

Cache Simulator - Assignment 3 COL216

Karthik Manikandan
2023CS10298

Yash Shindekar
2023CS10592

April 2025

Contents

1	Design Decisions	2
1.1	Bus Handling	2
1.2	Cache Coherence	2
1.3	Core Timing	2
2	Implementation Details	2
2.1	Overview of MESI Protocol	2
2.2	Main Classes and Data Structures	3
2.2.1	Cache Class (Cache.cpp)	3
2.2.2	Core Class (Core.cpp)	4
2.2.3	Bus Class (Bus.cpp)	4
2.2.4	Simulator Class (Simulator.cpp)	5
2.3	Key Data Structures	5
2.4	Workflow and Key Functions	6
2.4.1	Cache::accessCache()	7
2.4.2	Cache::handleReadMiss()	7
2.4.3	Bus::busRd()	7
2.4.4	Bus::busRdX()	7
2.4.5	Cache::handleWriteMiss()	8
2.4.6	Simulator::run()	8
3	Experimental Results	8
3.1	Default Configuration Analysis	8
3.2	Analysis of Consistency Among Metrics	8
3.2.1	Consistent Metrics	8
4	Cache Parameter Analysis	8
4.1	Methodology	8
4.2	Impact of Cache Size	9
4.3	Impact of Associativity	9
4.4	Impact of Block Size	10
5	Interesting Traces	10
5.1	Conflicts	10
5.2	False Sharing	11

1 Design Decisions

1.1 Bus Handling

- If the bus is busy, the processor cannot handle snooping calls. Requests are prioritized based on the lowest core ID.
- Invalidation operations require a free bus but do not consume bus bandwidth or cycles.

1.2 Cache Coherence

- If a value is modified after being requested by another cache, a silent upgrade from **Exclusive** to **Modified** state occurs in the modifying cache.
- The value in the other cache temporarily remains the old (now stale) value, leading to a short-lived incoherence.

1.3 Core Timing

- During block transfers between caches in BusRd ,for the requestor the 2N cycles are counted as execution cycles and it halts and the halt is also for 2N cycles.
- During a writeback in BusRd of core which has MODIFIED value, that core halts for 100 cycles and not idle. Then cache to cache transfer occurs as described above.
- During a writeback in BusRdx of core which has MODIFIED value that core halts for 100 cycles and this is idle time. Then the memory fetch of the requesting core is under execution cycles of that core for a total halt of 200 cycles.
-

2 Implementation Details

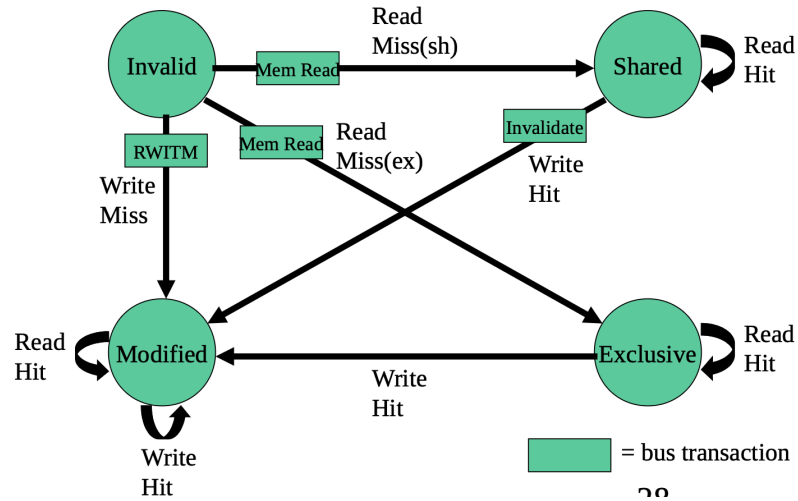
2.1 Overview of MESI Protocol

The cache simulator implements the MESI (Modified, Exclusive, Shared, Invalid) protocol, a widely used cache coherence protocol in multiprocessor systems. The protocol maintains coherence through four states:

- **Modified (M)**: The cache line is present only in the current cache and is dirty (modified).
- **Exclusive (E)**: The cache line is present only in the current cache and is clean (matches memory).
- **Shared (S)**: The cache line may be present in other caches and is clean.
- **Invalid (I)**: The cache line is invalid.

