 713-731, July, 1992.
- [25] D.B. Fite, J.E. Murray, D.P. Manley, M.M. McKeon, E.H. Fite, R.M. Salett, and T. Fossum, "Branch Prediction," U.S. Patent #5142634, assigned to Digital Equipment Corporation, Filed Feb. 3, 1989, Issued Aug. 25, 1992.
- [26] N. Suzuki, "Microprocessor Having Branch Prediction Function," U.S. Patent #5327536, assigned to NEC Corporation, Filed May 22, 1991, Issued July 5, 1994.
- [27] C.H. Perleberg, and A.J. Smith, "Branch Target Buffer Design and Optimization," *IEEE Transactions on Computers*, Vol. 42, Issue 4, pp. 396-412, Apr., 1993.

[28] B. Calder, and D. Grunwald, "Fast and Accurate Instruction Fetch and Branch Prediction," 21st ISCA pp. 2-11, April 18-21, 1994.

Appendix A.

In a memory hierarchy with N cache levels (L1, L2, ..., L_N, and memory) and a miss cluster of size C, there are $\binom{C+N}{C}$ possible miss penalties. Let each level of the memory hierarchy be represented by a distinct (non-multiple) miss penalty (number), there are N of them. We represent this problem as sampling with replacement to determine the number of unique combinations (sums) using these numbers. It is sampling with replacement because there is an inexhaustible supply of miss latencies regardless of the cluster size. We use the notation $\begin{bmatrix} N \\ C \end{bmatrix}$ to denote N items choose C (the cluster size) for sampling with replacement. The problem can then be expressed as determining the number of unique sums from N items as we vary the number of picks from 0 to C. Note that i and C can be greater than N in (1a) because sampling is done with replacement.

$$\sum \begin{bmatrix} N \\ i \end{bmatrix} 0 \le i \le C = \begin{bmatrix} N \\ 0 \end{bmatrix} + \begin{bmatrix} N \\ 1 \end{bmatrix} + \begin{bmatrix} N \\ 2 \end{bmatrix} + \begin{bmatrix} N \\ 3 \end{bmatrix} +, ..., + \begin{bmatrix} N \\ C \end{bmatrix}$$
(1a)
From [23] $\begin{bmatrix} N \\ k \end{bmatrix} = \begin{pmatrix} N+k-1 \\ k \end{pmatrix}$ so we can rewrite (1a) as
 $\begin{pmatrix} N-1 \\ 0 \end{pmatrix} + \begin{pmatrix} N \\ 1 \end{pmatrix} + \begin{pmatrix} N+1 \\ 2 \end{pmatrix} + \begin{pmatrix} N+2 \\ 3 \end{pmatrix} ... \begin{pmatrix} N+C-1 \\ C \end{pmatrix}$ (2a)
Also, from [23] $\begin{pmatrix} N \\ k \end{pmatrix} = \begin{pmatrix} N-1 \\ k-1 \end{pmatrix} + \begin{pmatrix} N-1 \\ k-1 \end{pmatrix} + \begin{pmatrix} N-1 \\ k \end{pmatrix}$ (3a)
So we can combine the first two terms of (2a) and obtain
 $\begin{pmatrix} N+1 \\ 1 \end{pmatrix} + \begin{pmatrix} N+1 \\ 2 \end{pmatrix} + \begin{pmatrix} N+2 \\ 3 \end{pmatrix} ... \begin{pmatrix} N+C-1 \\ C \end{pmatrix}$ (4a)

Applying (3a) repeatedly the series collapses and the desired sum is produced $\binom{C+N}{C}$.

Appendix B.

We show that the correlation coefficient $\rho = 1 - a$. We represent hits as 0 and misses as 1 such that $X_i = M$ or $X_i = H$ is equivalent to $X_i = 1$ or $X_i = 0$, respectively. Let $\Pr[X_1 = 1] = p$. By definition the correlation coefficient

$$\rho = \frac{COV(X_1X_2)}{\sqrt{Var(X_1)Var(X_2)}}.$$

$$E[X_1X_2] = \Pr[X_1X_2 = 1] =$$

$$\Pr[X_1 = 1, X_2 = 1] = p(1 - a + ap)$$

$$E[X_1] = E[X_2] = p.$$

-- -

$$COV(X_{1}, X_{2}) = p(1 - a + ap) - p^{2} =$$

 $(1 - a)[p(1 - p)]$

$$Var(X_1) = Var(X_2) = p(1-p)$$

So $\rho = \frac{(1-a)[p(1-p)]}{p(1-p)} = 1 - a$