

Report:

Memory Coalescing for GPUs

Assignment 5 GPU Architecture Course

Ramyad Hadidi - November 20, 2015

1 Introduction

In this assignment, trace generation is implemented in previously developed Harmonica ISA emulator ¹. Also, in emulator different mechanisms for memory coalescing is integrated. Therefore, with the help of the emulator a study on memory coalescing effects on cache has been done using Dinero IV² a Trace-Driven Uniprocessor Cache Simulator. In this report first we show how memory coalescing code was integrated, then a study about impact of cache line size on cache performance, and finally another study about memory mapping schemes and its impact on cache performance.

The benchmarks that we are using is *small* version of *vecsum* and *scan* applications with 8 threads/warp and 32 threads/warp³. For Dinero, we use one level cache with fixed properties except of cache line size in the first study. We have used a 32KB, 4 way set-associative cache with 8 banks.

2 Developed Code

In each execution step, the main class of *core_c* will call *step()* function, which it will call all warp in the core and execute them. In each warp after executing all threads, if threads encounter ST or LD, they will save their address to variable *m_coalMemAddr* and the number of them in *m_uniqueCoalMemAddr*. At the end of the execution of a warp, a *coalesce()* function will be called and address will be coalesced based on global configurations, such as cache line size and mapping scheme. Here is the code:

```
1 void warp_c::coalesce() {
2     unsigned maximumDistance = (1LL << CACHE_LINE_IN_BYTE/
3         WORD_SIZE_IN_BYTE) - 1;
4     set<Addr> memAddresses;
5     set<Addr> coalMemAddr;
6
7     //Get inputs
8     for (unsigned int i=0; i<m_uniqueCoalMemAddr ; i++)
9         memAddresses.insert( m_coalMemAddr[i] );
10
11     //Note: set will not insert duplicate element
12     for (set<Addr>::iterator itA = memAddresses.begin(); itA
        != memAddresses.end(); ) {
        bool coalFound = false;
```

¹https://github.com/ramyadhadidi/harmonica_emulator_proj

²<http://pages.cs.wisc.edu/~markhill/DineroIV/>

³16 threads/warp was not supplied for small version of traces

```

13     for (set<Addr>::iterator itB = memAddresses.begin(); itB
14         != memAddresses.end()); {
15         Addr firstElem = *itA >> RIGHT_BIT_SHIFT_COAL;
16         Addr secondElem = *itB >> RIGHT_BIT_SHIFT_COAL;
17         Addr mask = (1LL << RIGHT_BIT_SHIFT_COAL) - 1;
18         Addr lowerFirstElem = *itA & mask;
19         Addr lowerSecondElem = *itB & mask;
20         if ( (abs(secondElem - firstElem) < maximumDistance)
21             && (*itB != *itA) && (lowerFirstElem ==
22                 lowerSecondElem) ) {
23             coalFound = true;
24             memAddresses.erase(itB++);
25             coalMemAddr.insert(*itA);
26         }
27         else {
28             itB++;
29             m_core->statNotCoalesced ++;
30         }
31     }
32     //Unique item itslef, insert it
33     if (!coalFound)
34         coalMemAddr.insert(*itA);
35     memAddresses.erase(itA++);
36 }
37
38 //Restore outputs
39 int sizeOld = m_uniqueCoalMemAddr;
40 m_uniqueCoalMemAddr = coalMemAddr.size();
41 unsigned int i=0;
42 for (set<Addr>::iterator it = coalMemAddr.begin(); it !=
    coalMemAddr.end(); ++it, ++i)
    m_coalMemAddr[i] = *it;
    m_core->statCoalesced += sizeOld-m_uniqueCoalMemAddr;
}

```

First, the elements of `m_coalMemAddr` have been copied to a set. Set is used because it will not allow duplicate elements. Also, set will order the addresses which enable us to write a 2 for loop for coalescing and select the first item as the representative cache line if coalescing criteria passed. In the 2 for loop body we check if the 2 addresses has coalescing ability based on maximum distance (based on cache line size) and also if coalescing is based on higher value bits we check equality of lower bits. In addition I have used a statistics to show how many coalescing is occurring versus not coalesced accesses.

3 Cache Line Size Study

In this study coalescing based on lower bits is done and the impact of changing cache line size on different cache performance metrics has been studied. In figure 1 miss rate is shown for different warp size and cache line size. As we can see in all cases miss rate is decreasing. Also, an interesting point is between 8 warp and 32 warp version which shows 32 warp size is getting more benefit from cache line size 64 to 128 than 8 warp size. This is because 32 warp accesses has more covering. Figure 2 shows reduction in demand fetches and therefore pressure on memory system.

