

NC State University
Department of Electrical and Computer Engineering
ECE 463/521: Fall 2015 (Rotenberg)
Project #1: Cache Design, Memory Hierarchy Design

By
PARTH BHOGATE
(NCSU ID: 200108628)

NCSU Honor Pledge: "I have neither given nor received unauthorized aid on this test or assignment."

Student's electronic signature: Parth Bhogate
(sign by typing your name)

Course number: ECE 521

Graph #1: L1 Miss Rate vs. $\log_2(\text{L1 SIZE})$, without VC and L2

Cache Configuration:

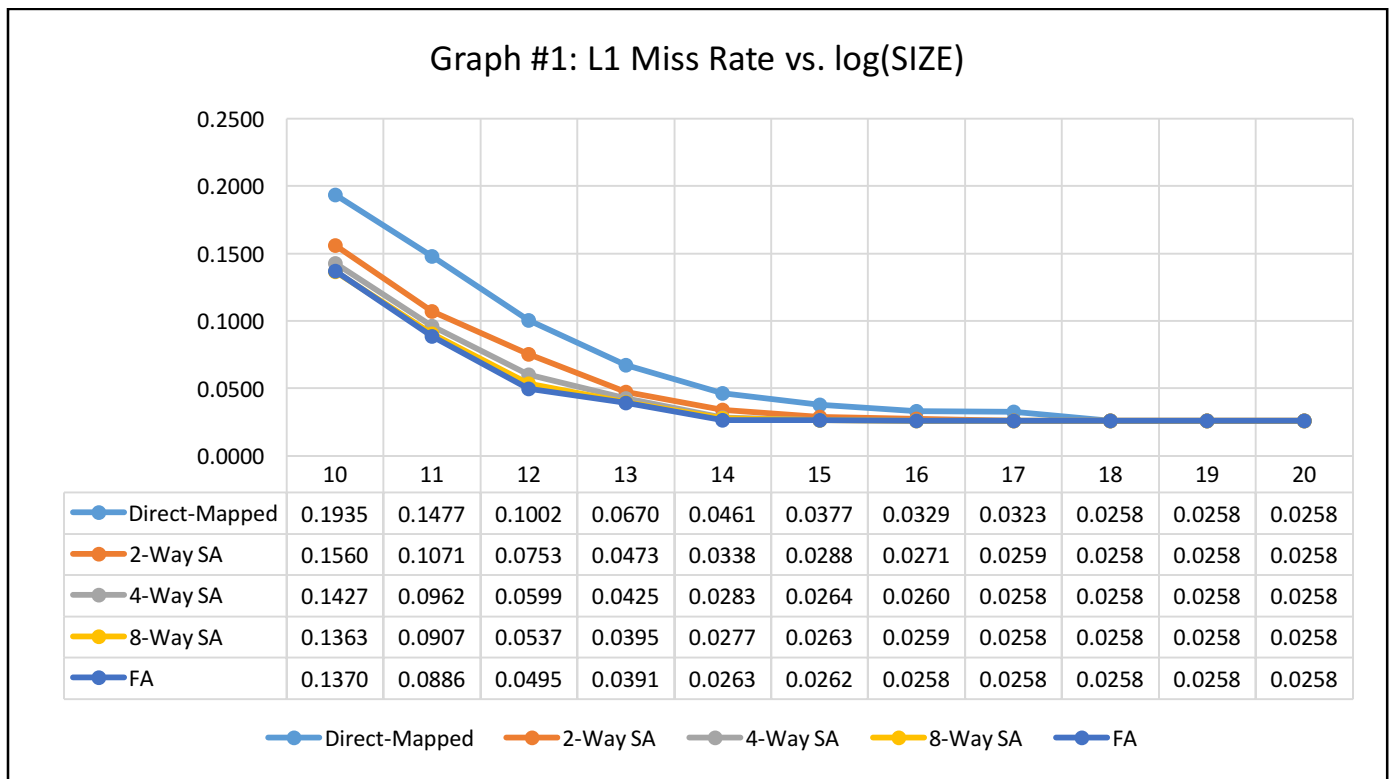
L1 Cache Size = 1KB, 2KB, 4KB, ..., 1MB

L1 Cache Associativity: Direct-mapped, 2-way SA, 4-way SA, 8-way SA,

Fully associative

No VC, No L2 Cache

The L1 miss rate is plotted against the cache size. Each line on the plot shows the miss rate for a particular associativity: direct-mapped, 2-way set associative, 4-way set associative, 8-way set associative and fully associative.



