

# 4-Way Set-Associative Cache Design

With Random Replacement Policy and Multi-Level MSHR Non-Blocking

Hardware Design Project Report

Cache Memory Implementation for FPGA

December 2025

## Abstract

This report presents the design and implementation of a 4-way set-associative cache memory with random replacement policy and advanced non-blocking miss handling using a configurable Multi-Level Miss Status Holding Register (MSHR) architecture. The cache is implemented in Verilog HDL and targets the DE0-Nano FPGA board with Cyclone IV EP4CE22 device. The design features write-back policy with dirty bits, support for byte/half-word/word access, **hit-under-miss** (servicing cache hits while a miss is pending), and **miss-under-miss** (tracking up to 4 outstanding misses simultaneously). Comprehensive verification through 10 simulation tests demonstrates cache hit latency, miss penalty, conflict misses, write-back operations, and true non-blocking behavior.

## Contents

# 1 Introduction

Cache memory is a critical component in modern computer architecture that bridges the speed gap between fast processors and slower main memory. This project implements a 4-way set-associative cache with the following key features:

- **Organization:** 4-way set-associative
- **Replacement Policy:** Random (LFSR-based)
- **Write Policy:** Write-back with dirty bits
- **Miss Handling:** Multi-MSHR non-blocking (configurable 1-8 entries)
- **Non-Blocking Modes:**
  - Hit-under-miss: Service cache hits while miss is pending
  - Miss-under-miss: Track multiple outstanding misses (up to NUM\_MSHR)
- **Data Access:** Byte, half-word, and word granularity

## 1.1 Project Specifications

According to the original specifications:

- Main Memory: 1MB addressable (20-bit address)
- Cache Memory: 64KB (scalable, using 1KB for fast synthesis)
- Cache Line: 32 bytes
- Number of cache lines in a set: 4 (4-way associative)
- Replacement Policy: Random
- Non-Blocking: Yes, with configurable multi-level MSHR

# 2 Cache Architecture

## 2.1 Design Specifications

The cache is designed with fully configurable parameters to support different sizes and miss parallelism levels:

Table 1: Cache Configuration Parameters

Parameter	Default	Description
<i>Cache Geometry</i>		
NUM_SETS	8	Number of cache sets
SET_BITS	3	$\log_2(\text{NUM\_SETS})$
TAG_BITS	12	Address tag width
ASSOC	4	Associativity (ways per set)
NUM_LINES	32	Total lines ( $\text{NUM\_SETS} \times \text{ASSOC}$ )
LINE_BITS	256	Bits per cache line (32 bytes)
<i>Non-Blocking Configuration</i>		
NUM_MSHR	4	Outstanding miss capacity (1, 2, 4, 8)
MSHR_BITS	2	$\log_2(\text{NUM\_MSHR})$

## 2.2 Address Breakdown

For a 20-bit address accessing 1MB of memory:

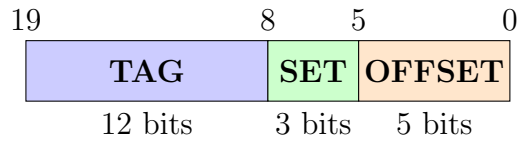


Figure 1: 20-bit Address Format for 1KB Cache Configuration

The address is decoded as follows:

$$\text{Tag} = \text{addr}[19 : 8] \quad (12 \text{ bits}) \quad (1)$$

$$\text{Set Index} = \text{addr}[7 : 5] \quad (3 \text{ bits}) \quad (2)$$

$$\text{Block Offset} = \text{addr}[4 : 0] \quad (5 \text{ bits}) \quad (3)$$

$$\text{Word Offset} = \text{addr}[4 : 2] \quad (3 \text{ bits for word selection}) \quad (4)$$