Figure 1 shows the state transitions in the MESI protocol.

MESI – locally initiated accesses



MESI – remotely initiated accesses

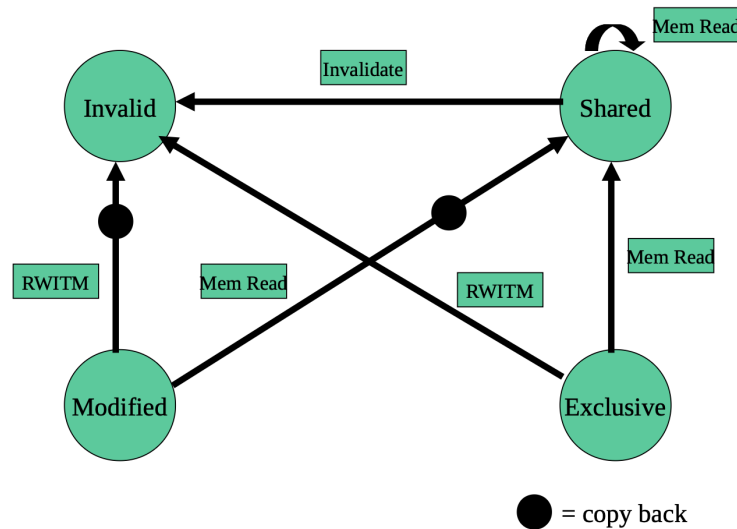


Figure 1: MESI Protocol State Transition Diagram

2.2 Main Classes and Data Structures

The simulator is built with an object-oriented approach using the following key classes:

2.2.1 Cache Class (Cache.cpp)

Manages cache lines using the LRU replacement policy and implements the MESI protocol state transitions.

```

1 class Cache {
2     int s, E, b; // Set index bits, associativity, block bits
3     std::vector<std::vector<CacheLine>> sets;
4     // Statistics counters
5     int readHits, readMisses, writeHits, writeMisses;
6     int writeBacks, idleCycles, evictions, invalidations;

```

```

7   size_t trafficBytes;
8   // Methods
9   int findLine(int setIndex, uint32_t tag);
10  int findReplacement(int setIndex, uint64_t cycle);
11  void updateLineOnHit(uint32_t setIndex, int lineIndex, uint64_t cycle);
12  void insertLine(int setIndex, int lineIndex, uint32_t tag, uint64_t cycle,
13                bool isWrite, CacheState initialState);
14  void accessCache(bool isWrite, uint32_t address, uint64_t cycle,
15                  int coreId, Bus& bus, std::vector<Core*>& cores);
16  void handleReadMiss(int coreId, uint64_t address, uint64_t cycle,
17                     Bus& bus, std::vector<Core*>& cores, uint64_t haltcycles);
18  void handleWriteMiss(int coreId, uint64_t address, uint64_t cycle,
19                      Bus& bus, std::vector<Core*>& cores, uint64_t haltcycles);
20  void busupdate(class Bus &bus);
21 };

```

The CacheLine structure represents individual cache lines with their state and metadata:

```

1 struct CacheLine {
2     bool valid = false;
3     uint32_t tag = 0;
4     uint64_t lastUsedCycle = 0;
5     CacheState state = INVALID;
6 };

```

2.2.2 Core Class (Core.cpp)

Represents processor cores with their associated caches and memory access traces.

```

1 class Core {
2 public:
3     int id;
4     Cache* cache;
5     std::vector<Request> trace;
6     size_t instPtr;
7     size_t prevInstr;
8     uint64_t nextFreeCycle;
9     // Execution statistics
10    int readCount, writeCount, execycles;
11
12    Core(int id, Cache* cache);
13    void loadTrace(const std::string& filename);
14 };

```

Each core loads its memory access trace from a file, which contains a sequence of read and write operations:

```

1 struct Request {
2     bool isWrite;
3     uint32_t address;
4
5     Request(bool isWrite, uint32_t address);
6 };