1. Trends in the graph:

From the graph, we can observe that the miss rate decreases (almost) exponentially for an increasing cache size. For higher cache sizes, the miss rates for the various cache associativities converge to approximately the same value. With a growing cache size, the conflict and capacity misses are nullified, and only the compulsory miss rate can still be observed in the graph.

For a given cache size, the miss rate decreases with an increase in cache associativity. This decrease is primarily due to the decrease in conflict miss rates.

2. Compulsory Miss Rate Estimation:

Compulsory misses occur when a particular cache block is referenced for the first time. Compulsory misses cannot be reduced by the standard techniques, except for software or hardware prefetching. The miss rate observed in a sufficiently large cache with full associativity can be approximated to be the compulsory miss rate. From the graph, the compulsory miss rate can be estimated to be **0.0258**.

3. Conflict Miss Rate Estimation:

Misses which occur in the cache under test but not in the fully-associative cache being used for comparison are classified as conflict misses. The difference between the miss rate for the fully-associative cache versus that for a N-Way set associative cache is the fraction of the miss rate that can be attributed to conflict misses. From the graph, the following values for conflict miss rate are obtained:

Assoc / Cache Size	1KB	2KB	4KB	8KB	16KB	32KB	64KB	128KB	256KB	512KB	1MB
Direct-Mapped	.056	.059	.051	.027	0.19	.011	.007	0.006	0.00	0.00	0.00
2-Way SA	.019	.018	.026	.008	.007	.003	.001	.0001	0.00	0.00	0.00
4-Way SA	.006	.007	.010	.003	.002	.0002	.0001	0.00	0.00	0.00	0.00
8-Way SA	--	.002	.004	0.00	.001	0.00	0.00	0.00	0.00	0.00	0.00

A fully-associative cache does not suffer from conflict misses; it suffers only compulsory and capacity misses.

Graph #2: AAT vs. log₂(L1 SIZE), without VC or L2

Cache Configuration:

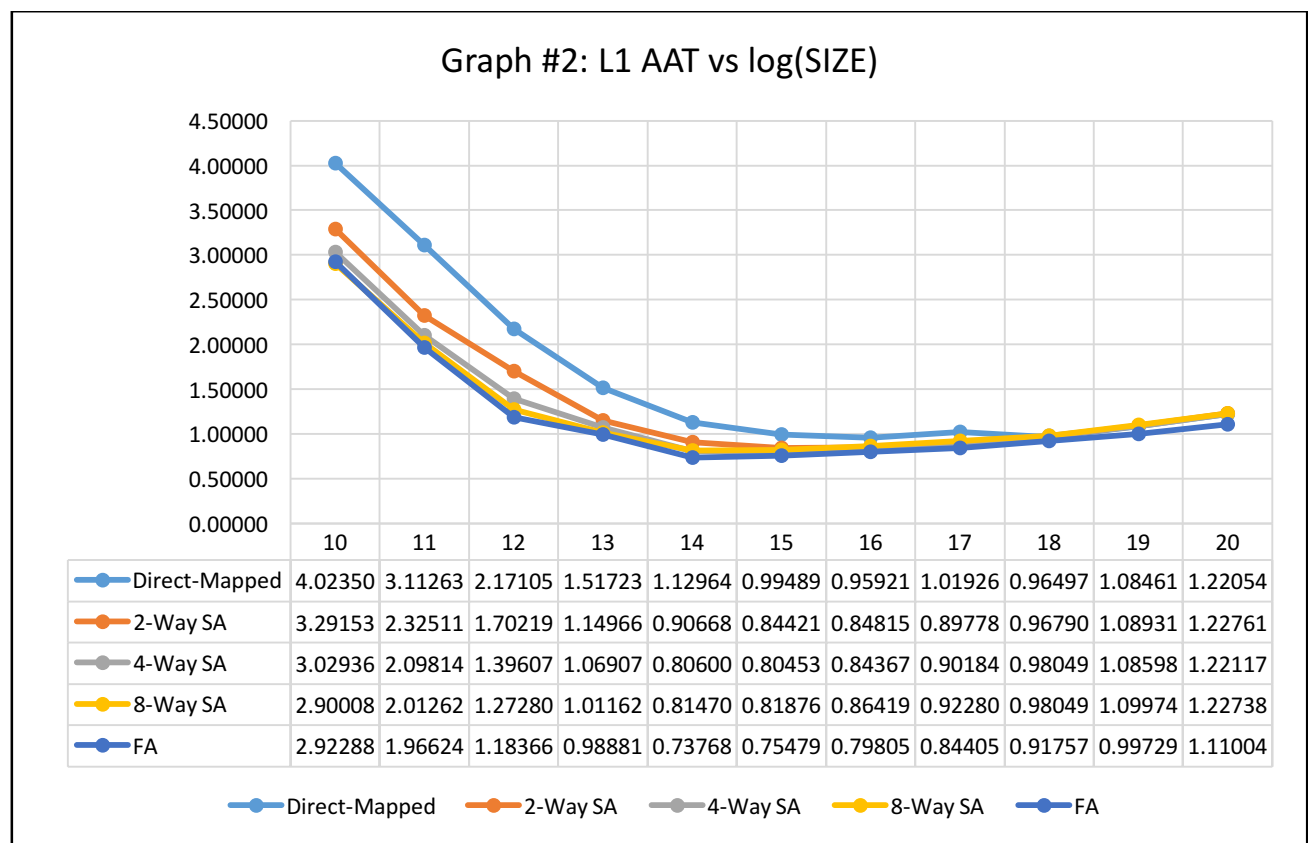
Blocksize = 32

L1 Cache Size = 1KB, 2KB, 4KB, ..., 1MB

L1 Cache Associativity: Direct-mapped, 2-way SA, 4-way SA, 8-way SA, Fully associative

No VC, No L2 Cache

The graph plots the Average Access Time (AAT) for L1 cache for different cache sizes and cache associativities. L2 Cache or Victim Cache is not included in the memory hierarchy.



1. Which configuration yields the lowest AAT?

From the graph, the lowest AAT observed is **0.737 ns**, which is for a **fully-associative L1 cache of size 16KB**. AAT depends upon the cache hit time, miss rate and the miss penalty. The miss penalty is the same for all the configurations. Hit time and the miss rate depend both on the associativity and the cache size. Hit time increases with the cache size, and the miss rate decreases with size. It is thus a trade-off between the size and associativity. With a cache of 16KB and full associativity, the trade-off gives the optimal value of the AAT.

In caches with full-associativity, the AAT increases with further increase in the cache size on account of increase in the hit time.

Graph #3: AAT vs. $\log_2(\text{L1 SIZE})$, with L2 Cache included

Cache Configuration:

Blocksize = 32

L1 Cache Size = 1KB, 2KB, 4KB, ..., 1MB

L1 Cache Associativity: Direct-mapped, 2-way SA, 4-way SA, 8-way SA,

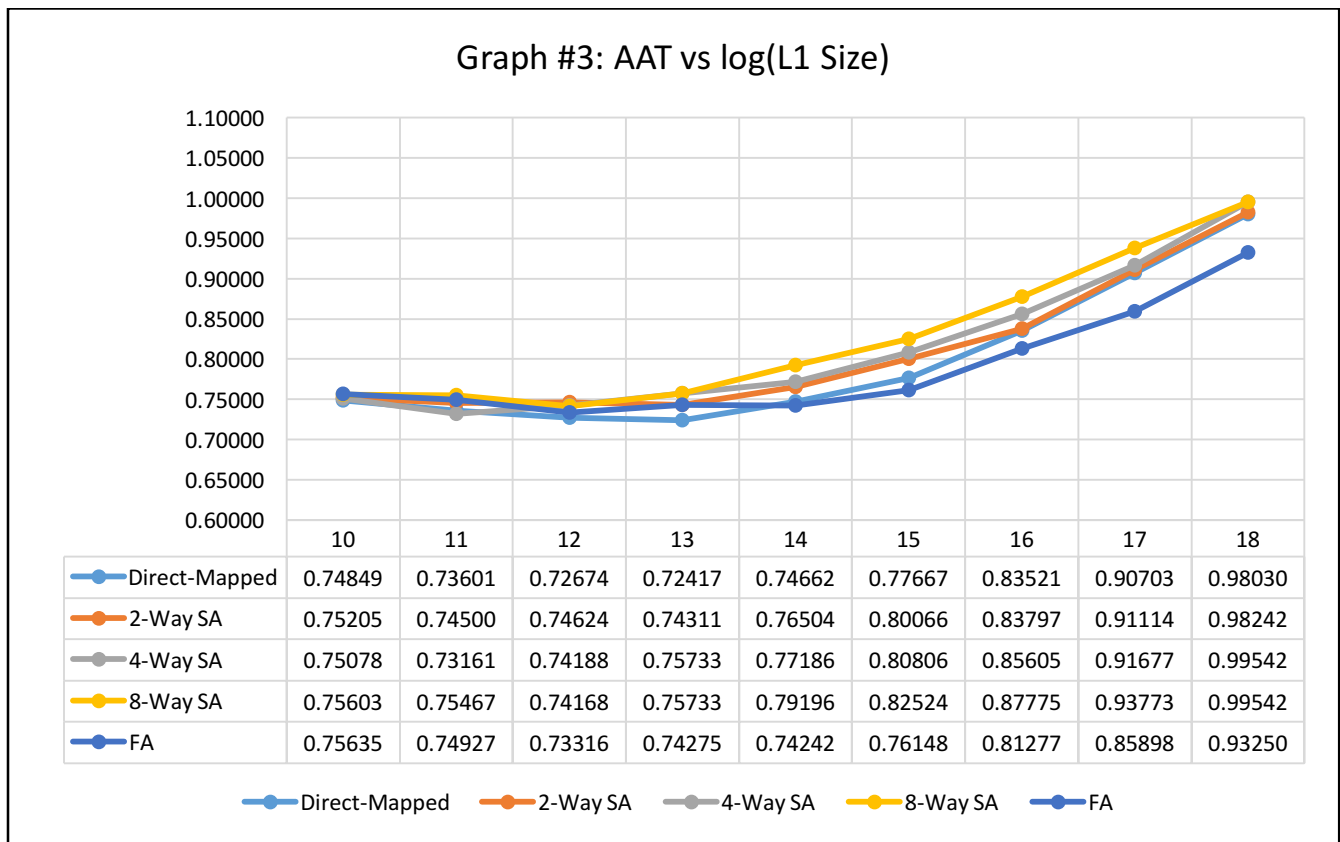
Fully associative

L2 Cache Size = 512KB

L2 Associativity: 8-way set associative

No VC

The graph plots the AAT versus the L1 cache size for different L1 cache sizes and associativities. L2 cache of fixed size and associativity is added to the memory hierarchy. Victim cache is not present.



1. With L2 added, which L1 cache configurations result in AATs close to the best AAT observed without L2 in the memory hierarchy?

The lowest AAT observed in Graph #2 was 0.737 ns. Looking at the AAT values in Graph #3 which are lesser than 0.774 ns, we observe that the AATs for the following configurations are within 5% of the lowest AAT in Graph #2:

(Yes implies that AAT is within 5% of the best AAT from Graph #2; No implies the AAT is not within a 5% range.)

Assoc/Cache Size	1KB	2KB	4KB	8KB	16KB	32KB	64KB	128KB	256KB
Direct-Mapped	Yes	Yes	Yes	Yes	Yes	No	No	No	No
2-Way SA	Yes	Yes	Yes	Yes	Yes	No	No	No	No
4-Way SA	Yes	Yes	Yes	Yes	Yes	No	No	No	No
8-Way SA	Yes	Yes	Yes	Yes	No	No	No	No	No
Fully-Associative	Yes	Yes	Yes	Yes	Yes	Yes	No	No	No

2. With L2 added to the memory hierarchy, which L1 cache configuration yields the lowest AAT?

The lowest AAT as seen in the graph is 0.724 ns, for a **direct-mapped L1 cache of 8KB** size. In Graph #2 which did not include a L2 cache, the lowest observed AAT was 0.737 ns; thus, the memory hierarchy with L2 included has a AAT which is **0.013 ns lower, which is 1.76%** lower than the one without L2. The presence of L2 reduces the miss penalty in the hierarchy and thus leads to a decrease in AAT.

3. Comparison of total areas required for optimal-AAT configuration with and without L2 cache

From Graph #2, the total area for optimal AAT configuration is as follows:

Area of 16KB fully-associative L1 cache without L2 cache:

Area of L1 Cache = **0.063446019 mm²**.

From Graph #3, the total area for optimal AAT configuration is:

Area of 8KB direct-mapped L1 cache with 512KB 8-way associative L2 cache -

Area of L1 Cache + Area of L2 Cache

= 0.053293238 + 2.640142073

Total memory area = **2.693435311 mm²**.

The increase in area when L2 cache is included in the memory hierarchy is:

2.693435311 - 0.063446019 = **2.629989292 mm²**

Percentage increase is $2.629989292 / 0.063446019 * 100 = \mathbf{4145.24\%}$.

Thus, for optimal AAT, the total area for L1-only cache configuration is much lower compared to a hierarchy with L2 included. For a small improvement in the best AAT, there is huge increase in the total memory hierarchy area. This clearly shows the trade-off between achieving an improvement in the AAT and the memory hierarchy area.