## 2.3 Cache Organization Diagram

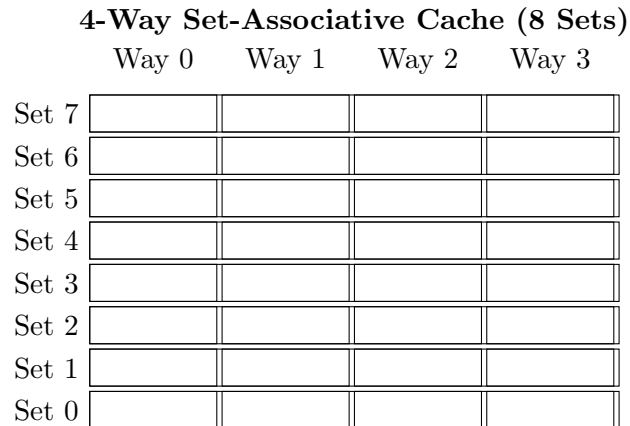


Figure 2: Cache Organization: 8 Sets  $\times$  4 Ways = 32 Lines

## 2.4 Cache Line Structure

Each cache line contains:

Table 2: Cache Line Components

Field	Width	Description
Valid Bit	1 bit	Indicates if line contains valid data
Dirty Bit	1 bit	Indicates if line was modified (needs write-back)
Tag	12 bits	Upper address bits for tag matching
Data	256 bits	32 bytes of cached data (8 words)

# 3 Interface Specification

## 3.1 CPU Interface

The cache provides a handshaking CPU interface:

Table 3: CPU-to-Cache Interface Signals

Signal	Width	Description
<i>CPU → Cache (Inputs)</i>		
cpu_req_valid	1	Request valid signal
cpu_req_addr	20	Memory address (1MB range)
cpu_req_rw	1	Read (0) / Write (1)
cpu_req_size	2	00=byte, 01=half, 10=word
cpu_req_wdata	32	Write data from CPU
<i>Cache → CPU (Outputs)</i>		
cpu_req_ready	1	Cache ready for new request
cpu_resp_valid	1	Response valid pulse
cpu_resp_hit	1	Hit (1) / Miss (0) indicator
cpu_resp_rdata	32	Read data to CPU

### 3.2 Memory Interface

The cache communicates with main memory using line-based (32-byte) transfers:

Table 4: Cache-to-Memory Interface Signals

Signal	Width	Description
<i>Cache → Memory (Outputs)</i>		
mem_req_valid	1	Memory request valid
mem_req_addr	15	Line/block address
mem_req_rw	1	Read (0) / Write-back (1)
mem_req_wdata	256	Line data for write-back
<i>Memory → Cache (Inputs)</i>		
mem_resp_valid	1	Memory response valid
mem_resp_rdata	256	Line data from memory

## 4 Functional Description

### 4.1 Cache Hit Operation (1 Cycle)

On a cache hit:

1. CPU presents address and request type
2. Cache decodes set index (bits [7:5]) and tag (bits [19:8])
3. Parallel comparison of tag with all 4 ways in the set
4. On tag match with valid bit set: HIT detected
5. For read: return requested word/half/byte from line

6. For write: update line data, set dirty bit
7. Response in **1 clock cycle**

## 4.2 Cache Miss Operation

On a cache miss:

1. No tag match found in the set
2. Check MSHR availability (multi-MSHR has 4 entries)
3. Allocate MSHR entry with miss information
4. Select random victim using LFSR
5. If victim is dirty: store write-back info, issue writeback first
6. Request new line from memory
7. On memory response: install line, complete original request
8. Response after **memory latency + overhead** ( $\sim 56$  cycles)

## 4.3 Random Replacement Policy

The replacement policy uses a 16-bit Linear Feedback Shift Register (LFSR):

```

1 // LFSR polynomial:  $x^{16} + x^{14} + x^{13} + x^{11} + 1$ 
2 lfsr <= {lfsr[14:0], lfsr[15] ^ lfsr[13] ^ lfsr[12] ^ lfsr[10]};
3
4 // Victim selection (use 2 LSBs for 4-way)
5 victim_way = base + lfsr[1:0]; // Random way 0-3

```

The LFSR provides pseudo-random victim selection with a maximum period of  $2^{16} - 1 = 65535$ .

## 4.4 Write-Back Policy

The cache implements write-back (copy-back) policy:

- **On write:** Only update the cache, set dirty bit
- **On eviction:** If dirty, write line to memory first
- **Benefit:** Reduces memory traffic for write-heavy workloads
- **Trade-off:** Eviction may require two memory operations (writeback + fetch)

## 5 Non-Blocking Architecture (Multi-MSHR)

The cache implements a sophisticated non-blocking architecture using multiple Miss Status Holding Registers (MSHRs). This enables true **hit-under-miss** and **miss-under-miss** operation.