```

2.2.3 Bus Class (Bus.cpp)

Handles cache coherence transactions between cores, implementing the communication protocols required by MESI.

```

1 class Cache {
2 public:
3     int s, E, b;
4     // LRU list to track usage order per set
5     typedef std::list<CacheKey> LRUList;
6
7     // Cache map: stores cache lines and references to their position in the LRU list
8     typedef std::pair<CacheLine, typename LRUList::iterator> CacheEntry;
9     typedef std::unordered_map<CacheKey, CacheEntry> CacheMap;
10
11    // One LRU list and one map per set
12    std::vector<LRUList> lruLists;
13    std::vector<CacheMap> cacheMaps;

```

```

14
15 // Statistics
16 uint64_t readHits;
17 uint64_t readMisses;
18 uint64_t writeHits;
19 uint64_t writeMisses;
20 uint64_t writeBacks;
21 uint64_t idleCycles;
22 uint64_t evictions;
23 uint64_t trafficBytes;
24 uint64_t invalidations;
25
26 Cache(int s, int E, int b);
27
28 // Core cache operations
29 void accessCache(bool isWrite, uint32_t address, uint64_t cycle, int coreId,
30                 class Bus& bus, std::vector<class Core*>& cores);
31
32 // Map-based cache operations
33 CacheLine* findLine(int setIndex, uint32_t tag);
34 std::pair<CacheKey, CacheLine*> findReplacement(int setIndex, uint64_t cycle);
35 void updateLRU(int setIndex, uint32_t tag, uint64_t cycle);
36 void insertLine(int setIndex, uint32_t tag, uint64_t cycle, bool isWrite, CacheState
    initialState);
37
38 // Bus and miss handling operations
39 void busupdate(class Bus& bus);
40 void handleReadMiss(int coreId, uint64_t address, uint64_t cycle, Bus& bus,
41                     std::vector<Core*>& cores, uint64_t haltcycles);
42 void handleWriteMiss(int coreId, uint64_t address, uint64_t cycle, Bus& bus,
43                      std::vector<Core*>& cores, uint64_t haltcycles);
44 };

```

2.2.4 Simulator Class (Simulator.cpp)

Coordinates the entire simulation, managing cores, the bus, and tracking global simulation state.

```

1 class Simulator {
2     std::vector<Core*> cores;
3     Bus bus;
4     uint64_t globalCycle;
5
6     void initialize(int numCores, int s, int E, int b);
7     void loadTraces(const std::vector<std::string>& traceFiles);
8     void run();
9     void printStats();
10    uint64_t computeMaxExecutionTime();
11 };

```

2.3 Key Data Structures

The simulator employs several critical data structures:

1. **Set-Associative Cache:** Implemented using an unordered Map.
2. **Memory Access Trace:** A vector of Request objects containing memory addresses and operation types (read/write).
3. **MESI State:** An enumeration representing the four possible states of a cache line.

