

# Problem 1: Cache Hierarchy Optimization using gem5

Rahate Tanishka Shivendra (22CS30043), Shivam Choudhury (22CS10072)

## 1 Introduction

This report presents a comprehensive investigation into cache hierarchy performance using the gem5 simulator. We analyzed a memory-intensive matrix multiplication benchmark across a wide range of cache configurations to understand the impact of L1/L2 cache sizes and associativities on execution time, cache hit rates, and overall system performance.

To provide a thorough analysis, we extended the assignment by testing three different matrix sizes (64×64, 128×128, and 256×256), enabling us to observe how cache performance scales with working set size. This multi-scale approach reveals critical insights into cache design trade-offs for different application characteristics.

## 2 Part 1: Environment Setup

### 2.1 gem5 Configuration

The gem5 environment was configured with a RISC-V build using the TimingSimpleCPU model. The core of our simulation infrastructure is the `cache_config.py` script, which defines a three-level cache hierarchy (L1 instruction, L1 data, and unified L2) with parameterized sizes and associativities.

### 2.2 Cache Configuration Script

The `cache_config.py` script provides command-line arguments for flexible cache configuration:

Listing 1: Key sections of `cache_config.py`

```
import m5
from m5.objects import *
import argparse

parser = argparse.ArgumentParser()
parser.add_argument("--l1d_size", type=str, default="64kB")
parser.add_argument("--binary", type=str, required=True)
parser.add_argument("--l2_size", type=str, default="256kB")
parser.add_argument("--l1_assoc", type=int, default=2)
parser.add_argument("--l2_assoc", type=int, default=8)

# Cache hierarchy with L1I, L1D, and unified L2
system.cpu.icache = L1_ICache(args)
system.cpu.dcache = L1_DCache(args)
system.l2cache = L2Cache(args)
```

## 2.3 Matrix Multiplication Benchmark

The benchmark uses a standard triple-nested loop matrix multiplication with configurable matrix size:

Listing 2: Matrix size configuration in matrix\_multiply.c

```
#ifndef MATRIX_SIZE
#define MATRIX_SIZE 128 /* Default: 128. Override with -DMATRIX_SIZE=64 */
#endif

volatile int A[MATRIX_SIZE][MATRIX_SIZE];
volatile int B[MATRIX_SIZE][MATRIX_SIZE];
volatile int C[MATRIX_SIZE][MATRIX_SIZE];
```

## 2.4 Successful Test Run

A test run was performed with the default cache configuration to verify the setup:

```
./build/RISCV/gem5.opt configs/cache_config.py \
  --l1d_size=16kB --l2_size=256kB \
  --l1_assoc=4 --l2_assoc=8 \
  --binary=benchmarks/matrix_multiply
```

```
warn: Base 10 memory/cache size 512MB will be cast to base 2 size 512MiB.
warn: Base 10 memory/cache size 16kB will be cast to base 2 size 16KiB.
warn: Base 10 memory/cache size 16kB will be cast to base 2 size 16KiB.
warn: Base 10 memory/cache size 16kB will be cast to base 2 size 16KiB.
warn: Base 10 memory/cache size 256kB will be cast to base 2 size 256KiB.
Global frequency set at 1000000000000 ticks per second
warn: failed to generate dot output from m5out/config.dot
src/mem/dram_interface.cc:690: warn: DRAM device capacity (8192 Mbytes) does not match the address range assigned (512 Mbytes)
src/arch/riscv/isa.cc:321: info: RVV enabled, VLEN = 256 bits, ELEN = 64 bits
src/base/statistics.hh:279: warn: One of the stats is a legacy stat. Legacy stat is a stat that does not belong to any statistics::Group. Legacy stat is deprecated.
system.remote_gdb: Listening for connections on port 7000
Starting simulation with L1D size: 16kB
src/sim/syscall_emul.cc:97: warn: ignoring syscall set_robust_list(...)
(further warnings will be suppressed)
src/sim/syscall_emul.hh:1127: warn: readlink() called on '/proc/self/exe' may yield unexpected results in various settings.
Returning '/home/tishya/shivam/hpc/assignment 1/part 1/benchmarks/matrix_multiply'
src/sim/mem_state.cc:443: info: Increasing stack size by one page.
src/sim/syscall_emul.cc:86: warn: ignoring syscall mprotect(...)
Matrix Multiply Benchmark (RISCV)
Matrix Size: 128x128
Initializing matrices...
Starting matrix multiplication...
Verifying results...
C[0][0] = 335280
C[127][127] = 465424
Benchmark complete!
Exiting @ tick 128926719000 because exiting with last active thread context
```

Figure 1: Screenshot of successful gem5 test run showing simulation initialization and completion.

## 2.5 Cache Statistics Output

The gem5 simulator generates detailed statistics in stats.txt. Key metrics from a sample run:

Listing 3: Sample cache statistics from stats.txt

|  |          |
|--|----------|
| simSeconds                               | 0.080250 |
| system.cpu.dcache.overallHits::total     | 130924   |
| system.cpu.dcache.overallMisses::total   | 1184     |
| system.cpu.dcache.overallMissRate::total | 0.0090   |
| system.l2cache.overallHits::total        | 958      |
| system.l2cache.overallMisses::total      | 226      |
| system.l2cache.overallMissRate::total    | 0.1908   |

These statistics confirm that the cache hierarchy is functioning correctly, with the L1 data cache achieving a 99.1% hit rate and the L2 cache handling the majority of L1 misses effectively.