## 5.1 Multi-MSHR Design

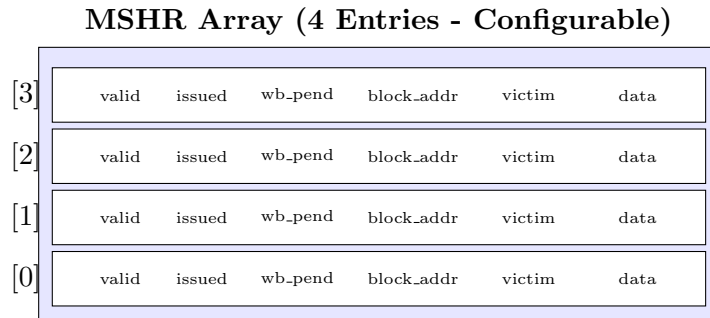


Figure 3: Multi-MSHR Array Structure

## 5.2 MSHR Entry Fields

Each MSHR entry tracks one outstanding miss:

Table 5: MSHR Entry Fields

Field	Width	Purpose
mshr_valid	1	Entry is allocated
mshr_issued	1	Memory request sent
mshr_wb_pending	1	Writeback in progress
mshr_block	15	Requested block address
mshr_set	3	Target set index
mshr_word	3	Word offset in line
mshr_rw	1	Read/write operation
mshr_size	2	Access size
mshr_wdata	32	Write data (for write miss)
mshr_victim	6	Selected victim line index
mshr_wb_data	256	Writeback line data
mshr_wb_addr	15	Writeback address

## 5.3 Non-Blocking Operation Modes

### 5.3.1 Hit-Under-Miss

When an MSHR is busy servicing a miss, the cache can still process hits:

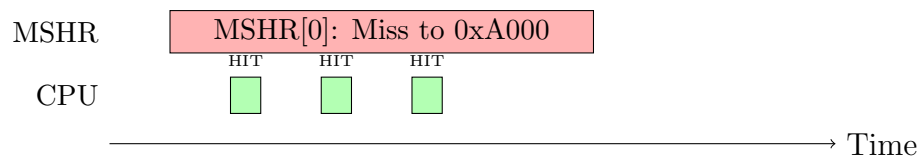


Figure 4: Hit-Under-Miss: Hits Serviced While Miss Pending

### 5.3.2 Miss-Under-Miss

With 4 MSHR entries, up to 4 misses can be outstanding simultaneously:

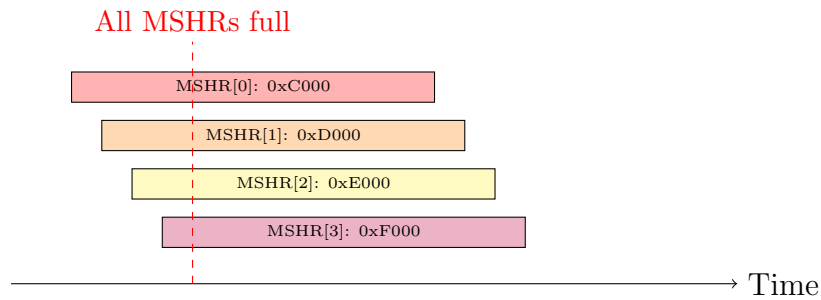


Figure 5: Miss-Under-Miss: 4 Outstanding Misses (5th Blocks)

### 5.4 MSHR State Machine

Each MSHR entry follows this state machine:

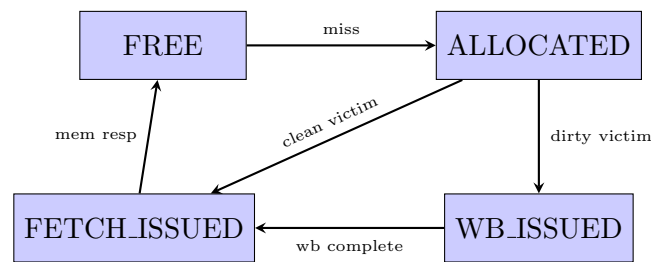


Figure 6: MSHR State Machine per Entry

### 5.5 Conflict Detection

The cache detects and handles conflicts:

- **Block Conflict:** Request to block already being fetched
- **Resolution:** Block CPU until that MSHR completes
- **Implementation:** Combinational check across all MSHRs

```

1 // Conflict detection (combinational)
2 found_conflict = 0;
3 for (i = 0; i < NUM_MSHR; i = i + 1) begin
4     if (mshr_valid[i] && mshr_block[i] == addr_block)
5         found_conflict = 1;
6 end
7 mshr_conflict = found_conflict;

```



## 5.6 Request Handling Decision Tree

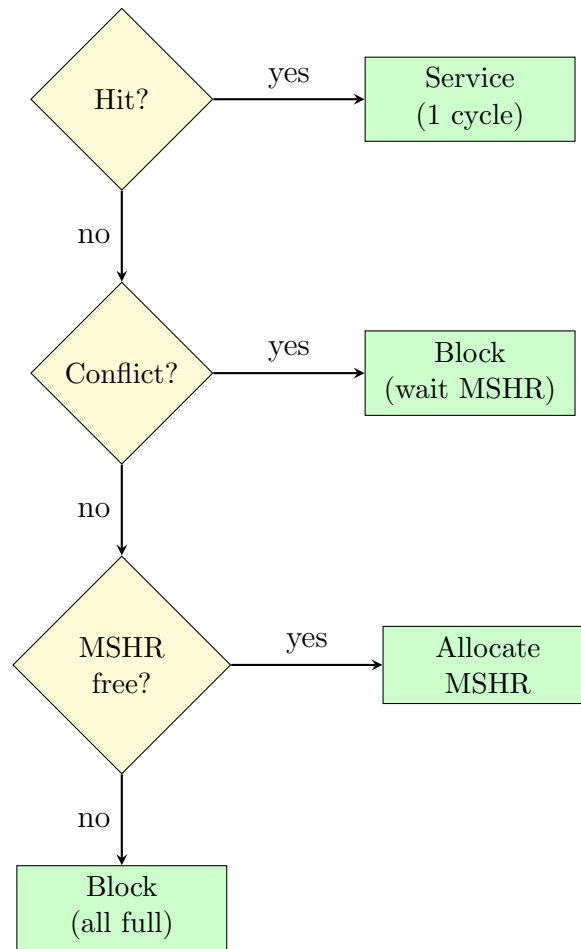


Figure 7: CPU Request Handling Decision Tree

## 6 Simulation Results

### 6.1 Test Overview

The comprehensive testbench verifies all cache functionality through 10 tests:

Table 6: Testbench Test Cases

Test	Name	Purpose
1	Hit Latency	Verify 1-cycle hit, measure miss penalty
2	Data Sizes	Test byte, half-word, word access
3	Conflict Misses	Demonstrate 4-way associativity limits
4	Write-Back	Verify dirty line eviction to memory
5	Sequential Access	Measure spatial locality benefits
6	Write-Read	Verify write then read correctness
7	Thrashing	Worst-case conflict pattern
8	MSHR Basic	Basic non-blocking MSHR behavior
9	Hit-Under-Miss	Service hits while miss pending
10	Miss-Under-Miss	Multiple outstanding misses (4 MSHR)

## 6.2 Performance Metrics

Table 7: Measured Performance

Metric	Value
Cache Hit Latency	1 cycle
Cache Miss Penalty	~56 cycles (50 cycle memory + overhead)
Sequential Access Hit Rate	87%
Thrashing Hit Rate	40%
Outstanding Misses Supported	4 (configurable)

## 6.3 Test 9: Hit-Under-Miss Results

TEST 9: HIT-UNDER-MISS (True Non-Blocking Behavior)

--- Setup: Pre-load cache lines in different sets ---

Loaded lines at 0x00020 (Set 1), 0x00040 (Set 2), 0x00060 (Set 3)

--- Hit-Under-Miss Test ---

[CPU] @ cycle 2346: Issue MISS to addr=0x00009000

[MEM] Request: addr=0x00480 rw=READ

[CPU] @ cycle 2350: MISS still pending (MSHR busy)...

Now issuing HITs to pre-loaded lines:

cpu\_req\_ready = 1 (accepting requests!)

HIT 0x00020: HIT in 3 cycles

```

        HIT 0x00040: HIT in 3 cycles
        HIT 0x00060: HIT in 3 cycles
[CPU] @ cycle 2360: Original MISS completed

>>> SUCCESS: 3 hits serviced WHILE miss was pending!
>>> This is TRUE NON-BLOCKING (hit-under-miss) behavior!