2.4 Workflow and Key Functions

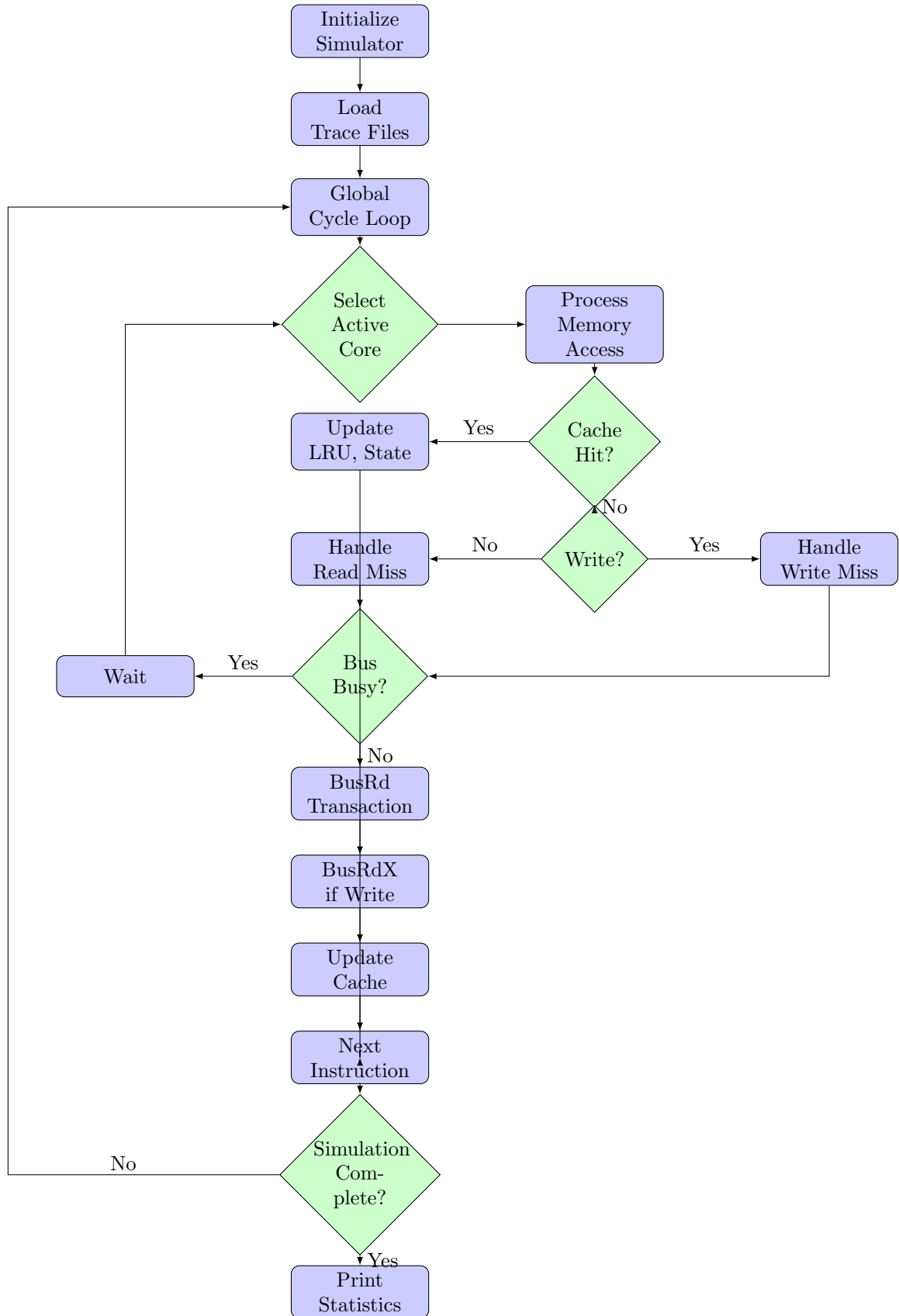


Figure 2: Simulation Workflow Diagram

The simulator follows the workflow shown in Figure 2, with several key functions handling the core operations:

2.4.1 **Cache::accessCache()**

This function handles both read and write operations to the cache:

1. Calculate set index and tag from the memory address
2. Check if the line is present in the cache (hit or miss)
3. For cache hits:
 - Update LRU information
 - For writes, handle coherence actions based on MESI state (upgrade to MODIFIED if needed)
 - For reads, simply update statistics
4. For cache misses:
 - Call `handleReadMiss()` or `handleWriteMiss()` based on operation

2.4.2 **Cache::handleReadMiss()**

This function handles read miss operations following the MESI protocol:

1. Check if the bus is busy; if so, increment idle cycles and wait
2. Find a replacement line using LRU policy
3. If the victim line is in MODIFIED state, write it back to memory
4. Perform a `busRd` operation to check for copies in other caches
5. Update the line's state based on the bus result:
 - SHARED if another cache has a copy
 - EXCLUSIVE if no other cache has a copy
6. Update statistics and proceed to the next instruction

2.4.3 **Bus::busRd()**

This function handles BusRead operations for cache coherence:

1. Increment bus transaction counter
2. Extract set index and tag from the address
3. Check if any other core has the line and in what state
4. Update the state in other caches if needed
5. Return the result indicating if data was found in other caches