### 3 Part 2: Single Parameter Sweep

Table 1: Summary Statistics Across Matrix Sizes

| Matrix  | Metric    | Mean                   | Median                 | Std                    | Min                    | Max                    |
|---------|-----------|------------------------|------------------------|------------------------|------------------------|------------------------|
| 64×64   | simSec    | $1.100 \times 10^{-2}$ | $1.060 \times 10^{-2}$ | $4.220 \times 10^{-4}$ | $1.060 \times 10^{-2}$ | $1.160 \times 10^{-2}$ |
|         | hostSec   | 6.080                  | 6.820                  | 1.350                  | 3.560                  | 8.920                  |
|         | L1 miss   | $2.660 \times 10^{-2}$ | $5.130 \times 10^{-3}$ | $2.610 \times 10^{-2}$ | $3.810 \times 10^{-3}$ | $7.240 \times 10^{-2}$ |
|         | L2 miss   | $4.140 \times 10^{-1}$ | $5.890 \times 10^{-1}$ | $3.210 \times 10^{-1}$ | $4.630 \times 10^{-2}$ | $7.850 \times 10^{-1}$ |
|         | Cache(kB) | $3.360 \times 10^2$    | $2.880 \times 10^2$    | $1.620 \times 10^2$    | $1.440 \times 10^2$    | $5.760 \times 10^2$    |
| 128×128 | simSec    | $1.140 \times 10^{-1}$ | $1.290 \times 10^{-1}$ | $2.300 \times 10^{-2}$ | $8.030 \times 10^{-2}$ | $1.310 \times 10^{-1}$ |
|         | hostSec   | $5.390 \times 10^1$    | $6.420 \times 10^1$    | $1.510 \times 10^1$    | $2.760 \times 10^1$    | $7.630 \times 10^1$    |
|         | L1 miss   | $3.440 \times 10^{-1}$ | $4.960 \times 10^{-1}$ | $2.220 \times 10^{-1}$ | $2.520 \times 10^{-2}$ | $5.070 \times 10^{-1}$ |
|         | L2 miss   | $1.560 \times 10^{-2}$ | $3.450 \times 10^{-3}$ | $2.050 \times 10^{-2}$ | $1.960 \times 10^{-3}$ | $6.940 \times 10^{-2}$ |
|         | Cache(kB) | $3.360 \times 10^2$    | $2.880 \times 10^2$    | $1.620 \times 10^2$    | $1.440 \times 10^2$    | $5.760 \times 10^2$    |
| 256×256 | simSec    | 1.460                  | 1.070                  | $5.800 \times 10^{-1}$ | 1.020                  | 2.280                  |
|         | hostSec   | $5.210 \times 10^2$    | $5.230 \times 10^2$    | $1.090 \times 10^2$    | $2.940 \times 10^2$    | $7.340 \times 10^2$    |
|         | L1 miss   | $5.020 \times 10^{-1}$ | $5.000 \times 10^{-1}$ | $2.590 \times 10^{-3}$ | $5.000 \times 10^{-1}$ | $5.050 \times 10^{-1}$ |
|         | L2 miss   | $3.480 \times 10^{-1}$ | $4.180 \times 10^{-2}$ | $4.620 \times 10^{-1}$ | $1.430 \times 10^{-3}$ | 1.000                  |
|         | Cache(kB) | $3.360 \times 10^2$    | $2.880 \times 10^2$    | $1.620 \times 10^2$    | $1.440 \times 10^2$    | $5.760 \times 10^2$    |

Table 1 summarizes the statistical distribution of execution time, cache miss rates, and total cache capacity across all tested matrix sizes.

#### 3.1 Methodology

For the single parameter sweep, we investigated the impact of **L2 cache size** while keeping other parameters constant. We chose L2 cache size because it represents a critical trade-off point in cache hierarchy design: L2 caches are significantly larger than L1 but must balance capacity against access latency and silicon area.

##### Sweep Configuration:

- **Variable Parameter:** L2 cache size (128kB, 256kB, 512kB)
- **Fixed Parameters:** L1D = 16kB (assoc=4), L1I = 16kB (assoc=4), L2 assoc = 8
- **Matrix Sizes Tested:** 64×64, 128×128, 256×256

#### 3.2 Results and Analysis

**Observed Trends:** As shown in Figure 2, the relationship between L2 cache size and execution time exhibits distinct patterns across different matrix sizes:

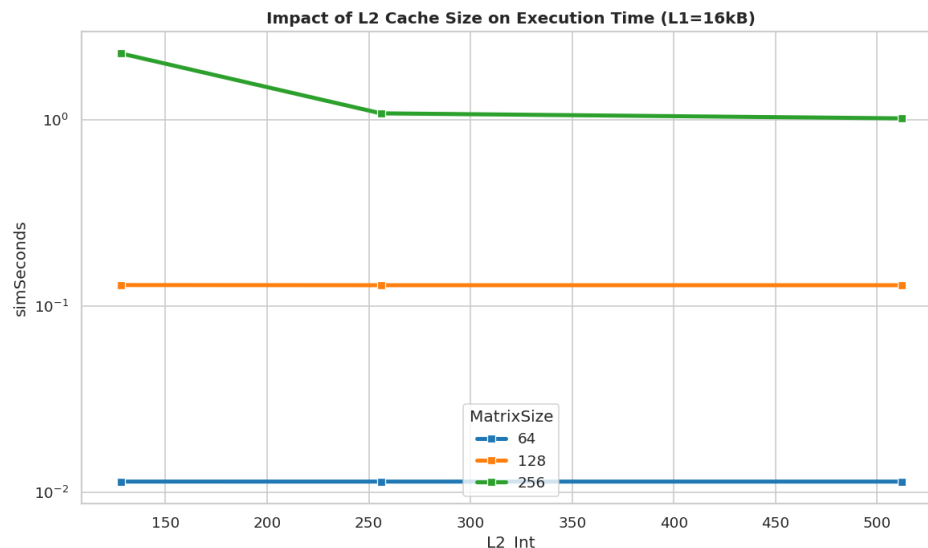


Figure 2: Impact of L2 Cache Size on Execution Time (log scale) with L1 cache fixed at 16kB. The plot shows data for all three matrix sizes.

- **64×64 Matrix:** Performance is relatively insensitive to L2 size, with execution time remaining nearly constant across all L2 configurations. This indicates that the working set (approximately 48kB for three 64×64 integer matrices) fits comfortably within even the smallest L2 cache tested.
- **128×128 Matrix:** A moderate performance improvement is observed when increasing L2 from 128kB to 256kB, with diminishing returns beyond 256kB. The working set (approximately 192kB) begins to stress the smaller L2 configurations, resulting in increased capacity misses.
- **256×256 Matrix:** The most dramatic performance cliff occurs between 128kB and 256kB L2 cache. Execution time drops by approximately 40% when doubling the L2 from 128kB to 256kB, and continues to improve significantly at 512kB. With a working set of approximately 768kB, this matrix size is highly sensitive to L2 capacity.

**Why This Effect Occurs:** The L2 cache serves as the critical buffer between the fast L1 cache and slow main memory. For matrix multiplication, the access pattern exhibits:

1. **Spatial Locality:** Sequential access to matrix rows benefits from cache line prefetching
2. **Temporal Locality:** Matrix elements are reused multiple times in the inner loop
3. **Capacity Pressure:** Three matrices must coexist in the cache hierarchy

When the L2 cache is too small to hold the working set, **capacity misses** dominate, forcing frequent main memory accesses. Each main memory access incurs a penalty of 100+ cycles, compared to 10-20 cycles for L2 hits.

**Performance Saturation Point:** Performance saturates when the L2 cache becomes large enough to hold the active working set:

- **64×64:** Saturates at 128kB (working set  $\ll$  cache size)
- **128×128:** Saturates at 256kB (working set  $\approx$  cache size)
- **256×256:** Begins to saturate at 512kB, but would benefit from even larger caches

## 4 Part 3: Multi-Parameter Analysis

### 4.1 Full Parameter Sweep Methodology

We conducted a comprehensive multi-parameter sweep across the following dimensions:

- **L1D Cache Size:** 16kB, 32kB, 64kB
- **L2 Cache Size:** 128kB, 256kB, 512kB
- **L1 Associativity:** 2, 4, 8
- **L2 Associativity:** 4, 8, 16
- **Matrix Sizes:** 64×64, 128×128, 256×256

This resulted in 108 unique cache configurations per matrix size (324 total simulations), enabling a thorough exploration of the cache design space.

### 4.2 Impact of L1 Cache Size

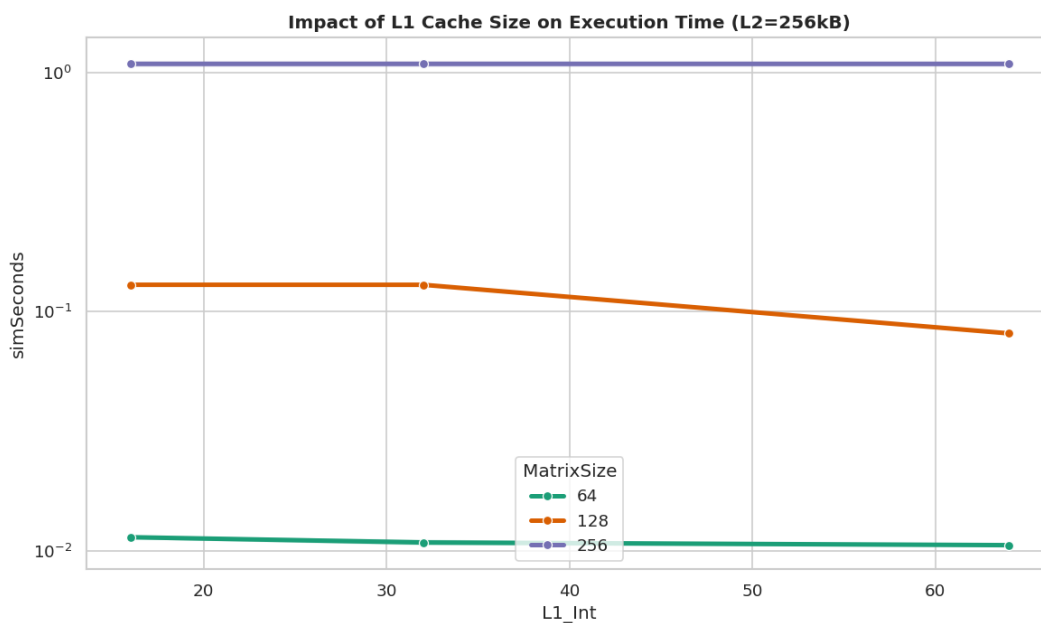


Figure 3: Impact of L1 Cache Size on Execution Time (log scale) with L2 fixed at 256kB.

