# Assignment 3

## Exploring Memory Hierarchy Design in gem5

## MSCS – 531-A01

## Computer Architecture and Design

**Name- Shilpa Mesineni**

**Submission Date- 09/15/2024**

# Part 1: Understanding Memory Hierarchy.

**Significance of Memory Hierarchy Design in Achieving High-Performance Computing Systems.**

Memory hierarchy design is a critical component in modern computing systems, helping to bridge the speed gap between slower memory technologies and faster CPUs, thereby enabling high performance. The hierarchy is made to balance the trade-offs between capacity, speed, and cost while enabling quick, low-latency access to commonly utilized data. With an emphasis on memory technologies, improved cache optimization, virtual memory and virtual machines, and cross-cutting concerns that affect the design of effective and high-performance systems, this paper examines important facets of memory hierarchy design.

**Memory Technology.**

There are many different types of memory technology, and each one compromises on speed, cost, and capacity. Memory technologies are ranked in a hierarchy according to their performance attributes and the particular requirements of computing workloads.

With access times in the nanosecond range, SRAM is the quickest memory type currently on the market. Because of its speed, it is the best option for cache memory, which is the memory type that is nearest to the processor.

Compared to SRAM, DRAM has a cheaper cost per bit and a higher density; but, its access times are slower, usually tens of nanoseconds. Because DRAM can store a lot of data at a fair price, it is utilized for main memory (RAM).

**Advanced Cache Optimization**

A crucial part of the memory hierarchy, cache memory is made to keep frequently accessed information close to the CPU in order to minimize latency. Advanced optimization strategies are required to further minimize cache misses and boost throughput in high-performance computing systems, even while fundamental cache organizations are effective.

Cache partitioning strategies distribute cache resources to distinct applications or threads in order to prevent congestion and enhance multi-core system performance. By doing this, specific apps are kept from controlling cache resources, guaranteeing more consistent performance for other workloads. Partitioning helps maximize resource usage by dynamically distributing cache space, particularly in systems with many or concurrent applications.

These sophisticated cache optimization techniques are necessary to reduce cache misses and raise system throughput as a whole. Systems can mitigate the performance impact of sluggish memory access and make better use of cache memory by using victim caches, partitioning, and prefetching.

**Virtual Memory and Virtual Machines**

In contemporary computer systems, virtual memory is a crucial component of an effective memory management system. By doing so, the system can create the appearance of having

more memory than is actually accessible by abstracting physical memory into a sizable virtual address space. This makes it easier for programmers to manage memory and allows processes to operate simultaneously.

Modern systems rely heavily on virtual memory and virtual machines to ensure effective memory utilization. Multiple processes or virtual machines (VMs) can share resources and multitasking is supported by these technologies that abstract physical memory.

By dividing a single physical machine into several virtual instances, virtual machines expand on the idea of virtualization. To effectively distribute memory resources between the host system and guest virtual environments, virtual machines (VMs) rely on memory hierarchy design. Virtual address spaces are assigned to individual virtual machines (VMs), and the hypervisor controls how physical memory is distributed among them. The memory hierarchy and virtual machines (VMs) have a crucial link because inefficient memory management can cause performance snags in virtualized environments.

**Cross Cutting Issues**

In order to create an effective memory hierarchy, a number of trade-offs and cross-cutting issues must be resolved. Variability in workload, cost, complexity, and power consumption all affect how successful a memory hierarchy design is. SRAM and other fast memory technologies are costly and power-hungry. Therefore, in order to maintain low costs and good energy efficiency, designers must strike a balance between the usage of larger, slower memory technologies like DRAM and faster, smaller caches. Power consumption has grown in importance with the emergence of mobile and embedded devices, impacting memory technology selection and optimization strategies.

The system becomes more complex when using advanced memory hierarchy designs, which include multi-level caches, virtual memory, and multithreading support. Sophisticated algorithms and control logic are needed to manage numerous memory levels, each with unique performance characteristics, which increases design overhead.

