# Constant Load Value Prediction

Prajyot Gupta
*Electrical and Computer Engineering*
*University of Wisconsin - Madison*
Madison, WI, USA
pgupta54@wisc.edu

Satvik Maurya
*Electrical and Computer Engineering*
*University of Wisconsin - Madison*
Madison, WI, USA
smaurya@wisc.edu

Robert Viramontes
*Electrical and Computer Engineering*
*University of Wisconsin - Madison*
Madison, WI, USA
rviramontes@wisc.edu

*Abstract*— **Locality is a crucial technique that drives performance and energy optimizations in the modern computing system, particularly as the performance of CPU cores and memory devices diverge. Value locality is a form of locality that occurs when a value at a memory address does not change [1]. Our project takes advantage of value locality by implementing constant value prediction. In constant value prediction, we add hardware that monitors the memory access stream and determines memory addresses where the value is highly stable. When a highly stable load is identified, our prediction hardware provides the value and allows the load instruction to skip accessing the memory hierarchy. We implement constant value prediction in the Minor in-order and DerivO3 out-of-order CPU models of the gem5 simulator to understand the impact of constant value prediction.**

**We conduct experiments utilizing two benchmark suites, SPEC2006 and Perfect. Our experiments explore several variations of the prediction hardware at different table sizes and under different replacement policies. An encouraging result of our simulations is a maximum of 50% reduction in data cache accesses for a relatively small hardware configuration. While this does not translate to significant performance increase via IPC, it may lead to energy savings by replacing portions of cache traffic with a lightweight prediction structure.**

## I. Introduction

The disparity between the increase in speed of a microprocessor and the access speed of an off-chip memory is one of the major aspects of microprocessor design. Loads that miss the last level cache can take hundreds of cycles to fetch the data, leading to long stalls in the processor pipelines. This problem is caused because the differential rate of change in speeds between the processor and memory is known as the "*Memory Wall*", wherein processor speed improvements range 35-55% and on the other hand, memory latency improvements amount to about 7% [2]. Processors make the use of data and instruction level parallelism techniques to hide the stalls. However, this solution still suffers from a fundamental bottleneck – limits in the off-chip data transfer bandwidth in case of accessing main memory. Modern-day workloads from several prevalent and emerging application domains, such as image processing, computer vision, machine learning, and data analytics, often have input datasets with significant repetition of values. This presents us an opportunity to leverage the inherent predictability of data values in these applications to tackle the bottleneck. This paper exploits this opportunity.

In our work, we explore the concept of load value locality specifically through the lens of constant value prediction. The concept of value locality was discussed by Lipasti et al. [1] as the "likelihood of the influence of previously seen program values on the newly occurring ones". In ECE 752, we discussed spatial and temporal locality, which attempted to improve the probability of finding a given memory location in the cache by reducing the distance between successive accesses to that location. Value locality is a different approach to tackle the same problem. Using this concept, the processor can speculate what value will be returned by a load instruction, and have it execute the instructions speculatively using this value. When the load completes, the actual value returned by it can be compared with the speculated value. If the speculation is correct, the processor increases the confidence of its predictions and continues; but if the guess is incorrect, the processor rolls back the computation and recomputes using the actual value obtained via the load instruction.

Contemporary work, like Lipasti [1], Seznec [5] show that value prediction can result in high accuracies for certain workloads using relatively simple prediction schemes. One important point to note in the work here is that when we talk about value prediction, we always discuss the prediction of value for a load operation. A store operation that updates a base register is quite unlikely to exhibit any locality. On the other hand, a load is very likely to exhibit locality because any value stored in memory will likely be accessed more than once.

In constant value prediction, we introduce hardware to monitor all load operations. Then, loads that have a very consistent value are predicted as *Constant* and skip accessing the memory system. We have re-implemented the ideas from [1] and run them on more modern workloads to see how the concepts of value locality hold for modern workloads. These workloads include the SPEC2006 [3] and Perfect [4] benchmark suites. We modify the gem5 in-order MinorCPU and out-of-order DerivO3 models [5] to incorporate our load value prediction hardware.

Our code is available on GitHub:
https://github.com/robertviramontes/ece752_loadvaluepredictor

## II. Implementation

Our load value prediction hardware follows the design of the hardware proposed in [1]. As such, we will keep the discussion of the hardware attributes brief and here focus on the specifics of our implementation methodology and lessons learned during implementation. We have implemented the load value prediction features in a single Load Value Prediction Unit, LVPU, that serves as the interface between the CPU

pipeline and the subunits that implement constant load value prediction, which are described below.

## A. Load Value Prediction Table

The load value prediction table, also referred to as LVPT, stores the value to predict for loads. It takes in the PC value of the load instruction and outputs the predicted data whose value was previously loaded by the instruction. Its architecture is similar to a directly mapped cache – it is indexed by the load instruction address and is untagged. Because of this, both constructive and destructive interference can occur between the load instructions with the same index value, and hence can contain aliases. It can have a history depth of *n,* which provides a mechanism for selecting the value to predict from a set of *n* unique values for a given instruction.

The LVPT stores a valid tag to indicate that the storage location contains valid data from runtime observations. This is cleared at processor startup and set true when an instruction indexes in to the row to update the stored value. If the LVPT determines that an instruction value is invalid, then the lookup in the load classification table is skipped because an invalid value is always *Strongly Unpredictable*.

## B. Load Classification Table

The load classification table, LCT, generates a classification for a load. The LCT is implemented as an array of saturating counters, which allowed us to base the software design on the similar Simple branch predictor available in gem5. As in [1], the saturating counters are indexed by the instruction address. The possible classifications in our LCT include *Strongly Unpredictable*, *Weakly Unpredictable*, *Predictable*, and *Constant*. A classification of *Strongly Unpredictable* is generated when all counter bits are 0, *Weakly Unpredictable* when the most significant bit is 0 and at least one other bit is 1, *Predictable* when the most significant bit is 1 and at least one other bit is 0, and *Constant* when all bits are 1. Although our current system implementation only utilizes a *Constant* classification, the remaining classification types are kept for statistic collection and future incorporation of modifications that can take advantage of other classifications.

The LCT is updated when the LVPU receives a request to verify a prediction as a memory transaction is committed. If the prediction was correct, the corresponding counter is incremented, and if incorrect, the counter is decremented. We choose this policy because it is straightforward to implement and reason about. However, we did discuss an alternate policy that would reset the counter to zero on a misprediction. The intuition behind this policy is that if a value is changing, particularly one that was previously *Constant*, this may represent a change in the instruction pattern where the previous value's predictability is no longer applicable. The tradeoff between these policies would be interesting to explore in future work, particularly in a system that utilizes the *Predictable* classification.

## C. Constant Verification Unit

The CVU is responsible for maintaining coherency between the value of a constant load stored in the LVPT and the actual value in memory. The CVU achieves this by using the following mechanism- When a load instruction graduates from being *Predictable* to *Constant* the LCT, a new entry is allocated in the Content Addressable Memory (CAM) of the CVU. This CAM is a fully associative cache in which every entry consists of the effective virtual address of a constant load and the corresponding index of that constant load in the LVPT. Thus, if a load address is present in the CAM of the CVU, that load is constant and the value in the LVPT is coherent with memory. The CVU then monitors all store instructions that reach the commit stage of the pipeline. If an entry in the CVU CAM has a memory address which matches that of a committing store instruction, the CVU invalidates that entry in its CAM which means that that load address is no longer coherent with the entry in the LVPT. Loads are classified in the frontend of the pipeline, but the verification of the classified loads takes place later. This is primarily because the effective address of the load instruction must be calculated before it can be verified whether the load value predicted by the LVPT is coherent with memory.

The CVU is also responsible for checking loads that are marked as *Predictable* by the LCT. *Predictable* loads execute in the pipeline normally and after they finish execution, the actual value of the load instruction is compared with the value predicted by the LVPT. If these values match, the LCT increments its saturating counter and if the values do not match, the LCT decrements its saturating counters.

Since the CVU consists of a fully associative CAM, the replacement policy for its entries, once the memory is full, is an important aspect to investigate. Fully associative memories are expensive, fully associative CAMs even more so. Hence, the number of entries available in the CVU CAM is limited which implies that the entries should be managed as efficiently as possible. Thus, we have explored five different replacement policies in the CVU CAM-

- FIFO: This is a simple queue-based policy where the oldest entry is ejected from the CAM.
- Most Recently Used (MRU): This policy ejects the most recently used entry from the CAM.
- Not-Most Recently Used (NMRU): This is a simpler policy than the MRU policy implementation-wise. This ejects any entry from the CAM which has not been used most recently.
- Least Recently Used (LRU): This policy ejects the least recently used entry from the CAM.
- Not-Most Recently Used (NLRU): Like the NMRU policy, this is a simpler implementation, and it ejects any entry from the CAM which is not the least recently used entry.

## D. Minor CPU Modifications

We modify the gem5 Minor CPU model to understand the impact of the constant value prediction in an in-order CPU. The

key modifications include modifications to the dynamic instruction class, the decode stage, the load store queue and the commit stage. We also disable the early memory reference option to ensure the Minor CPU processes memory references in program order.

The dynamic instruction is modified to contain additional metadata about the instruction, including the classification from the LCT and the value predicted by the LVPT. We also add the virtual address as a property of the dynamic instruction, which was already present in the DerivO3 dynamic instruction.

The decode stage is modified so that all load instructions check for a prediction in the LVPU. The dynamic instruction's metadata is updated with the classification and value that is returned from the LVPU at this stage.

Memory operations proceed normally into the load store queue. When the instruction enters the queue, its virtual address is available, so the instruction metadata is updated to carry the virtual address through the rest of the execution pipeline, similar to the approach already implemented in the DerivO3 CPU. When a load instruction is at the head of the queue, we check if it was classified as *Constant*. If it was, then we call the LVPU to verify with the CVU that the value is still valid. If the constant load is still valid, the memory request is immediately fulfilled with the predicted value from the dynamic instructions metadata. This allows the constant predictions to avoid memory access. When a store instruction reaches the head of the queue, its PC and virtual address are forwarded to the CVU to invalidate subsequent constant loads.

The commit stage is modified so that all load instructions, regardless of classification, are sent to update predictions in the LCT and LVPT.

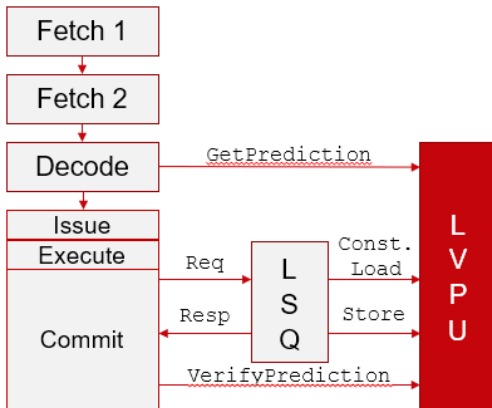The modifications to the Minor CPU model are summarized in Figure 1.



*Figure 1 Graphical depiction of modificatons to Minor CPU*

### E. DerivO3 CPU Modifications

The DerivO3 CPU has a pipeline which is summarized in Figure 2. The integration of the load value prediction unit in the out-of-order CPU supported in gem5 faced the following challenges-

- Register renaming: Because the out-of-order CPU renames architectural registers as arbitrary physical registers before issuing instructions to eliminate false dependencies, the load value prediction unit cannot forward values of constant loads before register renaming is done.
- Load-Store ordering violations: The out-of-order CPU also has a memory dependency checker, which enables load instructions to execute out-of-order with stores. This adds an additional overhead of correcting ordering violations in case a younger load is executed before an older store instruction. Thus, the load value prediction unit will have to be mindful of this nuance since a constant load may reach the execute stage (where the fun stuff happens, to be described below) before an older store. This implies that the safest way to ensure functional sanity is by tightly coupling the load value prediction unit with the load-store queue unit.

These challenges were thus overcome by the following implementation steps-

- Load instruction addresses looked up in the LVPT and LCT during dispatch – instructions are still in-order when they reach dispatch. At this stage, the destination registers of all instructions have been renamed. Once the LCT and LVPT return the classification and the predicted value respectively, the load instruction is pushed into the load-store queue and the instruction queue normally. The predicted load value and the classification (not predictable, predictable, constant) are attached to the instruction itself.
- Irrespective of their classification, all load instructions flow through the issue stage of the pipeline normally. After a load instruction starts execution, the CVU lookup for a load classified as constant starts as soon as the effective load address has been calculated. If the CVU verification is successful, the load instruction does not trigger a memory read and is sent to the writeback stage of the pipeline. For all other cases, the load instructions will access memory.
- At writeback, the constant value predicted by the LVPT is written into the destination register of the load instruction. This is safe to do here because all memory ordering violations would have been detected earlier. The writeback stage of the pipeline is also the stage where other loads update the LCT. Loads classified as predictable are verified by the CVU and the result of the verification is used to update the LCT.
- The invalidation of CVU entries by store instructions occurs after stores have executed and have been sent for commit. Since these stores will commit at some point in the future, any constant loads with the same effective memory address in the LVPT will be incoherent with memory.
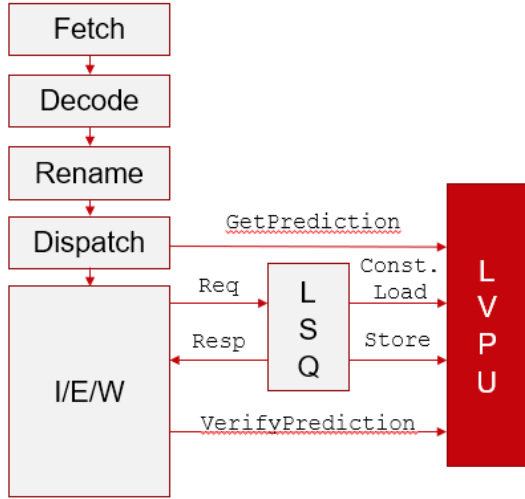
*Figure 2 Graphical depiction of DerivO3 modifications*

## III. RESULTS

In this section, we will present the results of our simulations and highlight some notable datapoints. We save our discussion of the results for section IV.

### A. Simulator Configuration Overview

Our implementation of LVPU is simulated in gem5 [5] using X86 ISA using in-order Minor and out-of-order DerivO3 CPU. The simulator configurations can be observed in *Table I* which are consistent with Seznec's implementation of value predictor in [6].

### B. LVPU Variations

In our work, we explore 8 variations of the LVPU hardware to understand the effects of various parameters on the system. These variations are described in TABLE *II*TABLE III. The parameters can be set via command line arguments to the gem5 system emulation script (`se.py`).

### C. Coverage and Accuracy

Two metrics of interest for load value prediction are the coverage and accuracy of the predictor [6]. The coverage is the number of loads that the predictor classifies as predictable. The accuracy is the number of loads classified as predictable that are correctly classified. Figure 3 shows the coverage and accuracy for selected benchmarks in the Minor CPU and Figure 4 shows this for the DerivO3 CPU.

In the minor CPU, the `lbm` benchmark has the highest coverage with a peak of 56.65% for Variation-0 while `sjeng` exhibits low value locality with a peak coverage of only 1.36%. The accuracy varies significantly, with `libquantum` reaching the highest accuracy of 99.90% in Variation-1. Predictably, the accuracy is worst for the `specrand` benchmark, which heavily relies on random number generation [3], with the minimum accuracy of 1.55%.