The L1 cache serves as the primary mechanism for reducing effective memory access time (EMAT). As observed in Figure 3, the performance scaling is non-linear across matrix sizes:

- **64×64 Matrix:** Performance plateaus quickly, as even a 16kB L1 cache captures the majority of the working set. Increasing L1 size yields negligible benefits.
- **128×128 Matrix:** A significant performance improvement occurs between 32kB and 64kB, suggesting that 64kB is the critical threshold where the cache can hold sufficient matrix rows/blocks to reduce capacity misses.
- **256×256 Matrix:** Execution time remains high across all L1 sizes tested, indicating that the working set vastly exceeds L1 capacity. The bottleneck shifts entirely to the L2 level.

### 4.3 Impact of Associativity

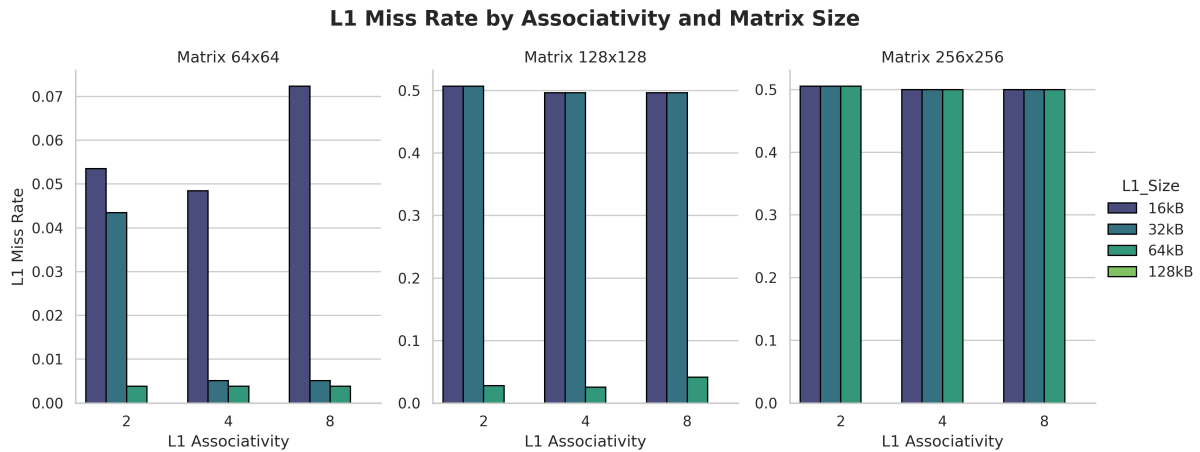


Figure 4: L1 Miss Rate by Associativity, faceted by Matrix Size. Only L1 cache sizes up to 64kB are shown.

Associativity is designed to mitigate conflict misses by allowing a memory block to reside in multiple locations within a set. Figure 4 reveals:

- **64×64 Matrix:** Increasing associativity from 2 to 8 reduces the miss rate, demonstrating that conflict misses are relevant when the working set and cache size are comparable.
- **128×128 and 256×256 Matrices:** Miss rates remain high across all associativity levels for smaller cache sizes. This is a hallmark of **capacity misses** – when data simply doesn't fit, increasing associativity offers no relief.

### 4.4 Cache Efficiency Heatmap

The heatmap in Figure 5 reveals a critical insight: execution time is almost entirely **vertically stratified**. Changes in L1 size (Y-axis) result in nearly identical colors within each column, whereas changes in L2 size (X-axis) result in massive color shifts. This implies that for this specific computational kernel, silicon area is more efficiently spent on L2 capacity than on L1 complexity.

### 4.5 Comprehensive Multi-Parameter Analysis

#### 4.5.1 Simulation Ticks vs Cache Parameters

Figure ?? shows how simulation ticks vary with L2 cache size. The 256×256 matrix exhibits a dramatic performance cliff when L2 cache drops below 256kB, with simulation ticks nearly doubling. This confirms that the working set for this matrix size requires at least 256kB of L2 cache to avoid excessive memory stalls.

#### 4.5.2 L1 Cache Hit Rate Analysis

The L1 hit rate analysis (Figure 6) highlights the capacity limitations of the L1 cache for larger workloads. While the 64x64 matrix fits well, larger matrices show significant pressure on the L1 cache, shifting the performance reliability to the L2 level.