Graph #4: L1 Miss Rate vs. $\log_2(\text{BLOCKSIZE})$, for fixed L1 Associativity

Cache Configuration:

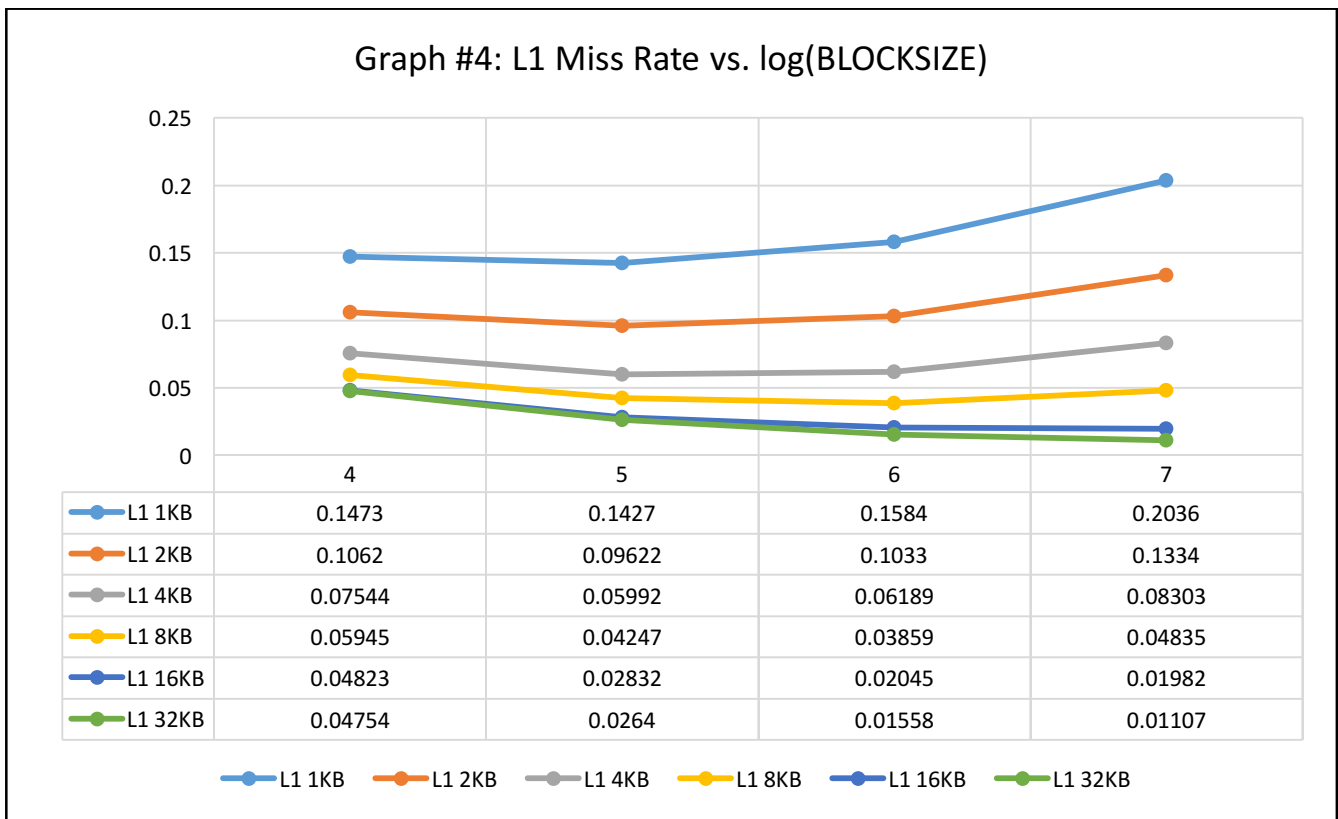
Blocksize = 32

L1 Cache Size = 1KB, 2KB, 4KB, ..., 1MB

L1 Cache Associativity: 4-way set associative

No VC, No L2 Cache

The graph plots the L1 miss rate versus the L1 blocksize. The associativity of the cache is fixed, and the miss rate is plotted versus different block sizes so as to study the effect of varying the cache block size. Each line in the graph depicts the miss rate for a particular L1 cache size.



1. Trends observed in the graph:

The cache blocksize influences the amount of spatial locality present in the cache. If the address trace has good spatial locality, a larger blocksize will help to reduce the miss rate. On the other hand, a large blocksize leads to fewer unique blocks being present in the cache. Thus, a larger blocksize leads to lesser variety of blocks in the cache. This might lead to an increase in miss rate.