*TABLE I* SIMULATOR CONFIGURATION

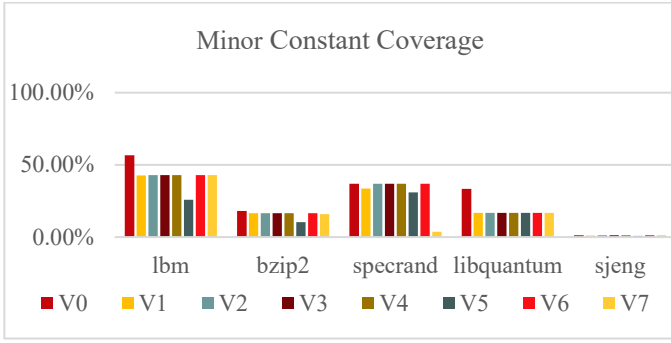| Front-End | **Fetch through Rename width:** 8 insts/cycle<br>**Branch Pred.:** TournamentBP, size 8K<br>**BTB:** 2-way 4k-entry, Tag size 16<br>**RAS:** size 16 |
|---|---|
| Cache | **L1 D$:** Size = 32kB, Associativity = 8;<br>**L1 I$:** Size = 32kB, Associativity = 8;<br>**L2$:** Size = 256kB, Associativity = 8; |
| Main Memory | **Type** = DDR3_2133_8x8, **Size**=16GB |

*TABLE II LVPU VARIATIONS*

| Var. | LCT | | LVPT | | CVU | |
|---|---|---|---|---|---|---|
| | Entries | Counter | Entries | Hist. depth | entries/ thread | Rep. policy |
| 0 | 512 | 2-bit | 1024 | 1 | 8 | FIFO |
| 1 | 512 | 4-bit | 1024 | 1 | 8 | FIFO |
| 2 | 1024 | 2-bit | 4096 | 16 | 32 | FIFO |
| 3 | 1024 | 2-bit | 4096 | 16 | 32 | LRU |
| 4 | 1024 | 2-bit | 4096 | 16 | 32 | NLRU |
| 5 | 1024 | 2-bit | 4096 | 16 | 32 | MRU |
| 6 | 1024 | 2-bit | 4096 | 16 | 32 | NMRU |
| 7 | 1024 | 4-bit | 4096 | 16 | 32 | FIFO |

In the DerivO3 CPU, we see the continuation of some general trends, such as `libquantum` demonstrating high value locality and `sjeng` demonstrating low locality. The highest coverage achieved is 49.81% for `libquantum` Variation-0 and 2-6 while the lowest coverage is 2.51% for bzip2 on Variation-1. The highest accuracy of ~100% is achieved on the `libquantum` benchmarks for Variation-1 and Variation-7 while the lowest is 1.19% for `specrand` on Variation-5.

### D. D$ accesses

L1 data cache accesses for our selected benchmarks, normalized to the number of L1 data cache accesses in the unmodified processor are shown in Figure 5 and Figure 6 for the Minor and DerivO3 CPU, respectively. For the Minor CPU, our simulations show that in a benchmark with high value locality like `libquantum`, the L1 data cache accesses are, in the best case, reduced to 50.0306% of the unmodified Minor CPU accesses in variations V2-V6. In the DerivO3 CPU for the `libquantum` benchmark, the best case reduces the L1 data caches accesses to 74.9090% of the unmodified CPU accesses for variations V2-V4.

Our data verifies that, even in the worst case with benchmarks with poor locality, the LVPU does not increase data accesses. This agrees with our theoretical understanding

*(a) Constant Prediction Coverage*



*(a) Constant Prediction Coverage*



*(b) Constant Prediction Accuracy*

*Figure 3 Constant Coverage and Accuracy for Minor CPU*



*(b) Constant Prediction Accuracy*

*Figure 4 Constant Coverage and Accuracy for DerivO3 CPU*

that constant prediction acts as a highly specialized cache but does not modify program execution.

### E. Instructions Per Cycle

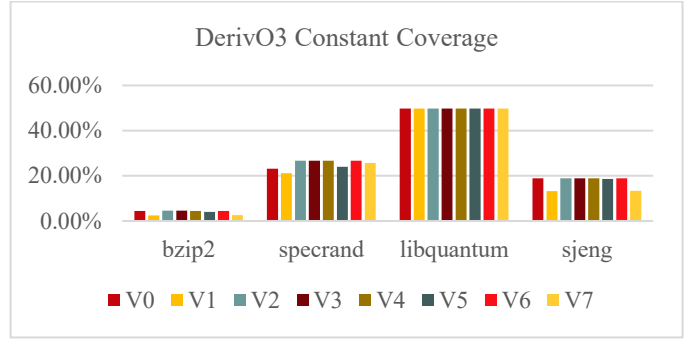We present results for IPC on selected benchmarks in Figure 7 for Minor CPU and Figure 8 for DerivO3 CPU. The IPC is collected from the default gem5 output and normalized to the unmodified processor to clearly show the impact of the change. We find that constant value prediction does not have a significant effect on the overall instructions per cycle, IPC, of a program. We observe that the IPC may have more potential for improvement in the out-of-order CPU than an in-order CPU. In the libquantum benchmark, the out-of-order CPU reaches a normalized IPC of 1.01 while the in-order reaches a normalized IPC of 1.006 even though the in-order skips more data accesses. We see a similar, though not as significant, trend in the specrand benchmark (peak 1.008 vs 1.007).
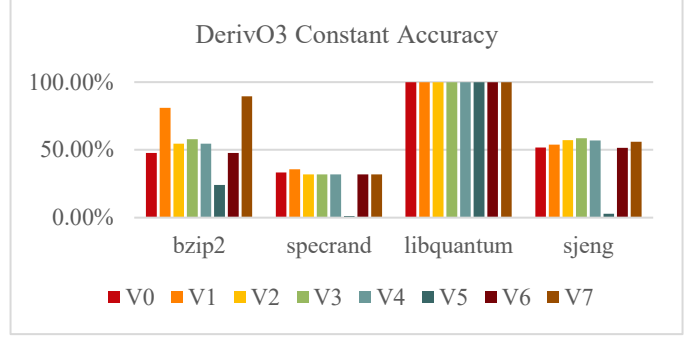
### F. Occurrence of Zeros and Ones

We collect statistics about the actual values that are predicted as constant, specifically tracking constant predictions that are 0 and 1. We collect these statistics at during the lookup process at the decode/dispatch stage for instructions classified as *Constant*, so these represent the values predicted for constant loads regardless of if they are committed or not.

Figure 9 shows the predicted values of Variation-0 in the Minor CPU. The sjeng benchmark has the highest proportion of values predicted as zeros and ones, at 21.9% and 7.6%, respectively, in Minor and 30.2% and 30.2%, respectively, in

DerivO3. However, for most benchmarks the "other" category of values represents a significant proportion of the values that are predicted as constant. Only the sjeng benchmark on the DerivO3 CPU demonstrates a pattern where zeros and ones constitute most of the constant values.

While our data shows that zeros and ones constitute only a small proportion of the constant loads, we would need further study to understand the frequency of individual values which may show that zero and one occur more frequently than any other individual value. However, our data does suggest that there is a benefit to general value prediction and limiting the predictions to zeros and ones would eliminate significant opportunity for value prediction.

### G. Additional Results

In addition to the SPEC2006 benchmark results described above, we also ran the following benchmarks from the Perfect suite [4] of image and signal processing benchmarks with our basic LVPU (LVP_V0)- sar-backprojection, sar-pfa-interp1, stap-inner-product, stap-outer-product, wami-change-detection.

These benchmarks did not show any significant change in the IPC, so we have omitted the IPC chart for the sake of brevity. However, the constant coverage and accuracy charts are shown in Figures 11 and 12 for the Minor and DerivO3 CPUs. The coverage for both CPUs is generally high (compared to SPEC06) for almost all the chosen workloads. However, the constant accuracy graphs in Figure 12 show that despite having a relatively high coverage, the accuracy of predictions was

*Figure 5 Minor CPU data cache accesses, normalized to unmodified CPU for select SPEC2006 benchmarks*



*Figure 6 DerivO3 CPU data cache accesses, normalized to unmodified CPU for select SPEC2006 benchmarks*



*Figure 7 Minor CPU IPC, normalized to unmodified CPU for select SPEC2006 benchmarks*



*Figure 8 DerivO3 CPU IPC, normalized to unmodified CPU for select SPEC2006 benchmarks*

mediocre. We attribute this to the fact that there are a lot of invalidations in the CVU CAM due to store instructions.

The decrease in data cache accesses was as expected from both CPUs as shown in Figure 13.

While we have not shown the results here, an analysis of the distribution of load *values* for the five benchmarks showed that the percentage of constant loads with values of one was very small (< 1%). However, for constant loads with values of zero, there was a slightly different trend, with the `wami-change-detection` workload reporting that about 25% of constant loads had a value of 0.



*Figure 9 Minor CPU constant value predictions categorized as either zero, one or other for select SPEC2006 benchmarks*



*Figure 10 DerivO3 CPU constant value predictions categorized as either zero, one or other for select SPEC2006 benchmarks*

Figure 11 Constant coverage for the perfect-suite workloads



Figure 12 Constant accuracy for the perfect-suite workloads



Figure 13 Normalized D$ accesses for the perfect-suite workloads

## IV. Discussion

### A. Effect of structure size

Our results show that the smallest version of the hardware, Variation-0, performs almost as well as the other variations at reducing data accesses. The values for the normalized data accesses of the `specrand` benchmark are reproduced in TABLE III and show that there is little variation as the hardware size changes. The standard deviation for the normalized accesses in the Minor CPU is $3.0 \times 10^{-6}$ and for the DerivO3 CPU is $1.2 \times 10^{-5}$, providing an empirical measurement of the small variation in data cache accesses across LVPU variations.

This is a useful result because one of the primary motivations of our predictor is to reduce the data cache accesses which may help reduce the power consumed in cache transactions. Our results show that large structures are not required to effectively implement constant prediction and suggest that smaller, more power-efficient structures are suitable. Further study would be required to determine a lower bound on the effective sizes and physical simulation to determine the actual energy cost of LVPU structures.

TABLE III NORMALIZED D$ ACCESSES FOR `specrand`

| Variation | Minor CPU | DerivO3 CPU |
|---|---|---|
| 0 | 0.500307173 | 0.749091022 |
| 1 | 0.500312572 | 0.749122977 |
| 2 | 0.500306173 | 0.749088838 |
| 3 | 0.500306173 | 0.749088838 |
| 4 | 0.500306173 | 0.749088838 |
| 5 | 0.500306173 | 0.749095587 |
| 6 | 0.500306173 | 0.749091022 |
| 7 | 0.500312972 | 0.749099159 |

### B. Effect of replacement policy

We see that overall, the least recently used (LRU) policy performs the best, even though the difference between this policy and other policies like NLRU is very slight. Since the NLRU policy is easier to implement in hardware in terms of area and complexity, this policy is the best for the CVU fully associative CAM.

The use of the MRU policy shows the most drastic degradation in performance and this degradation is more prominent in the out-of-order core than the in-order core. Using the MRU policy results in a direct violation of the principles of value and temporal locality because it discards the most local values. Thus, the use of this replacement policy serves as a proof-of-concept for our load value prediction scheme.

### C. Minimal impact on IPC

We found that constant value prediction had a very small impact on the IPC of the program, at most achieving a ~2.5% gain in IPC. This result is not surprising as the requirements for constant value prediction are multiple accesses to the same memory address with an unchanged value. This access pattern means this data is likely in the data cache, so we are only skipping a hit in the data cache which has a small (1 or 2-cycle) latency in our simulations. This analysis is based on our understanding of the systems and further study would be required to verify this analysis.

Furthermore, our implementation is also conservative as far as value forwarding for constant loads is concerned. The value of a prediction for a constant load is forwarded only during the writeback and commit stages for the O3 and Minor CPUs, respectively. If the value forwarding were done earlier and more aggressively, we might be able to see better results for the IPC, especially for the O3 CPU since there could be multiple dependent instructions waiting in the reservation stations for the load value.

Our results also show that the out-of-order CPU can improve more than the in-order CPU. We attribute this difference to the out-of-order CPUs ability to take advantage of the reduced pressure on the load queue by reordering instructions to take advantage of slots freed by constant prediction.

### D. Coverage and accuracy of in-order vs. out-of-order cores

When examining the coverage of the constant predictions, we observe and interesting difference between the in-order and out-of-order cores. For in-order, the coverage of bzip2 is generally higher than that of sjeng (15.90% vs 1.18% by arithmetic mean of coverage) while the out-of-order shows the opposite trend (3.94% vs 17.48% by arithmetic mean of coverage). When looking at the absolute values of the number of loads processed by the two CPUs for sjeng, we see that there is a 50x difference between the two CPUs. The opposite is true for bzip2. This difference in total loads processed is the reason we see different coverage for the two benchmarks across the two CPUs and we think that the number of loads being processed is different because of the different scheduling policies of the two CPUs. While the actual number of loads would still be the same, either CPU could get caught in a loop where a branch is mispredicted and numerous loads are predicted and later squashed. We cannot verify this at this point since we have not collected statistics about load instructions that we classified as constant and later squashed in the pipeline.

### V. FUTURE WORK

There are several improvements or variations we believe would be useful to explore in future works. One improvement we believe is possible is to optimize the timing of the value forwarding, making constant results available earlier in the pipeline. In our current implementation, the instructions are issued and executed as normal even though the value is already available in the decode stage. This would require more extensive modifications to the execution pipeline to allow the *Predictable* and *Constant* to appear as if they have executed normally while structures like the Load Store Queue. In fact, this was our original approach but were unable to fully complete this approach after significant programming and debugging effort. We also note that this optimization is limited by the availability of the virtual address, which is not computed until the execute phase.

A feature we would like to implement is the utilization of the *Predictable* classification. This would allow the pipeline to speculatively execute instructions after a load with a predicted value. We expect that this approach would lead to a more significant improvement in IPC, as instructions that previously stalled waiting for a load could proceed speculatively. Implementing this functionality was part of our original plan and we made steps towards this implementation, but we unfortunately ran out of time to debug the squashing mechanism on value mispredictions.

We also note that our implementation only considers a single-CPU model, which may not be realistic for many modern workloads. Because we are effectively implementing a specialized cache, we believe future works should explore the implications of cache coherence requirements on load value prediction. For instance, the CVU must consider stores that happen in other processors to matching memory addresses and the cost of additional snooping hardware must be outweighed by reduced data accesses. Our implementation also considers associating values and classifications with specific threads, but we do not extensively study the effect of multiple threads on our implementation.

### VI. CONCLUSIONS

In this report, we have investigated the principles of value locality and value prediction. By using a load value prediction unit, we can see a notable drop in memory accesses for a few benchmarks which could be useful for devices that require a high degree of power efficiency. Overall, the addition of the load value prediction unit does not degrade the IPC and it helps reduce memory accesses for frequently used constants. However, we believe that by implementing a more aggressive value forwarding scheme and by executing instructions based on speculative values, we might be able to see more improvement in the IPC and this is something we would like to pursue in the future.

### VII. REFERENCES

[1] M. H. Lipasti, C. B. Wilkerson and J. P. Shen, "Value locality and load value prediction," in *Proceedings of the seventh international conference on Architectural support for programming languages and operating systems*, New York, NY, 1996.

[2] G. Sohi, *Unit 8: Storage Hierarchy I: Caches,* Madison, WI, 2021.

[3] J. L. Henning, "SPEC CPU2006 Benchmark Descriptions," *SIGARCH Computer Architecture News,* vol. 34, no. 4, pp. 1-17, 2006.

[4] K. Barker, T. Benson, D. Campbell, E. David, R. Gioiosa, A. Hoisie, D. Kerbyson, J. Manzano, A. Marquez, L. Song, N. Tallent and A. Tumeo, "PERFECT (Power Efficiency Revolution For," December 2013. [Online]. Available: http://hpc.pnnl.gov/projects/PERFECT/.. [Accessed 4 May 2021].

[5] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. H. J. Basu, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill and D. A. Wood, "The Gem5 Simulator," *ACM SIGARCH Computer Architecture News,* pp. 1-7, May 2011.

[6] A. Perais and A. Seznec, "Practical data value speculation for future high-end processors," in *2014 IEEE 20th International Symposium on High Performance Computer Architecture (HPCA)*, Orlando, 2014.

