

## ECE552 Lab 5 Data Caches

Freddy Chen 997363124 Anthony Alayo 997487401

### Question 1:

mbq1.c:

- Allocate an array that's too big for the L1 data cache so there will be misses and prefetching
- Access elements with indices of increments BSIZE (right into the next block fetched) to see the next line prefetcher work well (0.21% miss rate)
- Access elements with indices of increments 2\*BSIZE (right over the next block fetched) to see the next line prefetcher fail (53.38% miss rate)

```
1 #include <stdio.h>
2
3 // Define variables based on the configuration of the l1 data cache config
4 // <name>:<nsets>:<bsize>:<assoc>:<repl>:<pref>
5 #define NSETS 64
6 #define BSIZE 64
7
8 int main(void) {
9     int arr_size = 10000 * NSETS;
10    char arr[arr_size];
11    int trash_sum = 0; // To ensure the compiler doesn't optimize the loop away
12
13    // If increment is BSIZE*2, we will try to access an element that's just over
14    // the block that got prefetched by the nextline prefetcher, and our miss
15    // rate becomes 53.38%.
16    // If increment is BSIZE, we will access an element in the prefetched block,
17    // and our miss rate becomes 0.21%.
18    int increment = BSIZE;
19    int i;
20    for (i = 0; i < arr_size; i = i + increment) {
21        trash_sum += arr[i];
22    }
23    printf("trash_sum %d\n", trash_sum);
24    return 0;
25 }
```

### Question 2:

mbq2.c:

- Allocate an array that's too big for the L1 data cache so there will be misses and prefetching
- Access elements with indices of different strides (BSIZE or 2\*BSIZE depending on the current index) each time to see the stride prefetcher fail (60.21% miss rate)
- Access elements with indices of constant strides (BSIZE) to see the stride prefetcher work well (0.45% miss rate)

```

1 #include <stdio.h>
2
3 // Define variables based on the configuration of the l1 data cache config
4 // <name>:<nsets>:<bsize>:<assoc>:<repl>:<pref>
5 #define NSETS 64
6 #define BSIZE 64
7
8 int main(void) {
9     int arr_size = 10000 * NSETS;
10    char arr[arr_size];
11    int trash_sum = 0; // To ensure the compiler doesn't optimize the loop away
12
13    int i;
14    for (i = 0; i < arr_size; i++) {
15        trash_sum += arr[i];
16
17        // If we offset i by a different stride each time, the stride prefetcher
18        // won't do a good job predicting which block to fetch, and our miss rate
19        // becomes 60.21%.
20        // If we offset i by a constant stride, the stride prefetcher does well
21        // and our miss rate becomes 0.45%.
22        i = i % 2 ? i + BSIZE : i + 2*BSIZE;
23        //i = i + 2*BSIZE;
24    }
25    printf("trash_sum %d\n", trash_sum);
26    return 0;
27 }

```

### Question 3:

Config	L1 Miss Rate	L2 Miss Rate	Average access time *
no prefetcher	0.0416	0.1140	1.89024
next line	0.0419	0.0838	1.770122
stride	0.0385	0.0578	1.60753

\* Average access time calculated using the following formula:

$$T_{AVG} = T_{access-L1Data} + L1 \text{ miss rate} * (T_{access-L2} + L2 \text{ miss rate} * T_{hit-Memory})$$

This formula is derived in class. We always have to look in the L1 data cache, and of those accesses that miss, we always have to look the L2 cache, and of those accesses that miss, we look to memory.

#### Question 4:



The miss rate for the stride prefetcher decreases monotonically as we increase the size of the RPT for the stride prefetcher. This makes sense because we have more entries in the RPT for making stride predictions and we evict less entries as we increase our size (less aliasing). As a result, our stride prefetcher keeps its states for more instructions, making it a better prefetcher.

We think we don't see any performance gains after an RPT size of 512 because it is enough to fit the working set for the entire duration of the compress benchmark. In other words, the additional space is not being used.

#### Question 5:

It would be useful to include statistics about the average accuracy and coverage of each prefetch. The accuracy statistics will let us know whether or not our prefetcher is bringing in a lot of useless blocks. The coverage statistics will let us know how many misses we saved as a direct result of prefetching a block.

#### Question 6:

Our open prefetcher is a hybrid prefetcher that uses the stride prefetcher first to see if we can make a prefetch, and if not, we fall back to looking into a data address miss queue for any history

of the address we missed on.