For smaller cache sizes, increasing the block size initially decreases the miss rate, but the cache pollution effect takes over quickly since the low cache size does not permit bringing in a wide variety of blocks. Thus, a **small cache size works better with a smaller blocksize**.

In the case of a larger cache size, the miss rate decreases with an increase in blocksize because of exploiting the spatial locality. Since the large cache can bring in a variety of different blocks along with sequential blocks, the miss rate goes on decreasing with an increasing blocksize. The cache pollution effect does not kick in (atleast until the blocksize is made a significant percent of the total cache size) for the case of a large cache size. Thus, **a large cache size works better with a larger blocksize**.

The effect on miss rate of an increase in blocksize depends upon the trade-off between the positive impact of the increase in spatial locality and the negative impact of possible cache pollution because of bringing in lesser variety of memory blocks. The overall effect of varying blocksize depends upon the total cache size. For a smaller cache size, the cache pollution effect begins to dominate, whereas for a larger cache size the spatial locality provides a greater positive impact than the negative effect of cache pollution.

Graph #5: AAT vs. $\log_2(\text{L1 SIZE})$, with different L2 Cache Sizes

Cache Configuration:

Blocksize = 32

L1 Cache Size = 1KB, 2KB, 4KB, ..., 1MB

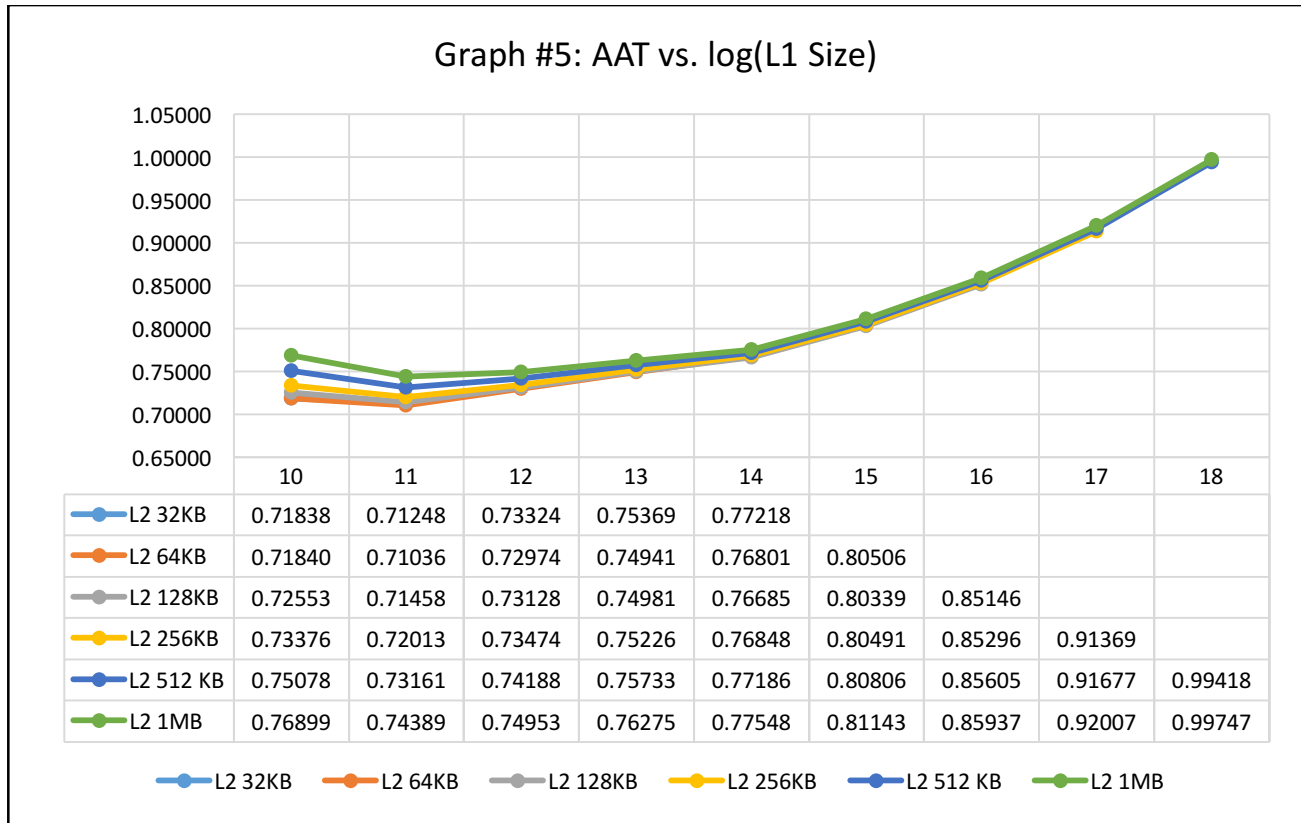
L1 Cache Associativity: 4-way set associative