# Appendix

We provide the code written to model the Load value prediction unit as well as the diffs that show how the LVPU was integrated into the DerivO3 and Minor CPU models. Section I shows the code for the load value prediction unit, Section II shows the integration of the LVPU with the DerivO3 CPU model, and Section III shows the integration of the LVPU with the Minor CPU model.

## I. LOAD VALUE PREDICTION UNIT

### A. Load Value Prediction Table

```
1.   /***************************************************************************/
2.   // Author    : Prajyot Gupta
3.   // Department   : Grad Student @ Dept. of Electrical & Computer Engineering
4.   // Contact      : pgupta54@wisc.edu
5.   // Project      : ECE 752
6.   /***************************************************************************/
7.
8.   #ifndef __CPU_LVP_LOADVALUEPREDICTIONTABLE_HH__
9.   #define __CPU_LVP_LOADVALUEPREDICTIONTABLE_HH__
10.
11.  #include "base/types.hh"
12.  #include "cpu/static_inst.hh"
13.  #include "sim/sim_object.hh"
14.
15.  #include "enums.hh"
16.  #include "params/LoadValuePredictionTable.hh"
17.
18.  #include "arch/types.hh"
19.  #include "base/logging.hh"
20.  #include "config/the_isa.hh"
21.
22.  /** Creating a default Load Value Prediction Table entry
23.   *  which will have below attributes
24.   *  tag    : Specifies the Opcode of the Load instruction
25.   *  taget : Specifies the Load value associated with the tag
26.   *  valid : Specifies if the value loaded is valid
27.   */
28.
29.  class LoadValuePredictionTable : public SimObject
30.  {
31.    private:
32.      struct LVPTEntry
33.      {
34.          LVPTEntry()
35.              : tag(0), target(0), valid(false)
36.          {}
37.
38.          /** The entry's tag. */
39.          Addr tag;
40.
41.          /** The entry's target. */
42.          RegVal target;
43.
44.          /** The entry's thread id. */
45.          ThreadID tid;
46.
47.          /** Whether or not the entry is valid. */
48.          bool valid;
49.      };
50.
51.    public:
52.      /** Creates a LVPT with the given number of entries, number of bits per
53.       *  tag, and instruction offset amount.
```

```cpp
54.        *   @param numEntries Number of entries for the LVPT.
55.        *   @param tagBits Number of bits for each tag in the LVPT.
56.        *   @param instShiftAmt Offset amount for instructions to ignore alignment.
57.        */
58.       LoadValuePredictionTable(const LoadValuePredictionTableParams *params);
59.
60.       void reset();
61.
62.       /** Looks up an address in the LVPT. Must call valid() first on the address.
63.        *   @param inst_PC The address of the branch to look up.
64.        *   @param tid The thread id.
65.        *   @return Returns the predicated load value.
66.        */
67.       RegVal lookup(ThreadID tid, Addr instPC, bool *lvptResultValid);
68.
69.       /** Checks if the load entry is in the LVPT.
70.        *   @param inst_PC The address of the branch to look up.
71.        *   @param tid The thread id.
72.        *   @return Whether or not the branch exists in the LVPT.
73.        */
74.       bool valid(Addr instPC, ThreadID tid);
75.
76.       /** Updates the LVPT with the latest predicted Load Value.
77.        *   @param inst_PC The address of the branch being updated.
78.        *   @param target_PC The predicted target data.
79.        *   @param tid The thread id.
80.        */
81.
82.       void update(Addr instPC, const RegVal target, ThreadID tid);
83.
84.       /** Returns the index into the LVPT, based on the branch's PC.
85.        *   @param inst_PC The branch to look up.
86.        *   @return Returns the index into the LVPT.
87.        */
88.
89.       unsigned getIndex(Addr instPC, ThreadID tid);
90.
91.       /** Returns the tag bits of a given address.
92.        *   @param inst_PC The branch's address.
93.        *   @return Returns the tag bits.
94.        */
95.       inline Addr getTag(Addr instPC);
96.
97.   private:
98.
99.       /** The actual LVPT declaration */
100.          std::vector<LVPTEntry> LVPT;
101.
102.          /** The number of entries in the LVPT. */
103.          const unsigned numEntries;
104.
105.          /** Depth of data history kept in the LVPT*/
106.          const unsigned historyDepth;
107.
108.          //prajyotg :: No sure if below will be used
109.          /** The index mask. */
110.          unsigned idxMask;
111.
112.          /** The number of tag bits per entry. */
113.          unsigned tagBits;
114.
115.          /** The tag mask. */
116.          unsigned tagMask;
117.
118.          /** Number of bits to shift PC when calculating index. */
119.          unsigned instShiftAmt;
```

```
120.
121.        /** Number of bits to shift PC when calculating tag. */
122.        unsigned tagShiftAmt;
123.
124.        /** Log2 NumThreads used for hashing threadid */
125.        unsigned log2NumThreads;
126.    };
127.
128.    #endif // __CPU_LVP_LOADVALUEPREDICTIONTABLE_HH__
129.
```

```
1.  /*****************************************************************************/
2.  // Author    : Prajyot Gupta
3.  // Department    : Grad Student @ Dept. of Electrical & Computer Engineering
4.  // Contact       : pgupta54@wisc.edu
5.  // Project       : ECE 752
6.  /*****************************************************************************/
7.
8.  #include "cpu/lvp/load_value_prediction_table.hh"
9.
10. #include "base/intmath.hh"
11. #include "base/trace.hh"
12. #include "debug/LVPT.hh"
13.
14. LoadValuePredictionTable::LoadValuePredictionTable(const LoadValuePredictionTableParams *params)
15.     : SimObject(params),
16.       numEntries(params->entries),
17.       historyDepth(params->historyDepth),
18.       idxMask(numEntries - 1),
19.       instShiftAmt(0)
20.
21. {
22.     DPRINTF(LVPT, "LVPT: Creating LVPT object.\n");
23.
24.     if (!isPowerOf2(numEntries)) {
25.         fatal("LVPT entries is not a power of 2!");
26.     }
27.
28.     LVPT.resize(numEntries);
29.
30.     DPRINTF(LVPT, "LVPT: Doing an initial reset \n");
31.     for (unsigned i = 0; i < numEntries; ++i) {
32.         LVPT[i].valid = false;
33.     }
34.
35.     idxMask = numEntries - 1;
36.     tagMask = (1 << tagBits) - 1;
37.     tagShiftAmt = instShiftAmt + floorLog2(numEntries);
38. }
39.
40. /* Reset API */
41. void
42. LoadValuePredictionTable::reset()
43. {
44.     DPRINTF(LVPT, "LVPT : Calling the Reset API \n");
45.     for (unsigned i = 0; i < numEntries; ++i) {
46.         LVPT[i].valid = false;
47.     }
48. }
49.
50. /* APIs to get index and tag*/
51. unsigned
52. LoadValuePredictionTable::getIndex(Addr instPC, ThreadID tid)
53. {
54.     // Need to shift PC over by the word offset.
```

```
55.        // Math: ((instPC >> instShiftAmt)^(tid<<(tagShiftAmt-instShiftAmt-log2NumThreads)))&idxMask;
56.        DPRINTF(LVPT, "LVPT : Getting Index \n");
57.        return ((instPC >> instShiftAmt)
58.                ^ (tid << (tagShiftAmt - instShiftAmt - log2NumThreads)))
59.                & idxMask;
60. }
61.
62. inline
63. Addr
64. LoadValuePredictionTable::getTag(Addr instPC)
65. {
66.        DPRINTF(LVPT, "LVPT : Getting Tag \n");
67.        return (instPC >> tagShiftAmt) & tagMask;
68. }
69.
70. /** Checks if the load entry is in the LVPT.i **/
71. bool
72. LoadValuePredictionTable::valid(Addr instPC, ThreadID tid)
73. {
74.        unsigned LVPT_idx = getIndex(instPC, tid);
75.
76.
77.        // Making sure index doesn't go out of bounds
78.        assert(LVPT_idx < numEntries);
79.
80.        // Check if: (a) LVPT entry is valid
81.        // (b) index matches
82.        // (c) tag matches
83.        if (LVPT[LVPT_idx].valid
84.            && LVPT[LVPT_idx].tid == tid) {
85.            return true;
86.        } else {
87.            DPRINTF(LVPT, "LVPT : Checking if LVPT entry is valid \n");
88.            return false;
89.        }
90. }
91.
92. // data = 0 represent invalid entry.
93. RegVal
94. LoadValuePredictionTable::lookup(ThreadID tid, Addr instPC, bool *lvptResultValid)
95. {
96.        unsigned LVPT_idx = getIndex(instPC, tid);
97.
98.        assert(LVPT_idx < numEntries);
99.
100.       if (valid(instPC, tid)) {
101.           DPRINTF(LVPT, "Found valid entry for tid: %d at pc %#x : %d \n",
102.               tid, instPC, LVPT[LVPT_idx].target);
103.           *lvptResultValid = true;
104.           return LVPT[LVPT_idx].target;
105.       } else {
106.           DPRINTF(LVPT, "Did not find valid entry for tid: %d at address %#x \n",
107.               tid, instPC);
108.           return 0;
109.       }
110.   }
111.
112.   void
113.   //prajyotg :: updated :: LoadValuePredictionTable::update(Addr instPC, const TheISA::PCState &target,
    ThreadID tid)
114.   LoadValuePredictionTable::update(Addr instPC, const RegVal target, ThreadID tid)
115.   {
116.        unsigned LVPT_idx = getIndex(instPC, tid);
117.        DPRINTF(LVPT, "LVPT : Updating the value in the LVPT at index %ld \n", LVPT_idx);
118.
119.        assert(LVPT_idx < numEntries);
```

```
120.
121.        LVPT[LVPT_idx].tid = tid;
122.        LVPT[LVPT_idx].valid = true;
123.        LVPT[LVPT_idx].target = target;
124.        LVPT[LVPT_idx].tag = getTag(instPC);
125.    }
126.
127.    LoadValuePredictionTable*
128.    LoadValuePredictionTableParams::create()
129.    {
130.        return new LoadValuePredictionTable(this);
131.    }
132.
```

## B. Load Classification Table

```
1.  /*
2.   * Robert Viramontes
3.   * Created March 25, 2021
4.   * Based on the 2bit_local branch predictor provided with the gem5 source.
5.   *
6.   *
7.   */
8.
9.  #ifndef __CPU_LVP_LOADCLASSIFICATIONTABLE_HH__
10. #define __CPU_LVP_LOADCLASSIFICATIONTABLE_HH__
11.
12. #include <vector>
13.
14. #include "base/sat_counter.hh"
15. #include "base/types.hh"
16. #include "cpu/static_inst.hh"
17. #include "sim/sim_object.hh"
18.
19. #include "enums.hh"
20. #include "params/LoadClassificationTable.hh"
21.
22. /**
23.  * Implements a local predictor that uses the PC to index into a table of
24.  * counters.  Note that any time a pointer to the bp_history is given, it
25.  * should be NULL using this predictor because it does not have any branch
26.  * predictor state that needs to be recorded or updated; the update can be
27.  * determined solely by the branch being taken or not taken.
28.  */
29. class LoadClassificationTable : public SimObject
30. {
31.   public:
32.     /**
33.      * Build load classfication table.
34.      */
35.     LoadClassificationTable(const LoadClassificationTableParams *params);
36.
37.     /**
38.      * Looks up the given instruction address in the LCT and returns
39.      * a LctResult that indicates unpredictable, predictable, or constant.
40.      * @param inst_addr The address of the instruction to look up.
41.      * @param bp_history Pointer to any bp history state.
42.      * @return Whether or not the branch is taken.
43.      */
44.     LVPType lookup(ThreadID tid, Addr inst_addr);
45.
46.     /**
47.      * Updates the branch predictor with the actual result of a branch.
48.      * @param tid ThreadID of the thread to predict.
49.      * @param inst_addr The address of the instruction to update.
```

```
50.      * @param prediction The prediction that was associated with this instruction on decode.
51.      * @param prediction_correct Whether or not the load prediction was correct.
52.      * @return Updated LVPType of the location
53.      */
54.     LVPType update(ThreadID tid, Addr inst_addr, LVPType prediction, bool prediction_correct);
55.
56.     // void squash(ThreadID tid, void *bp_history)
57.     // { assert(bp_history == NULL); }
58.
59.     /**
60.      * @brief Resets all saturating counters in LCT.
61.      */
62.     void reset();
63.
64.   private:
65.     /**
66.      *  Returns the unpredictable/predictable/constant prediction given
67.      *  the value of the counter.
68.      *  @param count The value of the counter.
69.      *  @return The prediction based on the counter value.
70.      */
71.     inline LVPType getPrediction(uint8_t &count);
72.
73.     /** Calculates the local index based on the PC. */
74.     inline unsigned getLocalIndex(Addr &PC);
75.
76.     /** Resets the ctr at the specified index. */
77.     inline void resetCtr(unsigned local_predictor_idx);
78.
79.     /** Size of the local predictor. */
80.     const unsigned localPredictorSize;
81.
82.     /** Number of bits of the local predictor's counters. */
83.     const unsigned localCtrBits;
84.
85.     /** Number of sets. */
86.     const unsigned localPredictorSets;
87.
88.     /** Array of counters that make up the local predictor. */
89.     std::vector<SatCounter> localCtrs;
90.     std::vector<ThreadID> localCtrThreads;
91.
92.     /** Mask to get index bits. */
93.     const unsigned indexMask;
94.
95.     /** Number of bits to shift instructions by for predictor addresses. */
96.     const unsigned instShiftAmt;
97.
98.     const bool invalidateConstToZero;
99. };
100.
101.  #endif // __CPU_LVP_LOADCLASSIFICATIONTABLE_HH__
102.
```

```
1.   /*
2.    * Copyright (c) 2004-2006 The Regents of The University of Michigan
3.    * All rights reserved.
4.    *
5.    * Redistribution and use in source and binary forms, with or without
6.    * modification, are permitted provided that the following conditions are
7.    * met: redistributions of source code must retain the above copyright
8.    * notice, this list of conditions and the following disclaimer;
9.    * redistributions in binary form must reproduce the above copyright
10.   * notice, this list of conditions and the following disclaimer in the
11.   * documentation and/or other materials provided with the distribution;
```

```
12.    * neither the name of the copyright holders nor the names of its
13.    * contributors may be used to endorse or promote products derived from
14.    * this software without specific prior written permission.
15.    *
16.    * THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS
17.    * "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT
18.    * LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR
19.    * A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT
20.    * OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL,
21.    * SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT
22.    * LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE,
23.    * DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY
24.    * THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT
25.    * (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE
26.    * OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
27.    */
28.
29.  #include "cpu/lvp/load_classification_table.hh"
30.
31.  #include "base/intmath.hh"
32.  #include "base/logging.hh"
33.  #include "base/trace.hh"
34.  #include "debug/LCT.hh"
35.
36.  LoadClassificationTable::LoadClassificationTable(const LoadClassificationTableParams *params)
37.      : SimObject(params),
38.        localPredictorSize(params->localPredictorSize),
39.        localCtrBits(params->localCtrBits),
40.        localPredictorSets(localPredictorSize / localCtrBits),
41.        localCtrs(localPredictorSets, SatCounter(localCtrBits, 0)),
42.        localCtrThreads(localPredictorSets, 0),
43.        indexMask(localPredictorSets - 1),
44.        invalidateConstToZero(params->invalidateConstToZero),
45.        instShiftAmt(0) // TODO what is correct???
46.  {
47.      if (!isPowerOf2(localPredictorSize)) {
48.          fatal("Invalid LCT local predictor size!\n");
49.      }
50.
51.      if (!isPowerOf2(localPredictorSets)) {
52.          fatal("Invalid number of LCT local predictor sets! Check localCtrBits.\n");
53.      }
54.
55.      DPRINTF(LCT, "LCT index mask: %#x\n", indexMask);
56.
57.      DPRINTF(LCT, "LCT size: %i\n",
58.              localPredictorSize);
59.
60.      DPRINTF(LCT, "LCT counter bits: %i\n", localCtrBits);
61.
62.      DPRINTF(LCT, "instruction shift amount: %i\n",
63.              instShiftAmt);
64.  }
65.
66.  LVPType
67.  LoadClassificationTable::lookup(ThreadID tid, Addr inst_addr)
68.  {
69.      unsigned local_predictor_idx = getLocalIndex(inst_addr);
70.
71.      if (tid != localCtrThreads[local_predictor_idx])
72.      {
73.          return LVP_STRONG_UNPREDICTABLE;
74.      }
75.
76.      DPRINTF(LCT, "Thread ID: %d", tid);
77.
```

```
78.     uint8_t counter_val = localCtrs[local_predictor_idx];
79.
80.     DPRINTF(LCT, "During lookup index %#x had value %i\n",
81.             local_predictor_idx, (int)counter_val);
82.
83.     return getPrediction(counter_val);
84. }
85.
86. LVPType
87. LoadClassificationTable::update(ThreadID tid, Addr inst_addr, LVPType prediction, bool
    prediction_correct)
88. {
89.     unsigned local_predictor_idx;
90.
91.     // No state to restore, and we do not update on the wrong
92.     // path.
93.     // if (squashed) {
94.     //     return;
95.     // }
96.
97.     // Update the local predictor.
98.     local_predictor_idx = getLocalIndex(inst_addr);
99.
100.      if (tid == localCtrThreads[local_predictor_idx])
101.      {
102.          if (prediction_correct) {
103.              DPRINTF(LCT, "Load classification updated as correct.\n");
104.              localCtrs[local_predictor_idx]++;
105.          } else {
106.              DPRINTF(LCT, "Load classification updated as incorrect.\n");
107.              if (prediction == LVP_CONSTANT && invalidateConstToZero) {
108.                  resetCtr(local_predictor_idx);
109.              } else {
110.                  localCtrs[local_predictor_idx]--;
111.              }
112.          }
113.      }
114.      else
115.      {
116.          // Destructive interference with a different thread, reset this index and update the thread
117.          localCtrThreads[local_predictor_idx] = tid;
118.          resetCtr(local_predictor_idx);
119.      }
120.
121.      uint8_t counter_val = localCtrs[local_predictor_idx];
122.      return getPrediction(counter_val);
123. }
124.
125. void
126. LoadClassificationTable::reset()
127. {
128.     auto zeroCounter = new SatCounter(localCtrBits, 0);
129.     for (int c = 0; c < localPredictorSets; c++)
130.     {
131.         localCtrs[c] = *zeroCounter;
132.     }
133.
134.     delete zeroCounter;
135. }
136.
137. void
138. LoadClassificationTable::resetCtr(unsigned local_predictor_idx)
139. {
140.     while (localCtrs[local_predictor_idx] != LVP_STRONG_UNPREDICTABLE) {
141.         localCtrs[local_predictor_idx]--;
142.     }
```

```
143.  }
144.
145.  inline
146.  LVPType
147.  LoadClassificationTable::getPrediction(uint8_t &count)
148.  {
149.      // If MSB is 0, value is unpredictable
150.      // If counter is saturated, value is constant
151.      // Otherwise, the value is predictable
152.      return (count == 0) ? LVP_STRONG_UNPREDICTABLE
153.          : (count >> (localCtrBits - 1)) == 0 ? LVP_WEAK_UNPREDICTABLE
154.          : count == (power(2, localCtrBits) -1) ? LVP_CONSTANT
155.          : LVP_PREDICTABLE;
156.  }
157.
158.  inline
159.  unsigned
160.  LoadClassificationTable::getLocalIndex(Addr &inst_addr)
161.  {
162.      return (inst_addr >> instShiftAmt) & indexMask;
163.  }
164.
165.  LoadClassificationTable*
166.  LoadClassificationTableParams::create()
167.  {
168.      return new LoadClassificationTable(this);
169.  }
170.
```