2.4.4 **Bus::busRdX()**

This function handles BusReadExclusive operations:

1. Similar to `busRd()`, but also invalidates copies in other caches
2. Returns the result indicating if data was found in other caches

2.4.5 Cache::handleWriteMiss()

This function handles write miss operations:

1. Similar to read miss handling but:
 - Performs `busRdX` instead of `busRd`
 - Sets the final state to `MODIFIED`
 - Invalidates copies in other caches

2.4.6 Simulator::run()

This is the main simulation loop that:

1. Increments the global cycle counter
2. Selects the next core to execute based on availability
3. Processes the current instruction from the trace
4. Updates simulation statistics
5. Continues until all traces are processed

3 Experimental Results

3.1 Default Configuration Analysis

We analyze the simulator with the default parameters:

- Cache size: 4KB per processor
- Associativity: 2-way set associative
- Block size: 32 bytes

3.2 Analysis of Consistency Among Metrics

3.2.1 Consistent Metrics

In our simulator, when multiple cores attempt to place transactions on the shared bus in the same cycle, priority is deterministically given to the lower-numbered core. This fixed arbitration scheme eliminates non-determinism in bus access order across simulation runs. As a result of this deterministic bus arbitration, all simulator outputs (e.g., cache miss rates, execution cycles, idle cycles, and bus traffic) remain identical across all 10 runs for the same application trace. No variation is observed in any metric, as the simulation behavior is fully repeatable under fixed input and parameters.

4 Cache Parameter Analysis

4.1 Methodology

To understand the impact of cache parameters on performance, we varied one parameter at a time while keeping the others at their default values:

- **Cache size:** 2KB, 4KB (default), 8KB, 16KB
- **Associativity:** 1-way (direct mapped), 2-way (default), 4-way, 8-way
- **Block size:** 16B, 32B (default), 64B, 128B

For each configuration, we computed the maximum execution time across all cores, representing the critical path that determines overall program completion time.

4.2 Impact of Cache Size

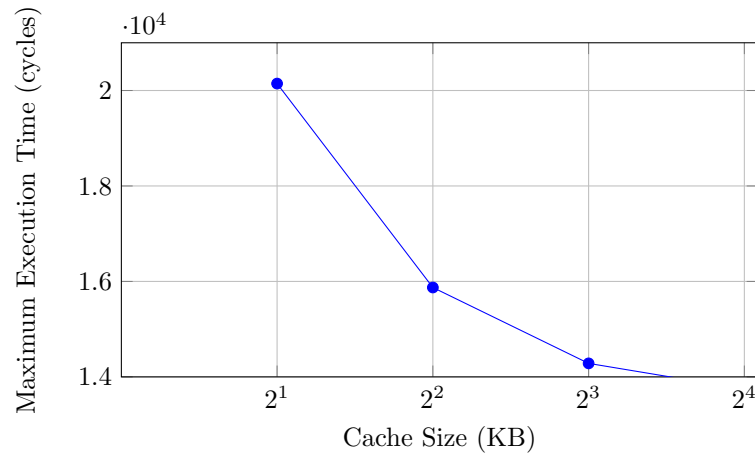


Figure 3: Maximum execution time vs. cache size

Observations:

- Increasing cache size significantly reduces execution time due to fewer conflict misses.
- The performance improvement follows a diminishing returns curve, with the largest gain seen when doubling from 2KB to 4KB.
- The improvement from 8KB to 16KB is less, suggesting that 8KB is close to the working set size of the application.
- Larger caches reduce bus contention by decreasing the frequency of cache line evictions and the resulting coherence traffic and so it plateaus at 16KB