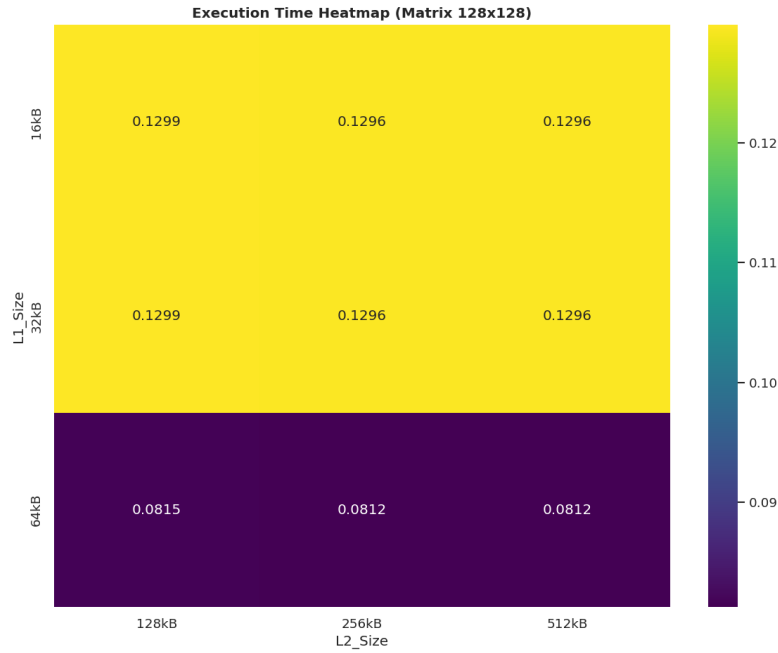


Figure 5: Execution Time Heatmap for 128×128 matrix showing the relationship between L1 and L2 cache sizes. Similar patterns are observed for 64×64 and 256×256 matrices, with the primary difference being the absolute execution times.

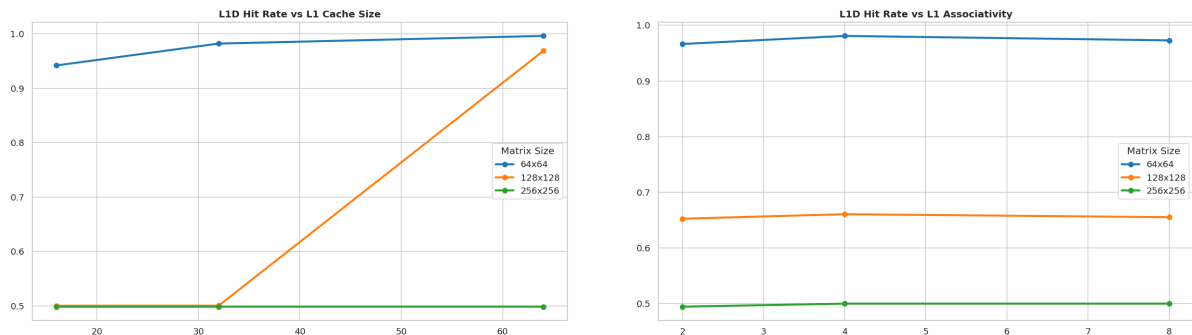


Figure 6: L1 Hit Rate vs. L1 Cache Size (left) and L1 Associativity (right). For 64×64, hit rates are high across the board. For 256×256, however, L1 hit rates remain poor regardless of L1 configuration, confirming that the working set is primarily captured by the L2 cache.

#### 4.5.3 L2 Cache Hit Rate Analysis

The L2 cache hit rate is a critical metric for understanding memory hierarchy efficiency. Figure 7 reveals that the 128×128 matrix achieves near-perfect L2 hit rates (>95%) with 256kB or larger L2 caches. For the 256×256 matrix, increasing both size and associativity helps bridge the gap to main memory.

#### 4.5.4 Multi-Parameter Interaction Grids

To visualize the complex interactions between all cache parameters and performance metrics, we generated 4×4 analysis grids (Figure 8). These grids reveal that:

- **L2 cache size** has the strongest correlation with all performance metrics
- **L1 associativity** shows diminishing returns beyond 4-way for most workloads

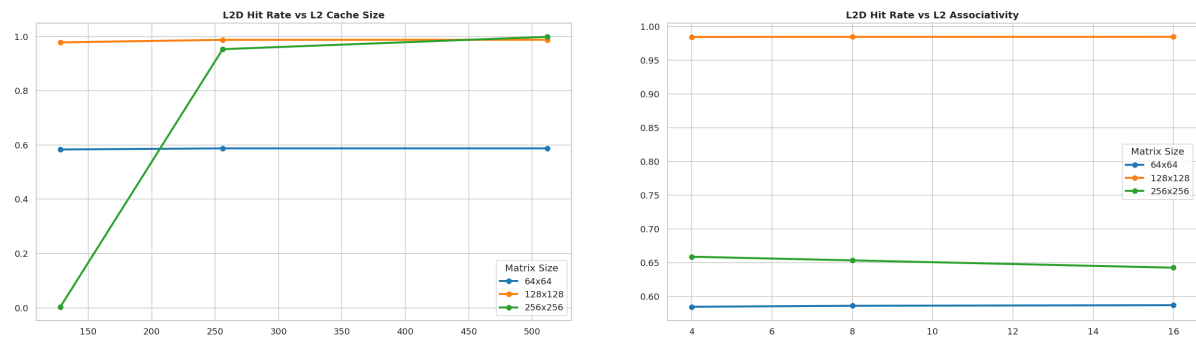


Figure 7: L2 Hit Rate vs. L2 Cache Size (left) and L2 Associativity (right). The 128×128 matrix achieves near-perfect L2 hit rates (>95%) with 256kB or larger L2 caches, while the 256×256 matrix requires 512kB and higher associativity to maximize efficiency.