## C. Constant Verification Unit

```
1.  /**
2.   * @defgroup   CONSTANT_VERIFICATION_UNIT constant verification unit
3.   *
4.   * @brief      This file the implements class of the constant verification unit
5.   *                         for a load value predictor.
6.   *                         The constant verification unit has the following functions:
7.   *                         - Keep track of load addresses which have been deemed "constant"
8.   *                         by the LCT
9.   *                         - Monitor all committed stores and invalidate a constant load
10.  *                         address if a store that is being committed has the same address
11.  *                         as the constant load
12.  *                         - Compare a load value for address that had been deemed
13.  *                         "predictable" by the LCT with the actual load value retrieved
14.  *                         from memory
15.  *                         - Provide feedback to the LCT for all predictable loads
16.  *                         - Possible implementation of this class could also include a
17.  *                         mechanism to pass load values from the commit stage back to the
18.  *                         LVPT (Since this will be needed if an instruction is to be
19.  *                         promoted from unpredictable to predictable and so on.)
20.  *
21.  * @author     Satvik
22.  * @date       2021
23.  */
24.
25. #ifndef __CPU_LVP_CONSTANT_VERIFICATION_UNIT__
26. #define __CPU_LVP_CONSTANT_VERIFICATION_UNIT__
27.
28. #include "sim/sim_object.hh"
29. #include "base/types.hh"
30. #include "params/ConstantVerificationUnit.hh"
31. #include "base/statistics.hh"
32.
33. #include <vector>
34. #include <map>
```

```cpp
35. #include <iterator>
36. #include <algorithm>
37. #include <chrono>
38.
39. // Replacement Policies
40. #define RP_FIFO 1        // CAM is implemented as a queue
41. #define RP_LRU  2                // The least recently used entry is kicked out
42. #define RP_NLRU 3                // The not least recently used entry is kicked out
43. #define RP_MRU  4        // The most recently used entry is kicked out
44. #define RP_NMRU 5        // The not most recently used entry is kicked out
45.
46. #define DURATION std::chrono::duration_cast<std::chrono::duration<double, std::milli>>
47. typedef std::chrono::high_resolution_clock::time_point TimePoint;
48.
49. /**
50.  * @brief      The CVU stores the load address and the index of the LVPT in a
51.  *             concatenated form in its CAM. This is searched by loads that are
52.  *             predicted as constant and by store addresses. Thus, the search of
53.  *             the CAM will be different for each case. The search for a load
54.  *             predicted as "constant" will be done on the entire entry in the
55.  *             CAM while a store address will only need to be compared with the
56.  *             load address stored in the CAM.
57.  */
58.
59. /**
60.  * @brief      Data stored in every entry of the CVU CAM
61.  */
62. struct CAMEntry {
63.    Addr lvpt_index;
64.    Addr pc;
65.    Addr load_address;
66.    TimePoint entry;
67.    TimePoint access;
68. };
69.
70. class ConstantVerificationUnit : public SimObject {
71. public:
72.    ConstantVerificationUnit(ConstantVerificationUnitParams *p);
73.
74.    ~ConstantVerificationUnit();
75.
76.    /**
77.     * @brief      A store address will be provided to the CVU which will then
78.     *                          search the CAM and invalidate every entry which matches the
79.     *                          input address. This will also generate the control signals
80.     *                          required for communicating with the LCT.
81.     *
82.     * @param[in]  address  The store address
83.     */
84.    void processStoreAddress(ThreadID tid, Addr address);
85.
86.    /**
87.     * @brief      Check if a load address classified as constant is present in
88.     *             the CVU CAM.
89.     *
90.     * @param[in]  pc         The PC for this load instruction
91.     * @param[in]  lvptIndex  The lvpt index corresponding to that load address
92.     * @param[in]  tid        Thread ID
93.     *
94.     * @return     True if the load address, LVPT index pair exist in the CAM
95.     *                          False otherwise
96.     */
97.    bool processLoadAddress(Addr load_address, Addr lvptIndex, ThreadID tid);
98.
99.    /**
100.             * @brief      Inserts info of a new constant load into the CVU CAM
```

```cpp
101.            *
102.            * @param[in]  pc         Instruction address
103.            * @param[in]  address    The load address
104.            * @param[in]  lvptIndex  The lvpt index
105.            * @param[in]  tid        The tid
106.            *
107.            * @return     True if the update is successful
108.            */
109.          bool updateConstLoad(Addr pc, Addr address, Addr lvptIndex, ThreadID tid);
110.
111.          /**
112.            * @brief      Replaces an entry in the CVU CAM with a new one according to
113.            *             a replacement policy
114.            *
115.            * @param[in]  new_entry  The new entry
116.            * @param[in]  tid        The tid
117.            */
118.          void replaceBlock(struct CAMEntry new_entry, ThreadID tid);
119.
120.          /**
121.            * @brief      Print stats
122.            */
123.          void regStats() override;
124.
125.  private:
126.          /**
127.            * The CVU Content Addressable Memory;
128.            * If a store address is found in this memory, the corresponding entry is
129.            * invalidated and the invalidation also triggers an update routine which
130.            * tells the LCT that this load address is no longer constant.
131.            * The vector stores an ordered pair (load address, LVPT index)
132.            */
133.          std::list<struct CAMEntry> _cvuCAM[64];
134.
135.          /**
136.            * Number of entries in the CVU CAM
137.            */
138.          uint32_t _numEntries;
139.
140.          /**
141.            * Number of loads marked "constant" which were incorrectly predicted.
142.            */
143.          Stats::Scalar _numConstantHits;
144.
145.          /**
146.            * Number of loads marked "constant" that were correctly predicted.
147.            */
148.          Stats::Scalar _numConstantMiss;
149.
150.          /**
151.            * Number of store addresses which hit in the CVU CAM
152.            */
153.          Stats::Scalar _numStoreHits;
154.
155.          /**
156.            * Number of store addresses which missed in the CVU CAM.
157.            */
158.          Stats::Scalar _numStoreMiss;
159.
160.          /**
161.            * Replacement policy of the CAM
162.            */
163.          uint8_t _replacementPolicy;
164.
165.          /**
166.            * Number of CAM blocks replaced.
```

```
167.            */
168.            Stats::Scalar _numReplacements;
169.
170.            /**
171.             * Total stores processed by the CVU
172.             */
173.            Stats::Formula _numStoreAccesses;
174.
175.            /**
176.             * Total Loads processed by the CVU
177.             */
178.            Stats::Formula _numLoadAccesses;
179. };
180.
181.
182. #endif
183.
```

```
1.  /**
2.   * @defgroup   CONSTANT_VERIFICATION_UNIT constant verification unit
3.   *
4.   * @brief      This file implements constant verification unit.
5.   *
6.   * @author     Satvik
7.   * @date       2021
8.   */
9.  #include "cpu/lvp/constant_verification_unit.hh"
10.
11. #include "base/intmath.hh"
12. #include "base/logging.hh"
13. #include "base/trace.hh"
14. #include "debug/CVU.hh"
15.
16. ConstantVerificationUnit::ConstantVerificationUnit(ConstantVerificationUnitParams *params) :
17.     SimObject(params), _numEntries(params->entries), _numConstantHits(0),
18.     _numConstantMiss(0), _numStoreHits(0), _numStoreMiss(0),
19.     _replacementPolicy(params->replacementPolicy), _numReplacements(0)
20. {}
21.
22. ConstantVerificationUnit::~ConstantVerificationUnit() {}
23.
24. void ConstantVerificationUnit::processStoreAddress(ThreadID tid,
25.
                Addr address) {
26.     DPRINTF(CVU, "[TID %d]: Store address: 0x%x being searched in CVU CAM\n", tid, address);
27.     // Only the load address needs to be compared with the store address
28.     bool found = false;
29.     auto itr = _cvuCAM[(uint16_t)tid%64].begin();
30.     while (itr != _cvuCAM[(uint16_t)tid%64].end()) {
31.             struct CAMEntry temp = *itr;
32.             if (temp.load_address == address) {
33.                     DPRINTF(CVU, "[TID %d]: Found store address: 0x%x in CVU CAM\n", tid, address);
34.                     itr = _cvuCAM[(uint16_t)tid%64].erase(itr);
35.                     ++_numStoreHits;
36.                     found = true;
37.                     if (itr == _cvuCAM[(uint16_t)tid%64].end()) break;
38.             }
39.             itr++;
40.     }
41.     if (!found) {
42.             DPRINTF(CVU, "[TID %d]: Address 0x%x not found in CVU CAM\n", tid, address);
43.             ++_numStoreMiss;
44.     }
45. }
46.
```

```
47. bool ConstantVerificationUnit::processLoadAddress(Addr load_address,
48.
            Addr lvptIndex,
49.
            ThreadID tid) {
50.   // Both load address and LVPT index have to be searched
51.   auto itr = _cvuCAM[(uint16_t)tid%64].begin();
52.   while (itr != _cvuCAM[(uint16_t)tid%64].end()) {
53.           struct CAMEntry temp = *itr;
54.           if (temp.load_address == load_address && temp.lvpt_index == lvptIndex) {
55.                   DPRINTF(CVU, "[TID %d] Load address: 0x%x matched in CVU CAM\n", tid,
    load_address);
56.                   ++_numConstantHits;
57.                   // Update the end time of this entry (for replacement purposes)
58.                   itr->access = std::chrono::high_resolution_clock::now();
59.                   return true;
60.           }
61.           itr++;
62.   }
63.   ++_numConstantMiss;
64.   DPRINTF(CVU, "[TID %d] Load address: 0x%x not found in CVU CAM\n", tid, load_address);
65.   return false;
66. }
67.
68. bool ConstantVerificationUnit::updateConstLoad(Addr pc, Addr address,
69.
    Addr lvptIndex, ThreadID tid) {
70.   DPRINTF(CVU, "[TID %d]: Adding load address: 0x%x with PC: 0x%x to CVU CAM\n", tid, address, pc);
71.   struct CAMEntry temp;
72.   temp.pc = pc;
73.   temp.lvpt_index = lvptIndex;
74.   temp.load_address = address;
75.   temp.entry = std::chrono::high_resolution_clock::now();
76.   temp.access = temp.entry;
77.   if(_cvuCAM[(uint16_t)tid%64].size() >= _numEntries) {
78.           // CAM is full, replace a block according to a set replacement policy
79.           DPRINTF(CVU, "[TID %d]: No space for load address: 0x%x in CVU CAM\n", tid, address);
80.           this->replaceBlock(temp, tid);
81.           return true;
82.   }
83.   _cvuCAM[(uint16_t)tid%64].push_back(temp);
84.   return true;
85. }
86.
87. void ConstantVerificationUnit::replaceBlock(struct CAMEntry new_entry, ThreadID tid) {
88.   _numReplacements++;
89.   if(_replacementPolicy == RP_FIFO) {
90.           // Kick the oldest entry from the CAM
91.           auto itr = _cvuCAM[(uint16_t)tid%64].begin();
92.           auto itr_replace = itr;
93.           TimePoint current_time = std::chrono::high_resolution_clock::now();
94.           double time_diff = DURATION(current_time - itr->entry).count();
95.           itr++;
96.           while(itr != _cvuCAM[(uint16_t)tid%64].end()) {
97.                   if(time_diff < DURATION(current_time - itr->entry).count()) {
98.                           itr_replace = itr;
99.                           time_diff = DURATION(current_time - itr->entry).count();
100.                          }
101.                          itr++;
102.                  }
103.                  DPRINTF(CVU, "Entry with PC: 0x%x being ejected from CVU CAM\n", itr_replace->pc);
104.                  _cvuCAM[(uint16_t)tid%64].erase(itr_replace);
105.                  _cvuCAM[(uint16_t)tid%64].push_back(new_entry);
106.                  return;
107.          }
108.          if(_replacementPolicy == RP_LRU) {
```

```cpp
109.                    // Kick the least recently used (compare the access)
110.                    auto itr = _cvuCAM[(uint16_t)tid%64].begin();
111.                    auto itr_replace = itr;
112.                    TimePoint current_time = std::chrono::high_resolution_clock::now();
113.                    double time_diff = DURATION(current_time - itr->access).count();
114.                    itr++;
115.                    while(itr != _cvuCAM[(uint16_t)tid%64].end()) {
116.                            if(time_diff < DURATION(current_time - itr->access).count()) {
117.                                    itr_replace = itr;
118.                                    time_diff = DURATION(current_time - itr->access).count();
119.                            }
120.                            itr++;
121.                    }
122.                    DPRINTF(CVU, "Entry with PC: 0x%x being ejected from CVU CAM\n", itr_replace->pc);
123.                    _cvuCAM[(uint16_t)tid%64].erase(itr_replace);
124.                    _cvuCAM[(uint16_t)tid%64].push_back(new_entry);
125.                    return;
126.            }
127.            if(_replacementPolicy == RP_NLRU) {
128.                    // Kick the not least recently used (compare the access)
129.                    auto itr = _cvuCAM[(uint16_t)tid%64].begin();
130.                    auto itr_replace = itr;
131.                    TimePoint current_time = std::chrono::high_resolution_clock::now();
132.                    double time_diff = DURATION(current_time - itr->access).count();
133.                    itr++;
134.                    while(itr != _cvuCAM[(uint16_t)tid%64].end()) {
135.                            if(time_diff < DURATION(current_time - itr->access).count()) {
136.                                    itr_replace = itr;
137.                                    time_diff = DURATION(current_time - itr->access).count();
138.                            }
139.                            itr++;
140.                    }
141.                    itr = _cvuCAM[(uint16_t)tid%64].begin();
142.                    while(itr == itr_replace && itr != _cvuCAM[(uint16_t)tid%64].end()) {
143.                            itr++;
144.                    }