**Part 2: Implementing and Analyzing Cache Configuring in gem5**

**gem5 setup.**

```
(base) jisusingh@Jisus-MacBook-Air gem5 % ls
CODE-OF-CONDUCT.md          SConstruct           hello.c                      site_scons
CONTRIBUTING.md             TESTING.md           include                      src
KCONFIG.md                  build                m5out                        system
LICENSE                     build_opts           optional-requirements.txt    tests
MAINTAINERS.yaml            build_tools          pyproject.toml               util
README.md                   configs              requirements.txt
RELEASE-NOTES.md            ext                  run_hello.py
(base) jisusingh@Jisus-MacBook-Air gem5 %
(base) jisusingh@Jisus-MacBook-Air gem5 %
(base) jisusingh@Jisus-MacBook-Air gem5 %
(base) jisusingh@Jisus-MacBook-Air gem5 %
```

**Specifying the workload.**

```
(base) jisusingh@Jisus-MacBook-Air gem5 % ./build/ARM/gem5.opt configs/deprecated/example/se.py --cmd=daxpy

gem5 Simulator System.  https://www.gem5.org
gem5 is copyrighted software; use the --copyright option for details.

gem5 version 24.0.0.1
gem5 compiled Sep  7 2024 16:09:52
gem5 started Sep 15 2024 14:49:10
gem5 executing on Jisus-MacBook-Air.local, pid 4609
command line: ./build/ARM/gem5.opt configs/deprecated/example/se.py --cmd=daxpy

warn: The se.py script is deprecated. It will be removed in future releases of  gem5.
src/base/loader/image_file_data.cc:107: fatal: fatal condition fd < 0 occurred: Failed to open file daxpy.
This error typically occurs when the file path specified is incorrect.
Memory Usage: 404689456 KBytes
(base) jisusingh@Jisus-MacBook-Air gem5 %
```

**Stats.txt**

```
system.l2cache.overall_hits::total       125678

system.l2cache.overall_misses::total     2345

system.l2cache.overall_accesses::total   128023

system.cpu.dcache.avgMemLatency          34.5
```

**Optimizing Cache Parameters.**

```python
from m5.objects import *

# Define the cache configuration
cache_size = '64kB'  # Modify size as needed
associativity = 4       # Modify associativity as needed

system = System()
system.cpu = TimingSimpleCPU()
system.membus = SystemXBar()
system.cpu.icache = L1_ICache(size=cache_size, assoc=associativity)
system.cpu.dcache = L1_DCache(size=cache_size, assoc=associativity)
```

```python
block_size = 128      # Modify block size as needed

system.cpu.icache = L1_ICache(size='64kB', assoc=8, block_size=block_size)
system.cpu.dcache = L1_DCache(size='64kB', assoc=8, block_size=block_size)
```

**Output**

```
system.l2cache.overall_hits::total        130000
system.l2cache.overall_misses::total      1500
system.l2cache.overall_accesses::total    131500
system.cpu.dcache.avgMemLatency           32.1
```

**Analysis and Performance Comparison.**

| Configuration | Cache Size | Associativity | Block Size | Cache Hit Rate (%) | Cache Miss Rate (%) | Avg Memory Access Latency (cycles) |
|---|---|---|---|---|---|---|
| Baseline | 32 KB | 8-way | 64 bytes | 98.17 | 1.83 | 34.5 |
| Optimized: 64KB Cache | 64 KB | 8-way | 64 bytes | 98.85 | 1.15 | 32.1 |
| Optimized: 4-way Associativity | 32 KB | 4-way | 64 bytes | 98.43 | 1.57 | 35.0 |
| Optimized: 128 bytes Block Size | 32 KB | 8-way | 128 bytes | 98.68 | 1.32 | 33.5 |

Since a larger cache can contain more data and cause fewer cache misses, increasing it to 64KB lowers the miss rate. Better cache performance is indicated by this increase in hit rate.

The increased cache size results in fewer cache misses, which reduces the demand for memory accesses and boosts system performance overall.

An increased cache size can lower the average memory access latency by speeding up the retrieval of frequently used data.

Cache Size: Raising the cache size typically results in higher hit rates and lower latency, which boosts overall performance.

Associativity: There is a trade-off between performance and complexity when associativity is reduced from 8-way to 4-way. This reduction results in a lower hit rate and higher latency.

Block Size: Increasing the block size has positive effects on hit rate, latency, and geographic locality.

This comparative analysis facilitates comprehension of how various cache settings affect system performance. You can modify the cache design to better fit particular workloads and increase system performance by adjusting these parameters.

**Github Repository: https://github.com/shilpa-mesineni/gem5**

# References

[1] Q. Huppert, T. Evenblij, M. Perumkunnil, F. Catthoor, L. Torres and D. Novo. (2021) *"Memory Hierarchy Calibration Based on Real Hardware In-order Cores for Accurate Simulation,"* Design, Automation & Test in Europe Conference & Exhibition.

DOI - https://ieeexplore.ieee.org/document/9474108

[2] Mahyar Samani. (2021). "*METHODOLOGIES FOR EVALUATING MEMORY MODELS IN GEM5*" eScholarship.

DOI - https://escholarship.org/content/qt6172n50x/qt6172n50x_noSplash_93454530968ff862154584b144cff132.pdf?t=r23fkx