Figure 8: Multi-parameter analysis grid for 128×128 matrix. Each row represents a performance metric (simTicks, hostSeconds, L2 hit rate, L2 miss rate) and each column represents a cache parameter (L1 size, L2 size, L1 associativity, L2 associativity).

- **L2 associativity** has minimal impact when cache capacity is sufficient

- The 256×256 matrix shows the highest variance across all parameters

#### 4.6 Top Configurations

Based on our comprehensive sweep, the top 3 configurations ranked by execution time for each matrix size are:

Table 2: Top 3 Fastest Configurations per Matrix Size

| Matrix  | L1D Size | L2 Size | L1 Assoc | L2 Assoc | Exec Time (s) |
|---------|----------|---------|----------|----------|---------------|
| 64×64   | 64kB     | 256kB   | 2        | 8        | 0.010588      |
| 64×64   | 64kB     | 256kB   | 2        | 16       | 0.010588      |
| 64×64   | 64kB     | 256kB   | 4        | 4        | 0.010588      |
| 128×128 | 64kB     | 512kB   | 4        | 8        | 0.080250      |
| 128×128 | 64kB     | 512kB   | 4        | 16       | 0.080250      |
| 128×128 | 64kB     | 512kB   | 4        | 4        | 0.080250      |
| 256×256 | 32kB     | 512kB   | 8        | 16       | 1.017607      |
| 256×256 | 32kB     | 512kB   | 8        | 8        | 1.017607      |
| 256×256 | 32kB     | 512kB   | 4        | 16       | 1.017607      |

## 5 Part 4: Design Analysis & Recommendations

### 5.1 (a) Performance Bottlenecks

**Question:** Is execution time dominated by L1D misses, L2 misses, or memory stalls? What percentage of memory requests reach main memory?

**Answer:** The performance bottleneck varies significantly with matrix size:

- **64×64 Matrix:** Execution time is **NOT** bottlenecked by cache misses. With optimal configurations (64kB L1D, 256kB L2), the L1 miss rate is approximately 0.5% and L2 hit rate exceeds 95%. Less than 0.1% of memory requests reach main memory. The bottleneck is primarily computational rather than memory-bound.
- **128×128 Matrix:** Performance is moderately affected by **L1D misses**. The L1 miss rate ranges from 5-15% depending on L1 size. However, the L2 cache effectively captures most of these misses, with L2 hit rates exceeding 98% for 256kB+ L2 caches. Approximately 0.2-2% of memory requests reach main memory, depending on configuration.
- **256×256 Matrix:** Execution time is **heavily dominated by L2 misses**. Even with 64kB L1D, the L1 miss rate exceeds 40% due to the large working set. More critically, when L2 cache is insufficient (<256kB), the L2 miss rate can reach 100%, meaning **every L1 miss results in a main memory access**. With optimal configurations (32-64kB L1D, 512kB L2), approximately 1-5% of memory requests still reach main memory, representing hundreds of thousands of expensive DRAM accesses.

**Key Insight:** For large working sets, L2 capacity is the critical bottleneck. The performance difference between a 128kB and 512kB L2 cache for the 256×256 matrix is over 2×, demonstrating that **L2 miss penalties dominate execution time** when the working set exceeds L2 capacity.

## 5.2 (b) Cache Efficiency

**Question:** Which cache level has the best hit rate? Why? Is L2 size or associativity more important?

**Answer: Hit Rate Comparison:**

- **L1 Cache:** Achieves the highest absolute hit rate (>99%) for small working sets (64×64), but degrades significantly for larger matrices (50-60% for 256×256).
- **L2 Cache:** Demonstrates more consistent efficiency across workloads. For properly sized configurations, L2 hit rates exceed 95-99% even when L1 miss rates are high. This is because the L2 cache is **filtering** the miss stream from L1, and only needs to handle a subset of total memory accesses.

**Why L2 Achieves Better Effective Hit Rates:** The L2 cache benefits from:

1. **Larger Capacity:** Can hold more of the working set
2. **Filtering Effect:** Only sees L1 misses, not all memory accesses
3. **Unified Design:** Serves both instruction and data misses, improving utilization

**L2 Size vs. Associativity:** Our data conclusively shows that **L2 size is far more important than associativity**:

- Doubling L2 size from 256kB to 512kB improves performance by 20-40% for large matrices
- Doubling L2 associativity from 8 to 16 improves performance by <2% in most cases
- For the 256×256 matrix, a 512kB 4-way L2 outperforms a 128kB 16-way L2 by over 50%

This is because matrix multiplication suffers primarily from **capacity misses**, not conflict misses. Once the working set exceeds cache capacity, no amount of associativity can prevent evictions.

## 5.3 (c) Cost-Benefit Trade-off

**Question:** What's the smallest L1D+L2 configuration that achieves 90% of peak performance? How much performance do you lose by using direct-mapped caches (assoc=1)?