145.                    itr_replace = itr;
146.                    DPRINTF(CVU, "Entry with PC: 0x%x being ejected from CVU CAM\n", itr_replace->pc);
147.                    _cvuCAM[(uint16_t)tid%64].erase(itr_replace);
148.                    _cvuCAM[(uint16_t)tid%64].push_back(new_entry);
149.                    return;
150.            }
151.            if(_replacementPolicy == RP_MRU) {
152.                    // Kick the least recently used (compare the access)
153.                    auto itr = _cvuCAM[(uint16_t)tid%64].begin();
154.                    auto itr_replace = itr;
155.                    TimePoint current_time = std::chrono::high_resolution_clock::now();
156.                    double time_diff = DURATION(current_time - itr->access).count();
157.                    itr++;
158.                    while(itr != _cvuCAM[(uint16_t)tid%64].end()) {
159.                            if(time_diff > DURATION(current_time - itr->access).count()) {
160.                                    itr_replace = itr;
161.                                    time_diff = DURATION(current_time - itr->access).count();
162.                            }
163.                            itr++;
164.                    }
165.                    DPRINTF(CVU, "Entry with PC: 0x%x being ejected from CVU CAM\n", itr_replace->pc);
166.                    _cvuCAM[(uint16_t)tid%64].erase(itr_replace);
167.                    _cvuCAM[(uint16_t)tid%64].push_back(new_entry);
168.                    return;
169.            }
170.            if(_replacementPolicy == RP_NMRU) {
171.                    // Kick the not least recently used (compare the access)
172.                    auto itr = _cvuCAM[(uint16_t)tid%64].begin();
173.                    auto itr_replace = itr;
174.                    TimePoint current_time = std::chrono::high_resolution_clock::now();
```

```
175.                      double time_diff = DURATION(current_time - itr->access).count();
176.                      itr++;
177.                      while(itr != _cvuCAM[(uint16_t)tid%64].end()) {
178.                              if(time_diff > DURATION(current_time - itr->access).count()) {
179.                                      itr_replace = itr;
180.                                      time_diff = DURATION(current_time - itr->access).count();
181.                              }
182.                              itr++;
183.                      }
184.                      itr = _cvuCAM[(uint16_t)tid%64].begin();
185.                      while(itr == itr_replace && itr != _cvuCAM[(uint16_t)tid%64].end()) {
186.                              itr++;
187.                      }
188.                      itr_replace = itr;
189.                      DPRINTF(CVU, "Entry with PC: 0x%x being ejected from CVU CAM\n", itr_replace->pc);
190.                      _cvuCAM[(uint16_t)tid%64].erase(itr_replace);
191.                      _cvuCAM[(uint16_t)tid%64].push_back(new_entry);
192.                      return;
193.              }
194.  }
195.
196.  void ConstantVerificationUnit::regStats() {
197.          SimObject::regStats();
198.
199.          _numConstantHits.name(name() + ".constLoadHits")
200.                                        .desc("Number of loads marked constant that hit in
      the CVU CAM");
201.          _numConstantMiss.name(name() + ".constLoadMisses")
202.                                        .desc("Number of loads marked constant that missed in
      the CVU CAM");
203.          _numStoreHits.name(name() + ".numStoreHits")
204.                                        .desc("Number of stores that hit in the CVU CAM");
205.          _numStoreMiss.name(name() + ".numStoreMisses")
206.                                        .desc("Number of stores that missed in the CVU CAM");
207.          _numReplacements.name(name() + ".numCAMReplacements")
208.                                        .desc("Number of CVU CAM entries replaced");
209.          _numStoreAccesses.name(name() + ".numStoreAccesses")
210.                                        .desc("Number of store instructions which accessed
      the CVU");
211.          _numLoadAccesses.name(name() + ".numLoadAccesses")
212.                                        .desc("Number of loads marked constant which accessed
      the CVU");
213.
214.          _numLoadAccesses = _numConstantMiss + _numConstantHits;
215.          _numStoreAccesses = _numStoreHits + _numStoreMiss;
216.  }
217.
218.  ConstantVerificationUnit*
219.  ConstantVerificationUnitParams::create()
220.  {
221.      return new ConstantVerificationUnit(this);
222.  }
223.
```

```
1.  diff --git a/gem5/src/cpu/o3/O3CPU.py b/gem5/src/cpu/o3/O3CPU.py
2.  index 51d9121..864d54e 100644
3.  --- a/gem5/src/cpu/o3/O3CPU.py
4.  +++ b/gem5/src/cpu/o3/O3CPU.py
5.  @@ -46,6 +46,7 @@ from m5.objects.BaseCPU import BaseCPU
6.    from m5.objects.FUPool import *
7.    from m5.objects.O3Checker import O3Checker
8.    from m5.objects.BranchPredictor import *
9.  +from m5.objects.LoadValuePredictionUnit import *
10.
11.  class FetchPolicy(ScopedEnum):
12.      vals = [ 'SingleThread', 'RoundRobin', 'Branch', 'IQCount', 'LSQCount' ]
13. @@ -127,7 +128,7 @@ class DerivO3CPU(BaseCPU):
14.      LQEntries = Param.Unsigned(32, "Number of load queue entries")
15.      SQEntries = Param.Unsigned(32, "Number of store queue entries")
16.      LSQDepCheckShift = Param.Unsigned(4, "Number of places to shift addr before check")
17. -    LSQCheckLoads = Param.Bool(True,
18. +    LSQCheckLoads = Param.Bool(False,
19.          "Should dependency violations be checked for loads & stores or just stores")
20.      store_set_clear_period = Param.Unsigned(250000,
21.              "Number of load/store insts before the dep predictor should be invalidated")
22. @@ -174,6 +175,12 @@ class DerivO3CPU(BaseCPU):
23.      branchPred = Param.BranchPredictor(TournamentBP(numThreads =
24.                                          Parent.numThreads),
25.                                  "Branch Predictor")
26. +    ##
27. +    ## SATVIK:
28. +    ## LVP object that comes from somewhere
29. +    ##
30. +    loadValPred = Param.LoadValuePredictionUnit(LoadValuePredictionUnit(),
31. +                                  "Load value predictor")
32.      needsTSO = Param.Bool(buildEnv['TARGET_ISA'] == 'x86',
33.                      "Enable TSO Memory model")
34.
35. diff --git a/gem5/src/cpu/o3/commit_impl.hh b/gem5/src/cpu/o3/commit_impl.hh
36. index 75d065f..61befa4 100644
37. --- a/gem5/src/cpu/o3/commit_impl.hh
38. +++ b/gem5/src/cpu/o3/commit_impl.hh
39. @@ -706,6 +706,9 @@ DefaultCommit<Impl>::tick()
40.              DPRINTF(Commit,"[tid:%i] Instruction [sn:%llu] PC %s is head of"
41.                  " ROB and ready to commit\n",
42.                  tid, inst->seqNum, inst->pcState());
43. +            if(inst->isLoad()) {
44. +                DPRINTF(Commit, "Inst ready to commit is a load:[sn:%llu]\n", inst->seqNum);
45. +            }
46.
47.          } else if (!rob->isEmpty(tid)) {
48.              const DynInstPtr &inst = rob->readHeadInst(tid);
49. @@ -715,6 +718,18 @@ DefaultCommit<Impl>::tick()
50.              DPRINTF(Commit,"[tid:%i] Can't commit, Instruction [sn:%llu] PC "
51.                      "%s is head of ROB and not ready\n",
52.                      tid, inst->seqNum, inst->pcState());
53. +            for(int i = 0; i < inst->numSrcRegs(); i++) {
54. +                if(!inst->isReadySrcRegIdx(i)) {
55. +                    auto ptr = inst->renamedSrcRegIdx(i);
56. +                    DPRINTF(Commit, "Src reg %d for inst[%llu]: 0x%x not ready\n", ptr->index(),inst->seqNum, inst->instAddr());
57. +                }
58. +            }
59. +            if(inst->isConstPredictionCorrect() && !inst->readyToCommit()) {
60. +                DPRINTF(Commit, "Head instruction[%d] wasn't ready in ROB and was a const load\n", inst->seqNum);
61. +                //inst->setExecuted();
62. +                //inst->setResultReady();
```

```
63. +                    //inst->setCanCommit();
64. +                }
65.              }
66.
67.              DPRINTF(Commit, "[tid:%i] ROB has %d insts & %d free entries.\n",
68. diff --git a/gem5/src/cpu/o3/cpu.cc b/gem5/src/cpu/o3/cpu.cc
69. index 01938f1..c28c25f 100644
70. --- a/gem5/src/cpu/o3/cpu.cc
71. +++ b/gem5/src/cpu/o3/cpu.cc
72. @@ -126,6 +126,7 @@ FullO3CPU<Impl>::FullO3CPU(DerivO3CPUParams *params)
73.                  params->activity),
74.
75.          globalSeqNum(1),
76. +        lvp(params->loadValPred),
77.          system(params->system),
78.          lastRunningCycle(curCycle())
79.  {
80. diff --git a/gem5/src/cpu/o3/cpu.hh b/gem5/src/cpu/o3/cpu.hh
81. index 0447275..cbf61ce 100644
82. --- a/gem5/src/cpu/o3/cpu.hh
83. +++ b/gem5/src/cpu/o3/cpu.hh
84. @@ -64,6 +64,8 @@
85.  #include "params/DerivO3CPU.hh"
86.  #include "sim/process.hh"
87.
88. +#include "cpu/lvp/load_value_prediction_unit.hh"
89. +
90.  template <class>
91.  class Checker;
92.  class ThreadContext;
93. @@ -426,6 +428,26 @@ class FullO3CPU : public BaseO3CPU
94.
95.      RegVal readArchFloatReg(int reg_idx, ThreadID tid);
96.
97. +    void tagLVPDestReg(PhysRegIdPtr reg) {
98. +        regFile.insertPredictedLoadRegister(reg);
99. +    }
100.  +
101.  +    bool checkLVPTag(PhysRegIdPtr src_reg) {
102.  +        return regFile.checkPredictedLoadRegister(src_reg);
103.  +    }
104.  +
105.  +    void removeLVPTag(PhysRegIdPtr reg) {
106.  +        regFile.clearPredictedLoadRegister(reg);
107.  +    }
108.  +
109.  +    void addToMispredictList(PhysRegIdPtr reg) {
110.  +        regFile.addToMispredictList(reg);
111.  +    }
112.  +
113.  +    bool searchMispredictList(PhysRegIdPtr src_reg) {
114.  +        return regFile.searchMispredictList(src_reg);
115.  +    }
116.  +
117.          const VecRegContainer& readArchVecReg(int reg_idx, ThreadID tid) const;
118.          /** Read architectural vector register for modification. */
119.          VecRegContainer& getWritableArchVecReg(int reg_idx, ThreadID tid);
120.  @@ -690,6 +712,16 @@ class FullO3CPU : public BaseO3CPU
121.          */
122.          Checker<Impl> *checker;
123.
124.  +    /**
125.  +     * SATVIK:
126.  +     * Create an object for the LVP here?
127.  +     * Makes sense since the LVP is a module that is shared across multiple
128.  +     * stages of the CPU, unlike the bpred which belongs in fetch
```

```
129.  +       * Any LVP API can then be called using inst->cpu->LVPAPI() since the
130.  +       * cpu object is passed to the dyn_inst class.
131.  +       */
132.  +     LoadValuePredictionUnit *lvp;
133.  +
134.        /** Pointer to the system. */
135.        System *system;
136.
137.  diff --git a/gem5/src/cpu/o3/dyn_inst.hh b/gem5/src/cpu/o3/dyn_inst.hh
138.  index 5f2a588..8064f79 100644
139.  --- a/gem5/src/cpu/o3/dyn_inst.hh
140.  +++ b/gem5/src/cpu/o3/dyn_inst.hh
141.  @@ -50,6 +50,7 @@
142.   #include "cpu/base_dyn_inst.hh"
143.   #include "cpu/inst_seq.hh"
144.   #include "cpu/reg_class.hh"
145.  +#include "cpu/lvp/load_value_prediction_unit.hh"
146.
147.   class Packet;
148.
149.  @@ -116,6 +117,16 @@ class BaseO3DynInst : public BaseDynInst<Impl>
150.        /** Number of destination misc. registers. */
151.        uint8_t _numDestMiscRegs;
152.
153.  +     /**
154.  +      * The classification of this instruction as per the LCT (only for loads)
155.  +      */
156.  +     LVPType _classification;
157.  +
158.  +     RegVal _predictedVal;
159.  +
160.  +     bool _predictionCorrect;
161.  +
162.  +     bool _specExecOnLoad;
163.
164.      public:
165.   #if TRACING_ON
166.  @@ -423,6 +434,123 @@ class BaseO3DynInst : public BaseDynInst<Impl>
167.            this->cpu->setCCReg(this->_destRegIdx[idx], val);
168.            BaseDynInst<Impl>::setCCRegOperand(si, idx, val);
169.        }
170.  +
171.  +     std::pair<LVPType, RegVal>
172.  +     predictLoad(ThreadID tid) {
173.  +         std::pair<LVPType, RegVal> temp = this->cpu->lvp->predictLoad(tid,
174.  +                                                     this->instAddr());
175.  +         _classification = temp.first;
176.  +         _predictedVal = temp.second;
177.  +         return temp;
178.  +     }
179.  +
180.  +     bool
181.  +     verifyConstLoad(ThreadID tid) {
182.  +         Addr lvpt_index = this->cpu->lvp->lookupLVPTIndex(tid,
183.  +                                                     this->instAddr());
184.  +         _predictionCorrect = this->cpu->lvp->processLoadAddress(tid,
185.  +                                         this->instAddr(), this->effAddr, lvpt_index);
186.  +         return _predictionCorrect;
187.  +     }
188.  +
189.  +     bool
190.  +     isConstPredictionCorrect() {
191.  +         return _predictionCorrect;
192.  +     }
193.  +
194.  +     RegVal
```

```
195.  +      getPredictedValue() {
196.  +          return _predictedVal;
197.  +      }
198.  +
199.  +      bool isConstLoad() {
200.  +          return _classification == LVP_CONSTANT;
201.  +      }
202.  +
203.  +      void
204.  +      tagLVPDestReg(int idx) {
205.  +          this->cpu->tagLVPDestReg(this->_destRegIdx[idx]);
206.  +      }
207.  +
208.  +      bool
209.  +      checkLVPTag(int idx) {
210.  +          return this->cpu->checkLVPTag(this->_srcRegIdx[idx]);
211.  +      }
212.  +
213.  +      PhysRegIdPtr
214.  +      getSrcRegPtr(int idx) {
215.  +          return this->_srcRegIdx[idx];
216.  +      }
217.  +
218.  +      PhysRegIdPtr
219.  +      getDestRegPtr(int idx) {
220.  +          return this->_destRegIdx[idx];
221.  +      }
222.  +
223.  +      void
224.  +      removeLVPTag(int idx) {
225.  +          this->cpu->removeLVPTag(this->_destRegIdx[idx]);
226.  +      }
227.  +
228.  +      void
229.  +      addToMispredictList(int idx) {
230.  +          this->cpu->addToMispredictList(this->_destRegIdx[idx]);
231.  +      }
232.  +
233.  +      bool
234.  +      searchMispredictList(int idx) {
235.  +          return this->cpu->searchMispredictList(this->_srcRegIdx[idx]);
236.  +      }
237.  +
238.  +      bool
239.  +      verifyPrediction(int idx) {
240.  +          RegVal temp = 0;
241.  +          auto ptr = this->_destRegIdx[idx];
242.  +          if(ptr->isIntPhysReg()) {
243.  +              temp = this->cpu->readIntReg(ptr);
244.  +          }
245.  +          else if(ptr->isFloatPhysReg()) {
246.  +              temp = this->cpu->readFloatReg(ptr);
247.  +          }
248.  +          else {
249.  +              return true;
250.  +          }
251.  +          this->removeLVPTag(idx);