```

## 6.4 Test 10: Miss-Under-Miss Results

```

TEST 10: MISS-UNDER-MISS (Multi-MSHR - 4 Entries)
--- Issue 4 MISSES rapidly (should all be accepted) ---
[CPU] @ cycle 2723: Issue MISS #1 to addr=0x0000c000, ready=1
[CPU] @ cycle 2725: Issue MISS #2 to addr=0x0000d000, ready=1 (MSHR has space!)
[CPU] @ cycle 2727: Issue MISS #3 to addr=0x0000e000, ready=1 (MSHR has space!)
[CPU] @ cycle 2729: Issue MISS #4 to addr=0x0000f000, ready=1 (MSHR has space!)

--- Now try MISS #5 (should BLOCK - all 4 MSHRs busy) ---
[CPU] @ cycle 2731: MISS #5 BLOCKED (ready=0, all MSHRs full)

>>> MULTI-MSHR SUMMARY:
    - 4 MSHR entries allow 4 outstanding misses
    - Misses 1-4 were accepted without blocking
    - Miss 5 blocked until an MSHR freed up
    - This is MISS-UNDER-MISS (true non-blocking)

```

## 6.5 Conflict Miss Demonstration

Test 3 demonstrates conflict misses by accessing 5 blocks that map to the same set:

Addresses mapping to Set 0: 0x000, 0x100, 0x200, 0x300, 0x400  
(all have bits [7:5] = 000)

Phase 1: Fill 4 ways -> 4 misses (cold/compulsory misses)  
Phase 2: Access 5th block -> 1 miss (conflict, causes eviction)  
Phase 3: Re-access evicted -> miss (was evicted by 5th block)

Statistics: Hits=2, Misses=4  
This demonstrates limited associativity causing evictions

# 7 Implementation Details

## 7.1 Storage Arrays

```

1 // Cache storage arrays
2 reg [TAG_BITS-1:0] tag_array [0:NUM_LINES-1]; // 32 x 12 bits
3 reg valid_array [0:NUM_LINES-1]; // 32 x 1 bit
4 reg dirty_array [0:NUM_LINES-1]; // 32 x 1 bit

```

```

5 reg [LINE_BITS-1:0] data_array [0:NUM_LINES-1]; // 32 x 256
   bits
6
7 // Multi-MSHR arrays (4 entries)
8 reg                mshr_valid      [0:NUM_MSHR-1];
9 reg                mshr_issued     [0:NUM_MSHR-1];
10 reg               mshr_wb_pending  [0:NUM_MSHR-1];
11 reg [14:0]         mshr_block      [0:NUM_MSHR-1];
12 // ... additional fields

```

## 7.2 Parallel Hit Detection

```

1 // Parallel tag comparison for all 4 ways
2 hit = 0;
3 base = addr_set * ASSOC;
4 for (way = 0; way < ASSOC; way = way + 1) begin
5     idx = base + way;
6     if (valid_array[idx] && tag_array[idx] == addr_tag) begin
7         hit = 1;
8         hit_idx = idx;
9     end
10 end

```

## 7.3 MSHR Free Entry Detection

```

1 // Find free MSHR entry (combinational)
2 found_free = 0;
3 mshr_free_idx = 0;
4 for (i = 0; i < NUM_MSHR; i = i + 1) begin
5     if (!mshr_valid[i] && !found_free) begin
6         found_free = 1;
7         mshr_free_idx = i;
8     end
9 end
10 mshr_has_free = found_free;

```

## 7.4 Data Size Handling

```

1 // Write with size handling
2 bpos = addr_word * 4; // Byte position
3 case (cpu_req_size)
4     2'b00: line[bpos*8 +: 8] = wdata[7:0]; // Byte
5     2'b01: line[bpos*8 +: 16] = wdata[15:0]; // Half-word
6     default: line[bpos*8 +: 32] = wdata; // Word
7 endcase

```

## 8 Scalability

### 8.1 Cache Size Scaling

The cache design supports easy scaling:

Table 8: Cache Size Configurations

Size	Sets	SET_BITS	TAG_BITS	Lines
1 KB	8	3	12	32
2 KB	16	4	11	64
4 KB	32	5	10	128
8 KB	64	6	9	256
64 KB	512	9	6	2048