L2 Cache Size = 32KB, 64KB, 128KB, 256KB, 512KB, 1MB

L2 Cache Associativity: 8-way set associative

No VC

The graphs looks at the AAT variation with L1 cache sizes for different L2 cache sizes included in the memory hierarchy. The associativity for both L1 and L2 caches is fixed. Victim cache is not included with L1 cache.



1. Which memory hierarchy configuration yields the lowest AAT?

The lowest AAT as obtained from the graph is for a **2KB 4-way set associative L1 cache along with a 64KB 8-way set associative L2 cache**. The lowest AAT value is found to be **0.710 ns**. The AAT in this hierarchy involves a complex trade-off between the L1 hit time

(which depends on L1 cache size), the L1 miss rate (depends on L1 size and associativity), the L2 hit time (depends on L2 size), and the L2 miss rate (again depends on L2 size and associativity). Moving towards bigger cache sizes or associativities leads to a decrease in miss rate but an increase in hit time. The optimal AAT is the one which has the balance between size and associativity for both the caches.

2. Which memory hierarchy configuration has the smallest total area, that yields an AAT within 5% of the best AAT?

Since the lowest AAT is 0.710 ns, thus we consider configurations whose AAT values are lesser than 0.7455 ns. Comparing the total cache areas for these memory configurations, the minimum area is found to be for a 4-way set associative **1KB L1 cache along with a 32KB 8-way set associative L2 cache**. The area for this configuration is **0.257 mm²**. Also, the AAT is **0.718 ns**, which is within 1.1% of the lowest AAT. Here, the trade-off is between greater area needed for a bigger cache and the decrease in miss rate observed on account of the greater size. The optimal-area configuration gives a reasonably good AAT for the least chip area.

Graph #6: AAT vs. $\log_2(\text{L1 SIZE})$, with Victim Cache

Cache Configuration:

Blocksize = 32

L1 Cache Size = 1KB, 2KB, 4KB, ..., 32KB

L1 Cache Associativity: Direct-mapped, 4-way, 8-way, 16-way set associative

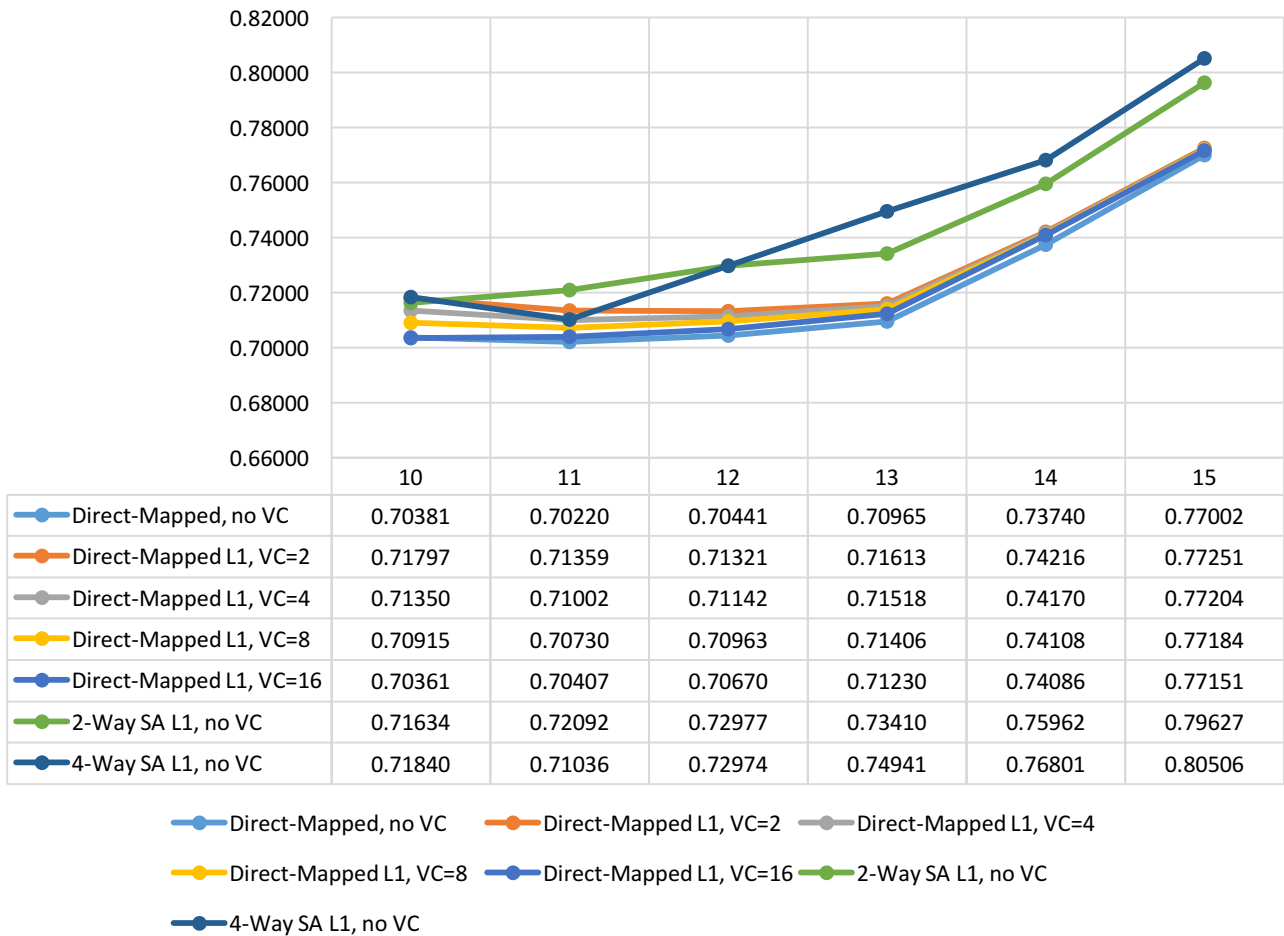
L2 Cache Size = 64KB

L2 Cache Associativity: 8-way set associative

Victim Cache: No. of entries is varied – 0, 2, 4, 8, 16

The graph shows the effect of adding a Victim Cache to augment the L1 cache. By varying the size and associativity of L1 cache and the number of memory blocks in the fully-associative victim cache, the effect on the Average Access Time (AAT) is studied in detail. A fixed L2 cache of size 64KB with 8-way set associativity is included in the memory hierarchy.

Graph #6: AAT vs. $\log_2(\text{L1 size})$



1. Trends in the graph:

Comparing the miss rates from the address trace data and the AAT's from the graph, we observe that the performance of L1 cache augmented with a Victim Cache is better than that of a set-associative L1. With a smaller cache size, the AAT for cache with Victim Cache is comparable to that of a set-associative cache. However, as the size of the cache is increased, the Victim Cache configuration outperforms the set-associative configuration.

The performance of a **1KB 2-way set associative L1 cache** is thus comparable to a 1KB L1 cache with a 2-entry or a 4-entry Victim Cache. As the size of the L1 cache increases, the Victim Cache configuration performs better than the set associative cache.

2. Which memory hierarchy configuration yields the best AAT?

The lowest AAT as seen from the graph is **0.702 ns**, for a **direct-mapped 2KB L1 cache with no Victim Cache**. Although the miss rate for this configuration is not the least among the various configurations, the overall AAT is lower because of no additional delay on account of VC access. Other comparable low AAT values are obtained with a 1KB, 2KB and 4KB L1 cache with 8 VC blocks. These configurations have a sufficiently low miss rate to offset the extra hit time due to presence of Victim Cache.

3. Which memory hierarchy has the smallest total area, that yields an AAT within 5% of the best AAT?

The best observed AAT is **0.702 ns**. Looking at total cache areas for configuration with an AAT within 5%, *i.e.* AAT upto 0.737 ns, the smallest total area is found to be for a **1KB 2-way set associative L1 cache with no Victim Cache, with 64KB 8-way set associative L2 Cache**. The area is found to be **0.369 mm²**. The AAT for this configuration is **0.716 ns**, which is within 2% of the lowest AAT. The trade-off here is between the extra area cost added to the memory hierarchy due to the presence of Victim Cache, and the decrease in the miss rate because of the presence of Victim Cache. The lowest area is found to be for a direct-mapped cache with no VC, which has lesser area since Victim Cache is absent and also a reasonably low AAT at the same time.