252.  +          if(!this->effAddrValid()) panic("Virtual address not valid yet");
253.  +          return this->cpu->lvp->verifyPrediction(this->threadNumber,
254.  +                          this->instAddr(), this->effAddr, temp, _predictedVal,
255.  +                          _classification);
256.  +      }
257.  +
258.  +      void
259.  +      lvpStoreAddressLookup() {
260.  +          if(!this->effAddrValid()) {
```

```
261.  +              panic("Virtual address of store not valid yet");
262.  +          }
263.  +          else {
264.  +              this->cpu->lvp->processStoreAddress(this->threadNumber, this->effAddr);
265.  +          }
266.  +      }
267.  +
268.  +      void
269.  +      speculativeExecOnLoad() {
270.  +          _specExecOnLoad = true;
271.  +      }
272.  +
273.  +      void
274.  +      resetSpeculativeExecOnLoad() {
275.  +          _specExecOnLoad = false;
276.  +      }
277.  +
278.  +      bool
279.  +      isExecOnSpecLoad() {
280.  +          return _specExecOnLoad;
281.  +      }
282.  +
283.  +      bool
284.  +      isSpeculatedLoad() {
285.  +          return (_classification == LVP_PREDICTABLE);
286.  +      }
287.    };
288.
289.    #endif // __CPU_O3_ALPHA_DYN_INST_HH__
290.  diff --git a/gem5/src/cpu/o3/dyn_inst_impl.hh b/gem5/src/cpu/o3/dyn_inst_impl.hh
291.  index 8a6a434..c1c1411 100644
292.  --- a/gem5/src/cpu/o3/dyn_inst_impl.hh
293.  +++ b/gem5/src/cpu/o3/dyn_inst_impl.hh
294.  @@ -106,6 +106,11 @@ BaseO3DynInst<Impl>::initVars()
295.
296.        _numDestMiscRegs = 0;
297.
298.  +      _classification = LVP_STRONG_UNPREDICTABLE;
299.  +      _predictedVal = 0;
300.  +      _predictionCorrect = false;
301.  +      _specExecOnLoad = false;
302.  +
303.    #if TRACING_ON
304.        // Value -1 indicates that particular phase
305.        // hasn't happened (yet).
306.  diff --git a/gem5/src/cpu/o3/iew_impl.hh b/gem5/src/cpu/o3/iew_impl.hh
307.  index 497c532..0435928 100644
308.  --- a/gem5/src/cpu/o3/iew_impl.hh
309.  +++ b/gem5/src/cpu/o3/iew_impl.hh
310.  @@ -57,6 +57,7 @@
311.    #include "debug/Activity.hh"
312.    #include "debug/Drain.hh"
313.    #include "debug/IEW.hh"
314.  +#include "debug/LVP.hh"
315.    #include "debug/O3PipeView.hh"
316.    #include "params/DerivO3CPU.hh"
317.
318.  @@ -957,6 +958,11 @@ template <class Impl>
319.    void
320.    DefaultIEW<Impl>::dispatchInsts(ThreadID tid)
321.    {
322.  +      /**
323.  +       * SATVIK:
324.  +       * LVP will forward the predicted load value to this API
325.  +       */
326.  +
```

```
327.          // Obtain instructions from skid buffer if unblocking, or queue from rename
328.          // otherwise.
329.          std::queue<DynInstPtr> &insts_to_dispatch =
330. @@ -964,6 +970,7 @@ DefaultIEW<Impl>::dispatchInsts(ThreadID tid)
331.              skidBuffer[tid] : insts[tid];
332.
333.          int insts_to_add = insts_to_dispatch.size();
334. +        DPRINTF(IEW, "Insts to dispatch: %d\n", insts_to_add);
335.
336.          DynInstPtr inst;
337.          bool add_to_iq = false;
338. @@ -1087,21 +1094,179 @@ DefaultIEW<Impl>::dispatchInsts(ThreadID tid)
339.                  DPRINTF(IEW, "[tid:%i] Issue: Memory instruction "
340.                          "encountered, adding to LSQ.\n", tid);
341.
342. -                // Reserve a spot in the load store queue for this
343. -                // memory access.
344. -                ldstQueue.insertLoad(inst);
345. +                /**
346. +                  * SATVIK:
347. +                  * A load classified as "constant" by the LCT need not be pushed in
348. +                  * the LDST queue?
349. +                  * -> In case the constant value has changed, this load will need to
350. +                  * be issued again. A CVU lookup can be initiated here and the
351. +                  * result of the lookup can be checked during writeback. This will
352. +                  * give more meaningful results since the lookup will take atleast
353. +                  * 1 cycle. In case the constant load is invalid, the load will be
354. +                  * reissued from writeback.
355. +                  *
356. +                  * Implementation note-
357. +                  * The PC can be sent to the LCT and LVPT for prediction during
358. +                  * fetch. However, it makes sense to store the predicted value and
359. +                  * forward it to other dependents only during dispatch. Since there
360. +                  * is already a delay of more than 1 cycle between fetch and
361. +                  * dispatch, it won't make a difference if the prediction and the
362. +                  * CVU lookup are done here itself. The CVU lookup for a constant
363. +                  * load and the prediction itself will happen for every load
364. +                  * instruction in all cases so it doesn't make a difference where
365. +                  * the two methods are called from.
366. +                  * Inserting the load in the ldstQueue for a load that is correctly
367. +                  * predicted as constant will waste a slot in the ldstQueue. It is
368. +                  * thus better to verify the constant load here and only add it to
369. +                  * the ldstQueue if the verification fails.
370. +                  *
371. +                  * TODO: Not too sure about the wakeDependents call below-
372. +                  */
373. +                std::pair<LVPType, RegVal> prediction = inst->predictLoad(tid);
374. +
375. +                if(prediction.first == LVP_CONSTANT) {
376. +                    // Trigger a CVU lookup of the lvpt index and the load address
377. +                    bool const_valid = false; //inst->verifyConstLoad(tid);
378. +                    if (!const_valid) {
379. +                        // This prediction failed
380. +                        // The CVU will have already incremented the misprediction
381. +                        // counter
382. +                        // Push this load instruction in the ldstQueue
383. +                        ldstQueue.insertLoad(inst);
384. +
385. +                        ++iewDispLoadInsts;
386. +
387. +                        add_to_iq = true;
388. +
389. +                        toRename->iewInfo[tid].dispatchedToLQ++;
390. +                    }
391. +                    else {
392. +                        // This load was predicted correctly
```

```
393.  +                          // A correct prediction for a constant load need not update
394.  +                          // the LCT.
395.  +                          // Write the predicted value to the allocated register and
396.  +                          // forward all values
397.  +                          // Mark this load as executed and ready to commit.
398.  +                          //toRename->iewInfo[tid].dispatchedToLQ++;
399.  +                          //toRename->iewInfo[tid].dispatched++;
400.  +                          //insts_to_dispatch.pop();
401.  +                          //inst->setIssued();
402.  +                          //inst->setExecuted();
403.  +                          //inst->setCanCommit();
404.  +                          //inst->setResultReady();
405.  +                          //instToCommit(inst);
406.  +                          //activityThisCycle();
407.  +                          //add_to_iq = false;
408.  +                          DPRINTF(LVP, "IEW: Const load [sn: %d] found by LVP: 0x%x\n", inst->seqNum, inst-
      >instAddr());
409.  +                          ldstQueue.insertLoad(inst);
410.  +
411.  +                          ++iewDispLoadInsts;
412.  +
413.  +                          add_to_iq = true;
414.  +
415.  +                          toRename->iewInfo[tid].dispatchedToLQ++;
416.  +
417.  +                          // Pass the load value to the destination register
418.  +                          /*if(inst->numDestRegs() == 1) {
419.  +                              auto ptr = inst->renamedDestRegIdx(0);
420.  +                              if(inst->isInteger()) {
421.  +                                  inst->setIntRegOperand(inst->staticInst.get(),
422.  +                                                  0, prediction.second);
423.  +                                  //instQueue.wakeDependents(inst);
424.  +                                  scoreboard->setReg(inst->renamedDestRegIdx(0));
425.  +                                  //DPRINTF(LVP, "LVP Const Inst[%llu]: ox%x getting reg %d as ready:
      %d\n", inst->seqNum, inst->instAddr(), ptr->index(), scoreboard->getReg(ptr));
426.  +                              }
427.  +                              else if(inst->isFloating()) {
428.  +                                  inst->setFloatRegOperandBits(inst->staticInst.get(),
429.  +                                                  0, prediction.second);
430.  +                                  //instQueue.wakeDependents(inst);
431.  +                                  scoreboard->setReg(inst->renamedDestRegIdx(0));
432.  +                                  //DPRINTF(LVP, "LVP Const Inst[%llu]: ox%x getting reg %d as ready:
      %d\n", inst->seqNum, inst->instAddr(), ptr->index(), scoreboard->getReg(ptr));
433.  +                              }
434.  +                              else {
435.  +                                  // This isn't supposed to happen
436.  +                              }
437.  +                          }
438.  +                          else {
439.  +                              // This isn't supposed to happen (except maybe for
440.  +                              // vectors)
441.  +                          }
442.  +                          */
443.  +                      }
444.  +                  }
445.  +              else if(prediction.first == LVP_PREDICTABLE) {
446.  +                  // Need to mark this instruction as predictable so that the CVU
447.  +                  // can verify later. -> this has been done during the
448.  +                  // predictLoad() call.
449.  +
450.  +                  // These loads will follow the normal execution flow: but the
451.  +                  // predicted value will be passed to all consumers.
452.  +                  // The destination register will also need to be tagged with the
453.  +                  // predictable flag so that instructions which consume this
454.  +                  // register are not flushed from the IQ.
455.  +                  if(inst->numDestRegs() == 1) {
```

```
456.  +                    // Tag the destination register for subsequent dependent
457.  +                    // instructions
458.  +                    //inst->tagLVPDestReg(0);
459.  +                    /*
460.  +                    if(inst->isInteger()) {
461.  +                        inst->setIntRegOperand(inst->staticInst.get(),
462.  +                                               0, prediction.second);
463.  +                        instQueue.wakeDependents(inst);
464.  +                        scoreboard->setReg(inst->renamedDestRegIdx(0));
465.  +                    }
466.  +                    else if(inst->isFloating()) {
467.  +                        inst->setFloatRegOperand(inst->staticInst.get(),
468.  +                                                 0, prediction.second);
469.  +                        instQueue.wakeDependents(inst);
470.  +                        scoreboard->setReg(inst->renamedDestRegIdx(0));
471.  +                    }
472.  +                    else {
473.  +                        // This isn't supposed to happen
474.  +                    }
475.  +                    */
476.  +                }
477.  +                else {
478.  +                    // This isn't supposed to happen (except maybe for
479.  +                    // vectors)
480.  +                }
481.
482.  -            ++iewDispLoadInsts;
483.  +                // Reserve a spot in the load store queue for this
484.  +                // memory access.
485.  +                ldstQueue.insertLoad(inst);
486.
487.  -            add_to_iq = true;
488.  +                ++iewDispLoadInsts;
489.
490.  -            toRename->iewInfo[tid].dispatchedToLQ++;
491.  +                add_to_iq = true;
492.  +
493.  +                toRename->iewInfo[tid].dispatchedToLQ++;
494.  +            }
495.  +            else {
496.  +
497.  +                // Reserve a spot in the load store queue for this
498.  +                // memory access.
499.  +                ldstQueue.insertLoad(inst);
500.  +
501.  +                ++iewDispLoadInsts;
502.  +
503.  +                add_to_iq = true;
504.  +
505.  +                toRename->iewInfo[tid].dispatchedToLQ++;
506.  +            }
507.           } else if (inst->isStore()) {
508.               DPRINTF(IEW, "[tid:%i] Issue: Memory instruction "
509.                       "encountered, adding to LSQ.\n", tid);
510.
511.               ldstQueue.insertStore(inst);
512.
513.  +            /**
514.  +             * SATVIK:
515.  +             * Stores could finish execution out of order with respect to other
516.  +             * loads. Processing the store address after execution has completed
517.  +             * could invalidate an older load which is constant.
518.  +             */
519.  +            //inst->lvpStoreAddressLookup();
520.  +
521.               ++iewDispStoreInsts;
```

```
522.
523.                    if (inst->isStoreConditional()) {
524.    @@ -1294,6 +1459,9 @@ DefaultIEW<Impl>::executeInsts()
525.                    } else if (inst->isLoad()) {
526.                        // Loads will mark themselves as executed, and their writeback
527.                        // event adds the instruction to the queue to commit
528.    +                   if(inst->isConstLoad() && !inst->strictlyOrdered() && !inst->isInstPrefetch())
529.    +                       inst->verifyConstLoad(inst->threadNumber);
530.    +
531.                        fault = ldstQueue.executeLoad(inst);
532.
533.                        if (inst->isTranslationDelayed() &&
534.    @@ -1309,6 +1477,26 @@ DefaultIEW<Impl>::executeInsts()
535.                        if (inst->isDataPrefetch() || inst->isInstPrefetch()) {
536.                            inst->fault = NoFault;
537.                        }
538.    +
539.    +                   /**
540.    +                    * SATVIK:
541.    +                    * Check the LVP prediction here
542.    +                   if(inst->isExecuted() && fault == NoFault) {