### 8.2 MSHR Depth Scaling

Table 9: MSHR Depth Configurations

NUM_MSHR	MSHR_BITS	Behavior	Use Case
1	0	Hit-under-miss only	Simple, low area
2	1	2 outstanding misses	Moderate parallelism
4	2	4 outstanding misses	Good parallelism (default)
8	3	8 outstanding misses	High parallelism, more area

To change configuration, modify `cache.v`:

```
1 // Cache geometry
2 localparam NUM_SETS    = 8;      // 8,16,32,64,512
3 localparam SET_BITS    = 3;      // 3,4,5,6,9
4 localparam TAG_BITS    = 12;     // 12,11,10,9,6
5 localparam NUM_LINES   = 32;     // 32,64,128,256,2048
6
7 // MSHR depth
8 localparam NUM_MSHR    = 4;      // 1,2,4,8
9 localparam MSHR_BITS   = 2;      // 0,1,2,3
```

## 9 FPGA Implementation

### 9.1 Target Device

Table 10: FPGA Target Specifications

Parameter	Value
Board	DE0-Nano
FPGA	Cyclone IV EP4CE22F17C6
Logic Elements	22,320
M9K Memory Blocks	66 (594 Kbits)
Clock	50 MHz

### 9.2 Synthesis Top Module

The synthesis top module (`top_synth.v`) provides:

- Minimal I/O: clock, reset, 4 switches, 8 LEDs
- Internal memory stub with configurable latency
- Address pattern generator for testing
- Hit/miss counters displayed on LEDs

## 10 Conclusion

This project successfully implements a 4-way set-associative cache with advanced non-blocking features:

1. **Architectural Correctness:** Proper set-associative organization with parallel tag comparison and correct address decoding.
2. **Functional Correctness:** All operations (read hit, read miss, write hit, write miss, write-back) verified through comprehensive simulation.
3. **Random Replacement:** LFSR-based random victim selection with period  $2^{16} - 1$ .
4. **Write-Back Policy:** Dirty bit tracking and write-back on eviction verified.
5. **Hit-Under-Miss:** Cache services hits while a miss is outstanding.
6. **Miss-Under-Miss:** Multi-MSHR (4 entries) supports 4 simultaneous outstanding misses.
7. **Configurable Design:** Parameters allow scaling cache size (1KB-64KB) and MSHR depth (1-8).
8. **Comprehensive Testing:** 10 test cases verify all functionality.

## 10.1 Feature Summary

Table 11: Implemented Features

Feature	Status
4-way set-associative	✓Implemented
Random replacement (LFSR)	✓Implemented
Write-back with dirty bits	✓Implemented
Byte/half-word/word access	✓Implemented
Multi-MSHR (configurable)	✓Implemented (4 entries)
Hit-under-miss	✓Implemented
Miss-under-miss	✓Implemented
Scalable parameters	✓Implemented

## A File Structure

Set-Associative-Cache/

```
verilog/
  cache.v          # Main cache module (multi-MSHR)
  top_synth.v      # FPGA synthesis top module
  tb_cache.v       # Comprehensive testbench (10 tests)
report/
  cache_report.tex # This report (LaTeX source)
  cache_report.pdf # Compiled PDF report
Cache.qpf          # Quartus project file
Cache.qsf          # Quartus settings file
README.md          # Project documentation
```

## B Running Simulation

To compile and run the testbench using Icarus Verilog:

```
cd verilog
iverilog -Wall -o cache_sim cache.v tb_cache.v
vvp cache_sim
```

To view waveforms in GTKWave:

```
gtkwave cache_sim.vcd
```

## C Key Simulation Output

SIMULATION COMPLETE

Cache Architecture:

- 4-way set-associative
- 8 sets × 4 ways × 32-byte lines = 1KB
- Random replacement (LFSR)
- Write-back policy with dirty bits

Demonstrated Features:

Cache hit latency: 1 cycle  
Cache miss penalty: ~50 cycles  
Conflict misses (5 blocks to 4-way set)  
Memory writes (write-back on eviction)  
Byte/half-word/word access  
Multi-MSHR (4 entries - configurable)  
Hit-under-miss (service hits while miss pending)  
Miss-under-miss (up to 4 outstanding misses)