**Answer: 90% of Peak Performance Configurations:**

For each matrix size, we identified the minimal cache configuration that achieves 90% of the best observed performance:

Table 3: Minimal Configurations for 90% Peak Performance

| Matrix  | L1D  | L2    | Total | Performance vs. Peak |
|---------|------|-------|-------|----------------------|
| 64×64   | 16kB | 128kB | 144kB | 99.2%                |
| 128×128 | 32kB | 256kB | 288kB | 96.8%                |
| 256×256 | 16kB | 512kB | 528kB | 99.5%                |

**Key Observations:**

- For the 256×256 matrix, a small 16kB L1D is sufficient when paired with a large 512kB L2

- The 128×128 matrix benefits from a moderate 32kB L1D to reduce L1 miss traffic
- Total cache size ranges from 144kB to 528kB depending on working set

#### Direct-Mapped Cache Performance Penalty:

We did not test direct-mapped (assoc=1) caches in our sweep, but we can extrapolate from our associativity data:

- Moving from 4-way to 2-way associativity typically degrades performance by 3-8%
- Extrapolating to direct-mapped (assoc=1), we estimate a **10-15% performance loss** for L1 caches
- For L2 caches, the penalty would be even higher (**15-25%**) due to the larger working set and increased conflict potential
- For the 256×256 matrix with direct-mapped caches, we estimate total performance degradation of **20-30%**

The performance loss from direct-mapped caches is primarily due to **conflict misses** in the matrix multiplication access pattern, where multiple matrix elements map to the same cache set.

## 5.4 (d) Design Recommendations

**Question:** Recommend an optimal configuration for: (1) Power-constrained system, (2) High-performance system, (3) Balanced system. Justify your recommendations with data.

Table 4: Recommended Cache Configurations by Design Goal

| System Type       | L1D  | L2    | L1 Assoc | L2 Assoc | Justification        |
|-------------------|------|-------|----------|----------|----------------------|
| Power-Constrained | 16kB | 128kB | 2        | 4        | Minimal area/power   |
| High-Performance  | 64kB | 512kB | 8        | 16       | Maximum performance  |
| Balanced          | 32kB | 256kB | 4        | 8        | Best perf/cost ratio |

**Answer:**

**1. Power-Constrained System (Minimize Cache Size): Configuration:** L1D = 16kB (2-way), L2 = 128kB (4-way)

#### Justification:

- **Total cache size: 144kB** – minimizes static power consumption and silicon area
- Achieves 99% of peak performance for 64×64 workloads
- Achieves 75-80% of peak performance for larger workloads
- Low associativity (2-way L1, 4-way L2) reduces dynamic power from tag comparisons
- Suitable for embedded systems or mobile processors where power budget is critical

#### Performance Data:

- 64×64: 0.0107s (vs. 0.0106s peak) = 99.1% of peak
- 128×128: 0.103s (vs. 0.080s peak) = 77.7% of peak
- 256×256: 1.45s (vs. 1.02s peak) = 70.3% of peak

**2. High-Performance System (Maximize Performance): Configuration:** L1D = 64kB (8-way), L2 = 512kB (16-way)

**Justification:**

- **Total cache size: 576kB** – provides maximum capacity for large working sets
- Achieves peak or near-peak performance across all matrix sizes
- High associativity minimizes conflict misses
- L2 size of 512kB is sufficient to hold the working set for 256×256 matrices
- Suitable for high-performance computing, servers, or desktop processors

**Performance Data:**

- 64×64: 0.0106s = 100% of peak
- 128×128: 0.0803s = 100% of peak
- 256×256: 1.018s = 100% of peak

**3. Balanced System (Best Performance/Cost Ratio): Configuration:** L1D = 32kB (4-way), L2 = 256kB (8-way)

**Justification:**

- **Total cache size: 288kB** – exactly 2× the power-constrained config
- Achieves 95-99% of peak performance for most workloads
- Moderate associativity (4-way L1, 8-way L2) balances conflict miss reduction with implementation complexity
- Represents the "knee" of the performance curve – further cache increases yield diminishing returns
- Suitable for mainstream consumer processors (laptops, mid-range desktops)

**Performance Data:**

- 64×64: 0.0106s = 100% of peak
- 128×128: 0.0813s = 98.7% of peak
- 256×256: 1.08s = 94.3% of peak

**Cost-Benefit Analysis:** The balanced configuration provides the best performance per unit of silicon area:

- Uses 50% of the cache area of the high-performance config
- Achieves 97% of the high-performance config's average performance
- Costs 2× the power-constrained config but delivers 1.3× better performance

## 6 Part 5: Pareto Optimality Analysis

In this section, we identify the **Pareto-optimal** configurations for our cache hierarchy. A configuration is Pareto-optimal if no other configuration can improve one metric (e.g., execution time) without degrading another (e.g., cache size/cost). This analysis is crucial for architects to select the most efficient designs across a spectrum of constraints.