543.    +                       if (inst->numDestRegs() == 1) {
544.    +                           bool verify = inst->verifyPrediction(0);
545.    +                           if(!verify && inst->isSpeculatedLoad()) {
546.    +                               // Add the destination register to a list of
547.    +                               // mispredicted registers so that the instns in the
548.    +                               // IQ know that they have to execute again.
549.    +                               //instQueue.addToMispredictList(inst, 0);
550.    +                           }
551.    +                       }
552.    +                       else {
553.    +                           // This shouldn't happen
554.    +                       }
555.    +                   }
556.    +                   */
557.    +
558.                    } else if (inst->isStore()) {
559.                        fault = ldstQueue.executeStore(inst);
560.
561.    @@ -1334,6 +1522,16 @@ DefaultIEW<Impl>::executeInsts()
562.                        activityThisCycle();
563.                    }
564.
565.    +                   /**
566.    +                    * SATVIK:
567.    +                    * At this point, if there is no error, this store instruction
568.    +                    * will commit at some point. So invalidating any CVU CAM
569.    +                    * entries that have this store address should be done here.
570.    +                    */
571.    +                   if(inst->isExecuted() && fault == NoFault) {
572.    +                       inst->lvpStoreAddressLookup();
573.    +                   }
574.    +
575.                        // Store conditionals will mark themselves as
576.                        // executed, and their writeback event will add the
577.                        // instruction to the queue to commit.
578.    @@ -1352,6 +1550,9 @@ DefaultIEW<Impl>::executeInsts()
579.                        inst->forwardOldRegs();
580.                    }
581.
582.    +           /**
583.    +            * Do this only if the source operands were not speculated.
584.    +            */
585.                inst->setExecuted();
586.
587.                instToCommit(inst);
```

```
588.    @@ -1488,6 +1689,41 @@ DefaultIEW<Impl>::writebackInsts()
589.            // are first sent to commit.  Instead commit must tell the LSQ
590.            // when it's ready to execute the strictly ordered load.
591.            if (!inst->isSquashed() && inst->isExecuted() && inst->getFault() == NoFault) {
592.  +             if(inst->isLoad()) {
593.  +                 if(inst->isConstPredictionCorrect() && !inst->strictlyOrdered() && !inst-
      >isInstPrefetch()) {
594.  +                     // Pass the load value to the destination register
595.  +                     //inst->setCanCommit();
596.  +                     //checkMisprediction(inst);
597.  +                     if(inst->numDestRegs() == 1) {
598.  +                         auto ptr = inst->renamedDestRegIdx(0);
599.  +                         if(ptr->isIntPhysReg()) {
600.  +                             inst->setIntRegOperand(inst->staticInst.get(),
601.  +                                                     0, inst->getPredictedValue());
602.  +                             //instQueue.wakeDependents(inst);
603.  +                             scoreboard->setReg(ptr);
604.  +                             DPRINTF(LVP, "LVP Const Inst[%llu]: ox%x setting reg %d as ready\n",
      inst->seqNum, inst->instAddr(), ptr->index());
605.  +                         }
606.  +                         else if(ptr->isFloatPhysReg()) {
607.  +                             inst->setFloatRegOperandBits(inst->staticInst.get(),
608.  +                                                     0, inst->getPredictedValue());
609.  +                             //instQueue.wakeDependents(inst);
610.  +                             scoreboard->setReg(inst->renamedDestRegIdx(0));
611.  +                             DPRINTF(LVP, "LVP Const Inst[%llu]: ox%x setting reg %d as ready\n",
      inst->seqNum, inst->instAddr(), ptr->index());
612.  +                         }
613.  +                         else {
614.  +                             // This isn't supposed to happen
615.  +                         }
616.  +                     }
617.  +                     else {
618.  +                         // This isn't supposed to happen (except maybe for
619.  +                         // vectors)
620.  +                     }
621.  +                 }
622.  +                 else {
623.  +                     if(inst->numDestRegs() == 1)
624.  +                         inst->verifyPrediction(0);
625.  +                 }
626.  +             }
627.                int dependents = instQueue.wakeDependents(inst);
628.
629.                for (int i = 0; i < inst->numDestRegs(); i++) {
630.  diff --git a/gem5/src/cpu/o3/inst_queue_impl.hh b/gem5/src/cpu/o3/inst_queue_impl.hh
631.  index ff5b3be..c878485 100644
632.  --- a/gem5/src/cpu/o3/inst_queue_impl.hh
633.  +++ b/gem5/src/cpu/o3/inst_queue_impl.hh
634.  @@ -765,6 +765,14 @@ InstructionQueue<Impl>::processFUCompletion(const DynInstPtr &inst, int fu_idx)
635.        instsToExecute.push_back(inst);
636.  }
637.
638.  +/**
639.  + * SATVIK:
640.  + * If an instruction reaches the head of the IQ, check if any of its source
641.  + * registers is tagged with the predicted load. If yes, do not erase the
642.  + * instruction from the IQ. Only after the CVU validates the load can the
643.  + * instruction be removed from the IQ.
644.  + */
645.  +
646.    // @todo: Figure out a better way to remove the squashed items from the
647.    // lists.  Checking the top item of each list to see if it's squashed
648.    // wastes time and forces jumps.
649.  @@ -908,6 +916,11 @@ InstructionQueue<Impl>::scheduleReadyInsts()
650.                // complete.
```

```
651.                    ++freeEntries;
652.                    count[tid]--;
653.  +                 /**
654.  +                  * SATVIK:
655.  +                  * This will be conditional- if the source registers are not
656.  +                  * speculated values, then clear the entry from the IQ.
657.  +                  */
658.                    issuing_inst->clearInIQ();
659.                } else {
660.                    memDepUnit[tid].issue(issuing_inst);
661.  diff --git a/gem5/src/cpu/o3/lsq_impl.hh b/gem5/src/cpu/o3/lsq_impl.hh
662.  index c4cb45e..2584860 100644
663.  --- a/gem5/src/cpu/o3/lsq_impl.hh
664.  +++ b/gem5/src/cpu/o3/lsq_impl.hh
665.  @@ -735,6 +735,12 @@ LSQ<Impl>::pushRequest(const DynInstPtr& inst, bool isLoad, uint8_t *data,
666.            inst->effSize = size;
667.            inst->effAddrValid(true);
668.
669.  +         //SATVIK
670.  +         // CHeck the CVU CAM for this load
671.  +         if(inst->isConstLoad()) {
672.  +             inst->verifyConstLoad(0);
673.  +         }
674.  +
675.            if (cpu->checker) {
676.                inst->reqToVerify = std::make_shared<Request>(*req->request());
677.            }
678.  diff --git a/gem5/src/cpu/o3/lsq_unit.hh b/gem5/src/cpu/o3/lsq_unit.hh
679.  index 3d6e3f0..4a94969 100644
680.  --- a/gem5/src/cpu/o3/lsq_unit.hh
681.  +++ b/gem5/src/cpu/o3/lsq_unit.hh
682.  @@ -683,6 +683,19 @@ LSQUnit<Impl>::read(LSQRequest *req, int load_idx)
683.            load_inst->recordResult(true);
684.        }
685.
686.  +     // SATVIK
687.  +     if(load_inst->isConstLoad() && load_inst->isConstPredictionCorrect()) {
688.  +         // This is basically forwarding values from a local resource
689.  +         // The data will be written during IEW WB
690.  +         // Here, we just schedule the writeback event of the load
691.  +         DPRINTF(LSQUnit, "Const load[%llu]: 0x%x scheduling a WB event\n", load_inst->seqNum,
        load_inst->instAddr());
692.  +         load_inst->memData = new uint8_t[MaxDataBytes];
693.  +         PacketPtr main_pkt = new Packet(req->mainRequest(), MemCmd::ReadReq);
694.  +         main_pkt->dataStatic(load_inst->memData);
695.  +         WritebackEvent *wb = new WritebackEvent(load_inst, main_pkt, this);
696.  +         cpu->schedule(wb, curTick());
697.  +         return NoFault;
698.  +     }
699.        if (req->mainRequest()->isLocalAccess()) {
700.            assert(!load_inst->memData);
701.            assert(!load_inst->inHtmTransactionalState());
702.  diff --git a/gem5/src/cpu/o3/lsq_unit_impl.hh b/gem5/src/cpu/o3/lsq_unit_impl.hh
703.  index 808a671..1e478e3 100644
704.  --- a/gem5/src/cpu/o3/lsq_unit_impl.hh
705.  +++ b/gem5/src/cpu/o3/lsq_unit_impl.hh
706.  @@ -617,6 +617,14 @@ LSQUnit<Impl>::executeLoad(const DynInstPtr &inst)
707.
708.        assert(!inst->isSquashed());
709.
710.  +     //if(inst->isConstPredictionCorrect() && !inst->strictlyOrdered() && !inst->isInstPrefetch()) {
711.  +     //     DPRINTF(LSQUnit, "Load PC %s was correctly predicted as constant sn:%llu", inst->pcState(),
        inst->seqNum);
712.  +     //     inst->setExecuted();
713.  +     //     iewStage->instToCommit(inst);
714.  +     //     iewStage->activityThisCycle();
```

```
715.  +      //     iewStage->checkMisprediction(inst);
716.  +      //     return NoFault;
717.  +      //}
718.         load_fault = inst->initiateAcc();
719.
720.         if (load_fault == NoFault && !inst->readMemAccPredicate()) {
721.  @@ -669,6 +677,20 @@ LSQUnit<Impl>::executeLoad(const DynInstPtr &inst)
722.             }
723.         }
724.
725.  +      /**
726.  +       * SATVIK:
727.  +       * We can call the CVU verification API here to check if the load was
728.  +       * judged "predictable" by the LCT and if so, whether the actual load value
729.  +       * is equal to the predicted value. If there's a mismatch,
730.  +       * all its dependents will have to be issued again, with the correct values
731.  +       *
732.  +       * If the load was deemed "predictable" -> reissue all dependent instructions
733.  +       * -> Not sure how to do this
734.  +       *
735.  +       * UPDATE: It's better to verify the load in the IEW stage since that will
736.  +       * give us direct access to the IQ
737.  +       */
738.  +
739.         return load_fault;
740.     }
741.
742.  @@ -1099,6 +1121,16 @@ LSQUnit<Impl>::writeback(const DynInstPtr &inst, PacketPtr pkt)
743.             return;
744.         }
745.
746.  +      // SATVIK
747.  +      if(inst->isConstLoad() && inst->isConstPredictionCorrect()) {
748.  +          DPRINTF(LSQUnit, "Load PC %s was correctly predicted as constant sn:%llu", inst->pcState(),
      inst->seqNum);
749.  +          inst->setExecuted();
750.  +          iewStage->instToCommit(inst);
751.  +          iewStage->activityThisCycle();
752.  +          iewStage->checkMisprediction(inst);
753.  +          return;
754.  +      }
755.  +
756.         if (!inst->isExecuted()) {
757.             inst->setExecuted();
758.
759.  diff --git a/gem5/src/cpu/o3/regfile.hh b/gem5/src/cpu/o3/regfile.hh
760.  index 922089c..624d262 100644
761.  --- a/gem5/src/cpu/o3/regfile.hh
762.  +++ b/gem5/src/cpu/o3/regfile.hh
763.  @@ -43,6 +43,7 @@
764.   #define __CPU_O3_REGFILE_HH__
765.
766.   #include <vector>
767.  +#include <list>
768.
769.   #include "arch/types.hh"
770.   #include "base/trace.hh"
771.  @@ -131,6 +132,14 @@ class PhysRegFile
772.       /** Mode in which vector registers are addressed. */
773.       VecMode vecMode;
774.
775.  +    /**
776.  +     * List containing physical register pointers for registers that contain a
777.  +     * predicted load.
778.  +     */
779.  +    std::list<PhysRegIdPtr> regsWithPredictedLoad;
```

```
780.  +
781.  +    std::list<PhysRegIdPtr> mispredictList;
782.  +
783.       public:
784.          /**
785.           * Constructs a physical register file with the specified amount of
786.  @@ -149,6 +158,69 @@ class PhysRegFile
787.           */
788.          ~PhysRegFile() {}
789.
790.  +    /**
791.  +     * SATVIK:
792.  +     * Implement a map and an API here which tags a PhysRegIdPtr with a
793.  +     * predicted load. Every time an instruction is to issue, this map will be
794.  +     * checked to see if any source operand of the instn uses a register that
795.  +     * has been tagged. If there's a match, the instruction is not popped from
796.  +     * the IQ.
797.  +     */
798.  +
799.  +    /**
800.  +     * @brief      Checks whether the given source register contains a value
801.  +     *             that has been predicted for a load but not verified yet
802.  +     *
803.  +     * @param[in]  src_reg  The source register
804.  +     *
805.  +     * @return     True if the value is a predicted one
806.  +     */
807.  +    bool checkPredictedLoadRegister (PhysRegIdPtr src_reg) {
808.  +        auto itr = regsWithPredictedLoad.begin();
809.  +        while (itr != regsWithPredictedLoad.end()) {
810.  +            if (*itr == src_reg) return true;
811.  +            itr++;
812.  +        }
813.  +        return false;
814.  +    }
815.  +
816.  +    /**
817.  +     * @brief      Inserts a physical register that contains a speculated value
818.  +     *
819.  +     * @param[in]  reg    The register
820.  +     */
821.  +    void insertPredictedLoadRegister (PhysRegIdPtr reg) {
822.  +        regsWithPredictedLoad.push_back(reg);
823.  +    }
824.  +
825.  +    /**
826.  +     * @brief      Clears a physical register pointer from the list
827.  +     *
828.  +     * @param[in]  reg    The register
829.  +     */
830.  +    void clearPredictedLoadRegister (PhysRegIdPtr reg) {
831.  +        auto itr = regsWithPredictedLoad.begin();
832.  +        while(itr != regsWithPredictedLoad.end()) {
833.  +            if (*itr == reg) break;
834.  +            itr++;
835.  +        }
836.  +        if (itr != regsWithPredictedLoad.end())
837.  +            regsWithPredictedLoad.erase(itr);
838.  +    }
839.  +
840.  +    void addToMispredictList(PhysRegIdPtr reg) {
841.  +        mispredictList.push_back(reg);
842.  +    }
843.  +
844.  +    bool searchMispredictList(PhysRegIdPtr src_reg) {
845.  +        auto itr = mispredictList.begin();
```

```
846.  +            while(itr != mispredictList.end()) {
847.  +                if(*itr == src_reg) return true;
848.  +                itr++;
849.  +            }
850.  +            return false;
851.  +        }
852.  +
853.         /** Initialize the free list */
854.         void initFreeList(UnifiedFreeList *freeList);
855.
856.
```