Each time data access is missed in the cache, the address gets enqueued in the data address miss queue, effectively forming a history of data addresses missed.

There is a slight modification for the stride prefetcher in that we only make a stride prefetch if the entry is in steady state. Otherwise, we go straight away into looking through our history of data address misses.

If we find our address in the data address miss queue, we prefetch the next address in the queue and hope that that will be the next data address accessed.

For our open ended prefetcher, we use the following parameters (as seen in cache.h):

- entries in RPT: 1024
- spots in data address miss queue: 2048

From those parameters, we can calculate how many bits are required for those two structures:

- RPT:  $1024 \text{ entries} \times ((32 - 10 - 3) \text{ tag bits} + 32 \text{ prev-addr bits} + 32 \text{ stride bits} + 2 \text{ state bits}) = 87040 \text{ bits} = 10880 \text{ bytes} = 10.625\text{kB}$  (smaller than our 16kB L1 data cache)
- data address miss queue:  $2048 * 32 \text{ addr bits} = 65536 \text{ bits} = 8192 \text{ bytes} = 8\text{kB}$  (half the size of our 16kB L1 data cache)

Here is a comparison of our new structures with the L1 data cache (estimated with CACTI):

Structure	Access time (ns)	Cache height x width (mm):
L1 Data Cache	0.499477	0.134653 x 0.210113
RPT	0.259995	0.26398 x 0.0682748
Data address miss queue	0.271737	0.135717 x 0.0978843

From the table above, we believe these parameters make this open prefetcher feasible.

mbq6.c:

- Allocate an array that's too big for the L1 data cache so there will be misses and prefetching
- Access elements with indices of different strides (BSIZE or 2\*BSIZE depending on the current index) each time to see the stride prefetcher fail
- Based on our strides, we will access a lower bound total of  $\text{arr\_size} / \text{BSIZE}$  data addresses. Or:  
$$\text{total data address accesses} = \text{arr\_size} / \text{BSIZE}$$
- If total data address accesses < size of our miss queue, then our open prefetcher performs wonderfully with a miss rate of 0.01%
- If total data address accesses > size of our miss queue, then our open prefetcher (and

stride prefetcher) will fail and get a miss rate of 100%.

```
1 #include <stdio.h>
2
3 // Define variables based on the configuration of the l1 data cache config
4 // <name>:<nsets>:<bsize>:<assoc>:<repl>:<pref>
5 #define NSETS 64
6 #define BSIZE 64
7
8 // The open prefetcher is a hybrid of the stride prefetcher and a data address
9 // miss queue. Each time a data access is missed in the cache, it gets enqueued
10 // in the data address miss queue.
11 // The open prefetcher first checks whether or not we can do a stride prefetch.
12 // We only do a stride prefetch if the entry is in the steady state.
13 // Otherwise, we check the data address miss queue and look for the address
14 // that we missed. If we find it, we return the next address in the queue and
15 // hope that that's the next data address being accessed.
16 int main(void) {
17     // Depending on how big arr_size is, the data that we access might not
18     // fit entirely into the data address miss queue used by our open predictor.
19     // If the size is 1000 * NSETS, then all of the data that we access will fit
20     // into our miss queue ((1000 * NSETS) / BSIZE = 1000 < 2048), and our miss
21     // rate will be 0.01%!
22     // If the size is 10000 * NSETS, then the data that we access won't fit into
23     // our miss queue ((10000 * NSETS) / BSIZE = 10000 > 2048), and our miss rate
24     // will be 100%.
25     int arr_size = 1000 * NSETS;
26     char arr[arr_size];
27     int trash_sum = 0; // To ensure the compiler doesn't optimize the loop away
28
29     // We have an outer loop here so that the open prefetcher is able to take
30     // advantage of the data address miss queue
31     int i, j;
32     for (j = 0; j < 100000; j++) {
33         for (i = 0; i < arr_size; i++) {
34             trash_sum += arr[i];
35
36             // If we offset i by a different stride each time, the stride prefetcher
37             // won't do a good job predicting which block to fetch. We will fall back
38             // to looking into the data address miss queue in this case.
39             i = i % 2 ? i + BSIZE : i + 2*BSIZE;
40         }
41     }
42     printf("trash_sum %d\n", trash_sum);
43     return 0;
44 }
```

### Statement of work:

Freddy Chen - *code, microbenchmarks, lab report* Anthony Alayo - *code, simulations, lab report*