4.3 Impact of Associativity

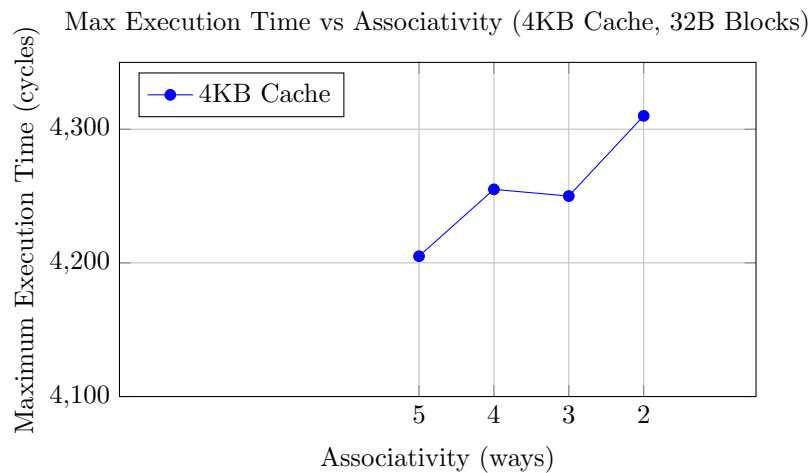


Figure 4: Maximum execution time vs. associativity

Observations:

- Moving from direct-mapped (1-way) to 2-way set associative yields a substantial increase in execution time due to more lines to compare in each set but however miss rate is improved.

4.4 Impact of Block Size

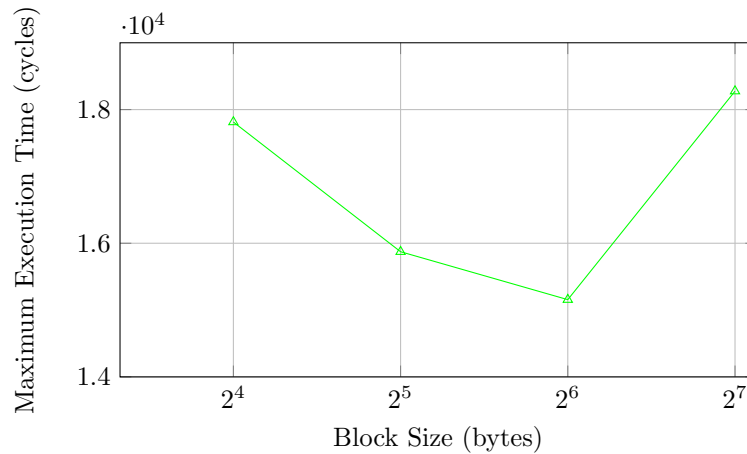


Figure 5: Maximum execution time vs. block size

Observations:

- Block size shows a U-shaped performance curve, with the best performance at 64 bytes.
- Small blocks (16B) lead to poor spatial locality utilization, causing more misses and increased execution time.
- Medium blocks (32B-64B) provide a good balance, with 64B offering optimal spatial locality for this workload.
- Large blocks (128B) actually degrade performance by:
 - Increasing miss penalty (more data to transfer)
 - Causing false sharing in multicore scenarios
 - Reducing the effective cache capacity by storing unnecessary data
 - Increasing bus traffic for coherence operations
- The performance degradation when going from 64B to 128B highlights the importance of selecting appropriate block sizes for multicore workloads.