## III. MODIFICATIONS FOR MINOR CPU

```
1.   diff --git a/gem5/src/cpu/minor/MinorCPU.py b/gem5/src/cpu/minor/MinorCPU.py
2.   index 1329dfb..29732bc 100644
3.   --- a/gem5/src/cpu/minor/MinorCPU.py
4.   +++ b/gem5/src/cpu/minor/MinorCPU.py
5.   @@ -47,6 +47,8 @@ from m5.objects.DummyChecker import DummyChecker
6.     from m5.objects.BranchPredictor import *
7.     from m5.objects.TimingExpr import TimingExpr
8.
9.   +from m5.objects.LoadValuePredictionUnit import LoadValuePredictionUnit
10.  +
11.    from m5.objects.FuncUnit import OpClass
12.
13.    class MinorOpClass(SimObject):
14.  @@ -273,7 +275,7 @@ class MinorCPU(BaseCPU):
15.        executeSetTraceTimeOnIssue = Param.Bool(False,
16.            "Set inst. trace times to be issue times")
17.
18.  -     executeAllowEarlyMemoryIssue = Param.Bool(True,
19.  +     executeAllowEarlyMemoryIssue = Param.Bool(False,
20.            "Allow mem refs to be issued to the LSQ before reaching the head of"
21.            " the in flight insts queue")
22.
23.  @@ -283,6 +285,8 @@ class MinorCPU(BaseCPU):
24.        branchPred = Param.BranchPredictor(TournamentBP(
25.            numThreads = Parent.numThreads), "Branch Predictor")
26.
27.  +     loadValPred = Param.LoadValuePredictionUnit(LoadValuePredictionUnit(), "Load value predictor")
28.  +
29.        def addCheckerCpu(self):
30.            print("Checker not yet supported by MinorCPU")
31.            exit(1)
32.   diff --git a/gem5/src/cpu/minor/SConscript b/gem5/src/cpu/minor/SConscript
33.   index 75cc940..2ef617a 100644
34.   --- a/gem5/src/cpu/minor/SConscript
35.   +++ b/gem5/src/cpu/minor/SConscript
36.   @@ -60,6 +60,7 @@ if 'MinorCPU' in env['CPU_MODELS']:
37.        DebugFlag('MinorCPU', 'Minor CPU-level events')
38.        DebugFlag('MinorExecute', 'Minor Execute stage')
39.        DebugFlag('MinorInterrupt', 'Minor interrupt handling')
40.  +     DebugFlag('MinorLoadPredictor', 'LVP in minor CPU')
41.        DebugFlag('MinorMem', 'Minor memory accesses')
42.        DebugFlag('MinorScoreboard', 'Minor Execute register scoreboard')
43.        DebugFlag('MinorTrace', 'MinorTrace cycle-by-cycle state trace')
44.   diff --git a/gem5/src/cpu/minor/cpu.cc b/gem5/src/cpu/minor/cpu.cc
45.   index a375e07..a051a51 100644
46.   --- a/gem5/src/cpu/minor/cpu.cc
47.   +++ b/gem5/src/cpu/minor/cpu.cc
48.   @@ -47,7 +47,8 @@
49.
50.    MinorCPU::MinorCPU(MinorCPUParams *params) :
51.        BaseCPU(params),
```

```
52. -     threadPolicy(params->threadPolicy)
53. +     threadPolicy(params->threadPolicy),
54. +     loadValuePredictor(params->loadValPred)
55.  {
56.       /* This is only written for one thread at the moment */
57.       Minor::MinorThread *thread;
58. diff --git a/gem5/src/cpu/minor/cpu.hh b/gem5/src/cpu/minor/cpu.hh
59. index 579a96b..24bdb6c 100644
60. --- a/gem5/src/cpu/minor/cpu.hh
61. +++ b/gem5/src/cpu/minor/cpu.hh
62. @@ -44,6 +44,8 @@
63.  #ifndef __CPU_MINOR_CPU_HH__
64.  #define __CPU_MINOR_CPU_HH__
65.
66. +#include "cpu/lvp/load_value_prediction_unit.hh"
67. +
68.  #include "cpu/minor/activity.hh"
69.  #include "cpu/minor/stats.hh"
70.  #include "cpu/base.hh"
71. @@ -186,6 +188,9 @@ class MinorCPU : public BaseCPU
72.       *  already been idled.  The stage argument should be from the
73.       *  enumeration Pipeline::StageId */
74.      void wakeupOnEvent(unsigned int stage_id);
75. +
76. +    public:
77. +        LoadValuePredictionUnit* loadValuePredictor;
78.  };
79.
80.  #endif /* __CPU_MINOR_CPU_HH__ */
81. diff --git a/gem5/src/cpu/minor/decode.cc b/gem5/src/cpu/minor/decode.cc
82. index b07ca4a..edee685 100644
83. --- a/gem5/src/cpu/minor/decode.cc
84. +++ b/gem5/src/cpu/minor/decode.cc
85. @@ -233,6 +233,15 @@ Decode::evaluate()
86.                  dynInstAddTracing(output_inst, parent_static_inst, cpu);
87.  #endif
88.
89. +                // If this is a load instruction, check the LVP unit to see if
90. +                // it can be predicted and tag the instruction as such
91. +                if(output_inst->isInst() && output_inst->isMemRef() && output_inst->staticInst-
    >isLoad())
92. +                {
93. +                    auto lvp_result = cpu.loadValuePredictor->lookup(tid, output_inst->pc.instAddr());
94. +                    output_inst->loadPredicted = lvp_result.taken;
95. +                    output_inst->loadPredictedValue = lvp_result.value;
96. +                }
97. +
98.                  /* Step to next sequence number */
99.                  decode_info.execSeqNum++;
100.
101.  diff --git a/gem5/src/cpu/minor/dyn_inst.hh b/gem5/src/cpu/minor/dyn_inst.hh
102.  index b90e277..586228e 100644
103.  --- a/gem5/src/cpu/minor/dyn_inst.hh
104.  +++ b/gem5/src/cpu/minor/dyn_inst.hh
105.  @@ -184,6 +184,25 @@ class MinorDynInst : public RefCounted
106.       /** Predicted branch target */
107.       TheISA::PCState predictedTarget;
108.
109.  +    /** Robert
110.  +     * Track if the load instruction was predicted and, if so, what the predicted value is.
111.  +     * Load prediction done in decode, acted upon by execute in scheduling and load value comparison.
    */
112.  +
113.  +    /** The classification of the load prediction for this instruction. Only valid for loads. */
114.  +    LVPType loadPredicted;
115.  +
```

```
116.  +      /** The register value to be forwarded in the load prediction, depending on how loadPredicted
      set. */
117.  +      RegVal loadPredictedValue;
118.  +
119.  +      /** Indicates that the processor executed this as a constant load, bypassing memory system */
120.  +      bool executedAsConstant = false;
121.  +
122.  +      /** Stores the virtual address of the instruction */
123.  +      Addr effAddr;
124.  +
125.  +      /** Indicates that the virtual address has been set and effAddr is valid */
126.  +      bool effAddrValid = false;
127.  +
128.         /** Fields only set during execution */
129.
130.         /** FU this instruction is issued to */
131.  diff --git a/gem5/src/cpu/minor/execute.cc b/gem5/src/cpu/minor/execute.cc
132.  index 45ca002..321ed33 100644
133.  --- a/gem5/src/cpu/minor/execute.cc
134.  +++ b/gem5/src/cpu/minor/execute.cc
135.  @@ -51,6 +51,7 @@
136.   #include "debug/ExecFaulting.hh"
137.   #include "debug/MinorExecute.hh"
138.   #include "debug/MinorInterrupt.hh"
139.  +#include "debug/MinorLoadPredictor.hh"
140.   #include "debug/MinorMem.hh"
141.   #include "debug/MinorTrace.hh"
142.   #include "debug/PCEvent.hh"
143.  @@ -374,6 +375,60 @@ Execute::handleMemResponse(MinorDynInstPtr inst,
144.                  static_cast<unsigned int>(packet->getConstPtr<uint8_t>()[0]));
145.             }
146.
147.  +          /** Robert
148.  +           * Check if the value loaded from memory was correct and update the
149.  +           * LVP prediction tables with the result. */
150.  +          if (is_load)
151.  +          {
152.  +              bool load_predicted_correctly = false;
153.  +              if (inst->loadPredicted != LVP_CONSTANT)
154.  +              {
155.  +                  // Get the data from the packet (borrowed from memhelpers.hh)
156.  +                  uint64_t mem;
157.  +
158.  +                  switch (packet->getSize()) {
159.  +                    case 1:
160.  +                      mem = packet->getConstPtr<uint8_t>()[0];
161.  +                      break;
162.  +                    case 2:
163.  +                      mem = packet->getConstPtr<uint16_t>()[0];
164.  +                      break;
165.  +                    case 4:
166.  +                      mem = packet->getConstPtr<uint32_t>()[0];
167.  +                      break;
168.  +                    case 8:
169.  +                      mem = packet->getConstPtr<uint64_t>()[0];
170.  +                      break;
171.  +                    default:
172.  +                      panic("Unhandled size in handleMemResponse.\n");
173.  +                  }
174.  +                  // DPRINTF(MinorLoadPredictor, "Verifying prediction: pc: %#x, predicted value: %#x,
      correct value: %#x\n",
175.  +                  //     inst->pc.instAddr(), inst->loadPredictedValue, mem);
176.  +                  load_predicted_correctly =
177.  +                      cpu.loadValuePredictor->verifyPrediction(
178.  +                          thread_id,
179.  +                          inst->pc.instAddr(),
```

```
180.  +                                  inst->effAddr,
181.  +                                  mem,
182.  +                                  inst->loadPredictedValue,
183.  +                                  inst->loadPredicted);
184.  +
185.  +                      /* Experimenting with predictable forwarding
186.  +                      if (!load_predicted_correctly && inst->loadPredicted == LVP_PREDICTABLE)
187.  +                      {
188.  +                          // TODO the instruction was incorrectly predicted, we need to rewind this
      instruction
189.  +                          // and any dependent instructions that used this value
190.  +                          DPRINTF(MinorLoadPredictor, "Predicted the wrong value, need to fix up!\n");
191.  +                          auto reg_id = thread->flattenRegId(inst->staticInst->destRegIdx(0));
192.  +                          thread->setIntRegFlat(reg_id.flatIndex(), mem);
193.  +                          BranchData &lvpResetTarget = *out.inputWire;
194.  +
195.  +                          updateBranchData(thread_id, BranchData::UnpredictedBranch, inst, inst-
      >pc.nextInstAddr(), lvpResetTarget);
196.  +                      }
197.  +                      */
198.  +                  }
199.  +              }
200.  +
201.              /* Complete the memory access instruction */
202.              fault = inst->staticInst->completeAcc(packet, &context,
203.                  inst->traceData);
204.  @@ -387,7 +442,9 @@ Execute::handleMemResponse(MinorDynInstPtr inst,
205.                  /* Stores need to be pushed into the store buffer to finish
206.                   *  them off */
207.                  if (response->needsToBeSentToStoreBuffer())
208.  +                {
209.                      lsq.sendStoreToStoreBuffer(response);
210.  +                }
211.              }
212.          } else {
213.              fatal("There should only ever be reads, "
214.  @@ -464,7 +521,7 @@ Execute::executeMemRefInst(MinorDynInstPtr inst, BranchData &branch,
215.              DPRINTF(MinorExecute, "Initiating memRef inst: %s\n", *inst);
216.
217.              Fault init_fault = inst->staticInst->initiateAcc(&context,
218.  -                inst->traceData);
219.  +                  inst->traceData);
220.
221.              if (inst->inLSQ) {
222.                  if (init_fault != NoFault) {
223.  @@ -930,6 +987,26 @@ Execute::commitInst(MinorDynInstPtr inst, bool early_memory_issue,
224.              bool completed_mem_inst = executeMemRefInst(inst, branch,
225.                  predicate_passed, fault);
226.
227.  +          /* Experiments with forwarding load values
228.  +          if (inst->staticInst->isLoad() && inst->loadPredicted == LVP_PREDICTABLE) {
229.  +              bool could_forward = true;
230.  +              switch (inst->flatDestRegIdx->classValue()) {
231.  +                  case IntRegClass:
232.  +                      thread->setIntRegFlat(inst->flatDestRegIdx->flatIndex(), inst-
      >loadPredictedValue);
233.  +                      break;
234.  +                  case FloatRegClass:
235.  +                      thread->setFloatRegFlat(inst->flatDestRegIdx->flatIndex(), inst-
      >loadPredictedValue);
236.  +                      break;
237.  +                  default:
238.  +                      could_forward = false;
239.  +              }
240.  +              if (could_forward)
241.  +                  scoreboard[inst->id.threadId].clearInstDests(inst, inst->isMemRef());
```

```
242.  +
243.  +            completed_mem_inst = could_forward;
244.  +        }
245.  +        */
246.  +
247.         if (completed_mem_inst && fault != NoFault) {
248.             if (early_memory_issue) {
249.                 DPRINTF(MinorExecute, "Fault in early executing inst: %s\n",
250.  @@ -1147,6 +1224,9 @@ Execute::commit(ThreadID thread_id, bool only_commit_microops, bool discard,
251.
252.                 lsq.popResponse(mem_response);
253.             } else {
254.  +            if(inst->executedAsConstant) {
255.  +                DPRINTF(MinorLoadPredictor, "Handling mem response of constant load %s.\n",
      *inst);
256.  +            }
257.                 handleMemResponse(inst, mem_response, branch, fault);
258.                 committed_inst = true;
259.             }
260.  diff --git a/gem5/src/cpu/minor/lsq.cc b/gem5/src/cpu/minor/lsq.cc
261.  index 106b51b..c583db3 100644
262.  --- a/gem5/src/cpu/minor/lsq.cc
263.  +++ b/gem5/src/cpu/minor/lsq.cc
264.  @@ -49,6 +49,7 @@
265.   #include "cpu/utils.hh"
266.   #include "debug/Activity.hh"
267.   #include "debug/MinorMem.hh"
268.  +#include "debug/MinorLoadPredictor.hh"
269.
270.   namespace Minor
271.   {
272.  @@ -960,7 +961,7 @@ LSQ::StoreBuffer::minorTrace() const
273.
274.   void
275.   LSQ::tryToSendToTransfers(LSQRequestPtr request)
276.  -{
277.  +{
278.       if (state == MemoryNeedsRetry) {
279.           DPRINTF(MinorMem, "Request needs retry, not issuing to"
280.               " memory until retry arrives\n");
281.  @@ -1042,6 +1043,13 @@ LSQ::tryToSendToTransfers(LSQRequestPtr request)
282.           return;
283.       }
284.       } else {
285.  +
286.  +        /* Send store instructions to LVP to invalidate CVU entries. */
287.  +        if (request->inst->staticInst->isStore() && request->request->hasVaddr()) {
288.  +            DPRINTF(MinorLoadPredictor, "Processing a store to CVU.\n");
289.  +            cpu.loadValuePredictor->processStoreAddress(request->inst->id.threadId, request->request-
      >getVaddr());
290.  +        }
291.  +
292.           /* Store.  Can it be sent to the store buffer? */
293.           if (bufferable && !request->request->isLocalAccess()) {
294.               request->setState(LSQRequest::StoreToStoreBuffer);
295.  @@ -1164,6 +1172,15 @@ LSQ::tryToSendToTransfers(LSQRequestPtr request)
296.               request->issuedToMemory = true;
297.           }
298.
299.  +        // deal with LVP stuff
300.  +        request->inst->executedAsConstant = request->inst->staticInst->isLoad()
301.  +            && request->inst->loadPredicted == LVP_CONSTANT
302.  +            && cpu.loadValuePredictor->processLoadAddress(request->inst->id.threadId, request->inst-
      >pc.instAddr(), request->request->getVaddr());
303.  +
304.  +        if(request->inst->executedAsConstant) {
```

```
305.  +              DPRINTF(MinorLoadPredictor, "Found instruction %s that can execute as constant\n",
      *request->inst);
306.  +          }
307.  +
308.              if (tryToSend(request)) {
309.                  moveFromRequestsToTransfers(request);
310.              }
311.  @@ -1191,7 +1208,25 @@ LSQ::tryToSend(LSQRequestPtr request)
312.               *  so the response can be correctly handled */
313.              assert(packet->findNextSenderState<LSQRequest>());
314.
315.  -          if (request->request->isLocalAccess()) {
316.  +          /* If the request can be executed as a constant, respond with the loadPredictedValue*/
317.  +          if(request->inst->executedAsConstant) {
318.  +              packet->setData((uint8_t *)&request->inst->loadPredictedValue);
319.  +
320.  +              DPRINTF(MinorLoadPredictor, "Forwarded values on constant load %s.\n", *request->inst);
321.  +
322.  +              request->stepToNextPacket();
323.  +              ret = request->sentAllPackets();
324.  +
325.  +              if (!ret) {
326.  +                  DPRINTF(MinorLoadPredictor, "Constant load access has another packet: %s\n",
327.  +                      *(request->inst));
328.  +              }
329.  +
330.  +              if (ret)
331.  +                  request->setState(LSQRequest::Complete);
332.  +              else
333.  +                  request->setState(LSQRequest::RequestIssuing);
334.  +          } else if (request->request->isLocalAccess()) {
335.                  ThreadContext *thread =
336.                      cpu.getContext(cpu.contextToThread(
337.                              request->request->contextId()));
338.  @@ -1654,6 +1689,12 @@ LSQ::pushRequest(MinorDynInstPtr inst, bool isLoad, uint8_t *data,
339.          inst->inLSQ = true;
340.          request->startAddrTranslation();
341.
342.  +      if (request->request->hasVaddr())
343.  +      {
344.  +          inst->effAddrValid = request->request->hasVaddr();
345.  +          inst->effAddr = request->request->getVaddr();
346.  +      }
347.  +
348.          return inst->translationFault;
349.      }
350.
```