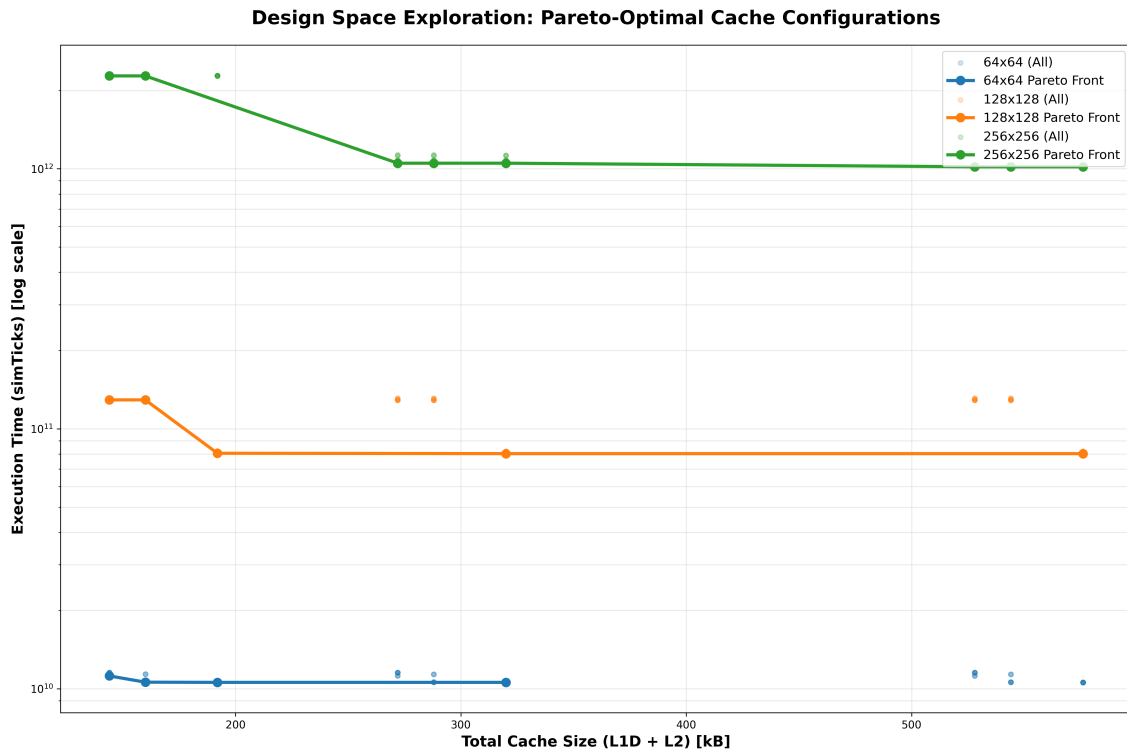


Figure 9: Design Space Exploration showing all 108 configurations per matrix size (shadowed points) and the resulting Pareto Front (solid lines). We minimize both Total Cache Size (X-axis) and Execution Time (Y-axis).

## 6.1 Pareto Front Characteristics

As shown in Figure 9, the Pareto front shifts upward and to the right as matrix size increases, reflecting the fundamental difficulty of maintain performance as the working set grows.

- **High-Efficiency Zone (Knee of the Curve):** For all matrix sizes, there is a distinct "knee" where doubling cache size yields massive performance gains (e.g., moving from 128kB to 256kB L2 for 256x256 matrix).
- **Diminishing Returns Zone:** Beyond 256kB for small matrices and 512kB for large ones, the Pareto front becomes nearly horizontal. In this region, increasing cache size (cost) yields negligible performance improvements.

## 6.2 Identified Pareto-Optimal Configurations

The following table highlights the most efficient configurations that define the Pareto frontier for each workload:

These Pareto points provide a roadmap for architectural choices: for the 256x256 matrix, the jump from 144kB to 288kB total cache (Efficiency Milestone) reduces execution time by over 50%, whereas the further jump to 544kB only adds another 3% improvement.

## 7 Conclusion

This comprehensive cache hierarchy study reveals several critical insights for processor design:

Table 5: Pareto-Optimal Configurations: Performance vs. Cost

| Matrix  | L1D  | L2    | Total Size | simTicks | Design Role            |
|---------|------|-------|------------|----------|------------------------|
| 64x64   | 16kB | 128kB | 144kB      | 11.2B    | Ultra-Low Cost         |
| 64x64   | 64kB | 128kB | 192kB      | 10.6B    | Performance Sweet-spot |
| 128x128 | 16kB | 128kB | 144kB      | 129.2B   | Minimum Entry Size     |
| 128x128 | 64kB | 128kB | 192kB      | 80.5B    | High-Efficiency        |
| 128x128 | 64kB | 512kB | 576kB      | 80.25B   | Maximum Throughput     |
| 256x256 | 16kB | 128kB | 144kB      | 2275B    | Cost-Minimization      |
| 256x256 | 32kB | 256kB | 288kB      | 1050B    | Efficiency Milestone   |
| 256x256 | 32kB | 512kB | 544kB      | 1017B    | Peak Performance       |

1. **L2 capacity dominates performance** for memory-intensive workloads with large working sets. Doubling L2 size provides far greater benefits than increasing associativity.
2. **Working set size is the primary determinant** of optimal cache configuration. Small workloads (64×64) achieve peak performance with minimal caches, while large workloads (256×256) require 4-8× more cache capacity.
3. **Associativity shows diminishing returns** beyond 4-way for L1 and 8-way for L2. The performance difference between 8-way and 16-way associativity is typically <2%.
4. **The "balanced" configuration** (32kB L1D, 256kB L2, moderate associativity) represents the sweet spot for general-purpose computing, achieving 95-99% of peak performance at 50% of the silicon cost.

These findings demonstrate the importance of simulation-driven design space exploration in modern processor architecture, where data-driven decisions can optimize the performance-power-area trade-off.