5 Interesting Traces

5.1 Conflicts

```
1 W 0x1000
2 W 0x1004
3 W 0x1008
4 W 0x2000
5 W 0x3000
6
7
8 Simulation Parameters:
9 Trace Prefix: bonus_tc/conflict
10 Set Index Bits: 6
11 Associativity: 2
12 Block Bits: 5
13 Block Size (Bytes): 32
14 Number of Sets: 64
15 Cache Size (KB per core): 4
16 MESI Protocol: Enabled
17 Write Policy: Write-back, Write-allocate
18 Replacement Policy: LRU
```

```

19 Bus: Central snooping bus
20
21 Core 0 Statistics:
22 Total Instructions: 5
23 Total Reads: 0
24 Total Writes: 5
25 Total Execution Cycles: 500
26 Idle Cycles: 202
27 Cache Misses: 5
28 Cache Miss Rate: 100.00%
29 Cache Evictions: 0
30 Writebacks: 5
31 Bus Invalidations: 2
32 Data Traffic (Bytes): 320
33
34 Core 1 Statistics:
35 Total Instructions: 5
36 Total Reads: 0
37 Total Writes: 5
38 Total Execution Cycles: 500
39 Idle Cycles: 504
40 Cache Misses: 5
41 Cache Miss Rate: 100.00%
42 Cache Evictions: 0
43 Writebacks: 5
44 Bus Invalidations: 5
45 Data Traffic (Bytes): 320
46
47 Core 2 Statistics:
48 Total Instructions: 5
49 Total Reads: 0
50 Total Writes: 5
51 Total Execution Cycles: 502
52 Idle Cycles: 908
53 Cache Misses: 3
54 Cache Miss Rate: 60.00%
55 Cache Evictions: 0
56 Writebacks: 5
57 Bus Invalidations: 3
58 Data Traffic (Bytes): 256
59
60 Core 3 Statistics:
61 Total Instructions: 5
62 Total Reads: 0
63 Total Writes: 5
64 Total Execution Cycles: 602
65 Idle Cycles: 1412
66 Cache Misses: 3
67 Cache Miss Rate: 60.00%
68 Cache Evictions: 1
69 Writebacks: 3
70 Bus Invalidations: 3
71 Data Traffic (Bytes): 192
72
73 Overall Bus Summary:
74 Total Bus Transactions: 16
75 Total Bus Traffic (Bytes): 576

```

Having this in all the 4 processor input files produces conflict misses as shown in the output

5.2 False Sharing

```

1 For proc0
2 W 0x1000
3 R 0x1000
4
5 For proc1
6 W 0x1001
7 R 0x1001
8
9 For proc2
10 W 0x1002

```

```

11 R 0x1002
12
13 For proc3
14 W 0x1003
15 R 0x1003
16
17 Output is
18
19 Simulation Parameters:
20 Trace Prefix: bonus_tc/app4
21 Set Index Bits: 6
22 Associativity: 2
23 Block Bits: 5
24 Block Size (Bytes): 32
25 Number of Sets: 64
26 Cache Size (KB per core): 4
27 MESI Protocol: Enabled
28 Write Policy: Write-back, Write-allocate
29 Replacement Policy: LRU
30 Bus: Central snooping bus
31
32 Core 0 Statistics:
33 Total Instructions: 1
34 Total Reads: 0
35 Total Writes: 1
36 Total Execution Cycles: 100
37 Idle Cycles: 0
38 Cache Misses: 1
39 Cache Miss Rate: 100.00%
40 Cache Evictions: 0
41 Writebacks: 1
42 Bus Invalidations: 0
43 Data Traffic (Bytes): 64
44
45 Core 1 Statistics:
46 Total Instructions: 1
47 Total Reads: 0
48 Total Writes: 1
49 Total Execution Cycles: 100
50 Idle Cycles: 100
51 Cache Misses: 1
52 Cache Miss Rate: 100.00%
53 Cache Evictions: 0
54 Writebacks: 1
55 Bus Invalidations: 1
56 Data Traffic (Bytes): 64
57
58 Core 2 Statistics:
59 Total Instructions: 1
60 Total Reads: 0
61 Total Writes: 1
62 Total Execution Cycles: 100
63 Idle Cycles: 201
64 Cache Misses: 1
65 Cache Miss Rate: 100.00%
66 Cache Evictions: 0
67 Writebacks: 1
68 Bus Invalidations: 1
69 Data Traffic (Bytes): 64
70
71 Core 3 Statistics:
72 Total Instructions: 1
73 Total Reads: 0
74 Total Writes: 1
75 Total Execution Cycles: 100
76 Idle Cycles: 302
77 Cache Misses: 1
78 Cache Miss Rate: 100.00%
79 Cache Evictions: 0
80 Writebacks: 0
81 Bus Invalidations: 1
82 Data Traffic (Bytes): 32
83

```

```
84 Overall Bus Summary:  
85 Total Bus Transactions: 4  
86 Total Bus Traffic (Bytes): 96
```

We can see that each core after the initial compulsory miss produces a invalidation to invalidate the other copies of the same cache line and hence can never have a cache hit even though true sharing as above testcase doesnt occur here