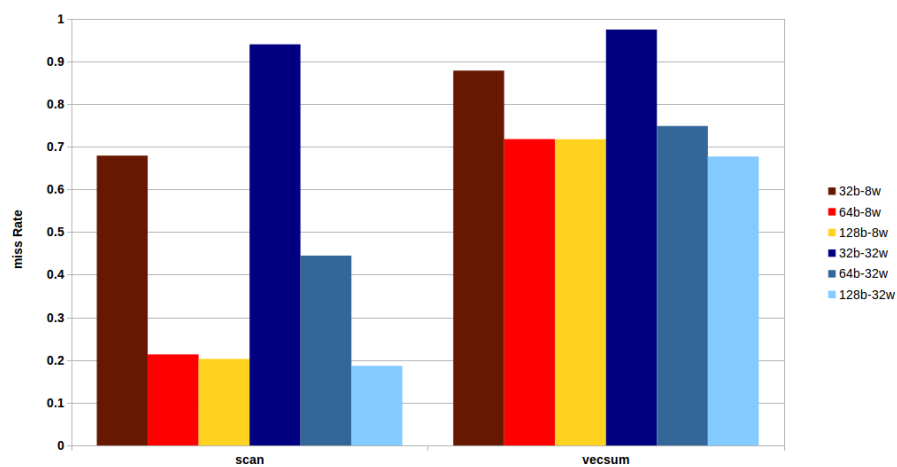


Figure 1: Miss Rate with Different Cache Line Size

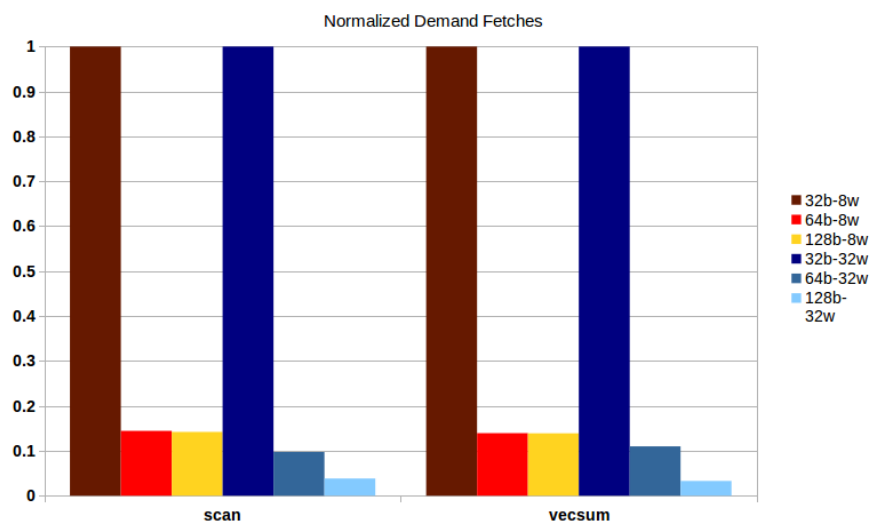


Figure 2: Normalized Demand Fetches

4 Mapping Scheme

For this section cache line size is 32 Bytes. The schemes are different in which bits they choose for coalescing. The main metric I use to select coalescing on which bits are more beneficial for an applications is the number of coalesced request versus number of non coalesced request. After a study on different bits for coalescing, for both applications I found mapping based on **bits 19 12** gives the most coalesced accesses. In figure 3 I have shown the cache performance metrics with the base line of coalescing with lower bits and no coalescing. Basic Coalescing miss rate is higher since it just reduces number of demand fetches but the ratio stays the same.

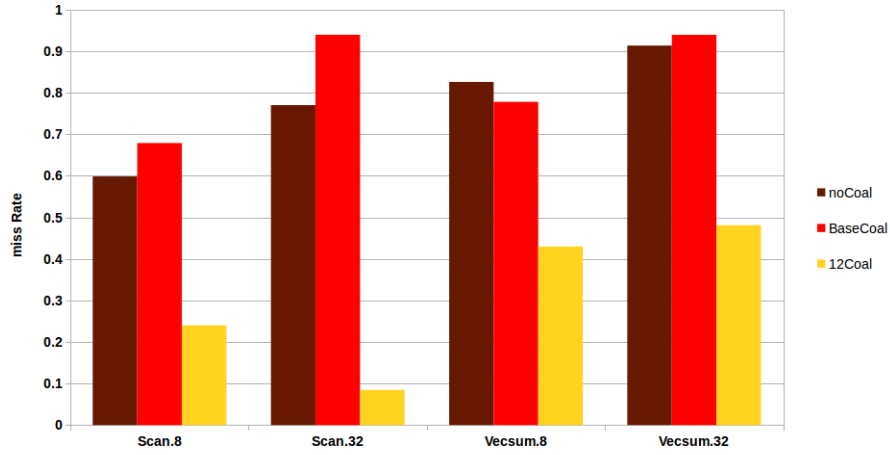


Figure 3: Mapping Schemes Miss Rates