

**NC State University**

**Department of Electrical and Computer Engineering**

**ECE 463/521: Spring 2016**

**Project #1: Cache Hierarchy**

**by**

**Ritesh Gajare**

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

Student's electronic signature: \_\_\_\_Ritesh Gajare\_\_\_\_  
(sign by typing your name)

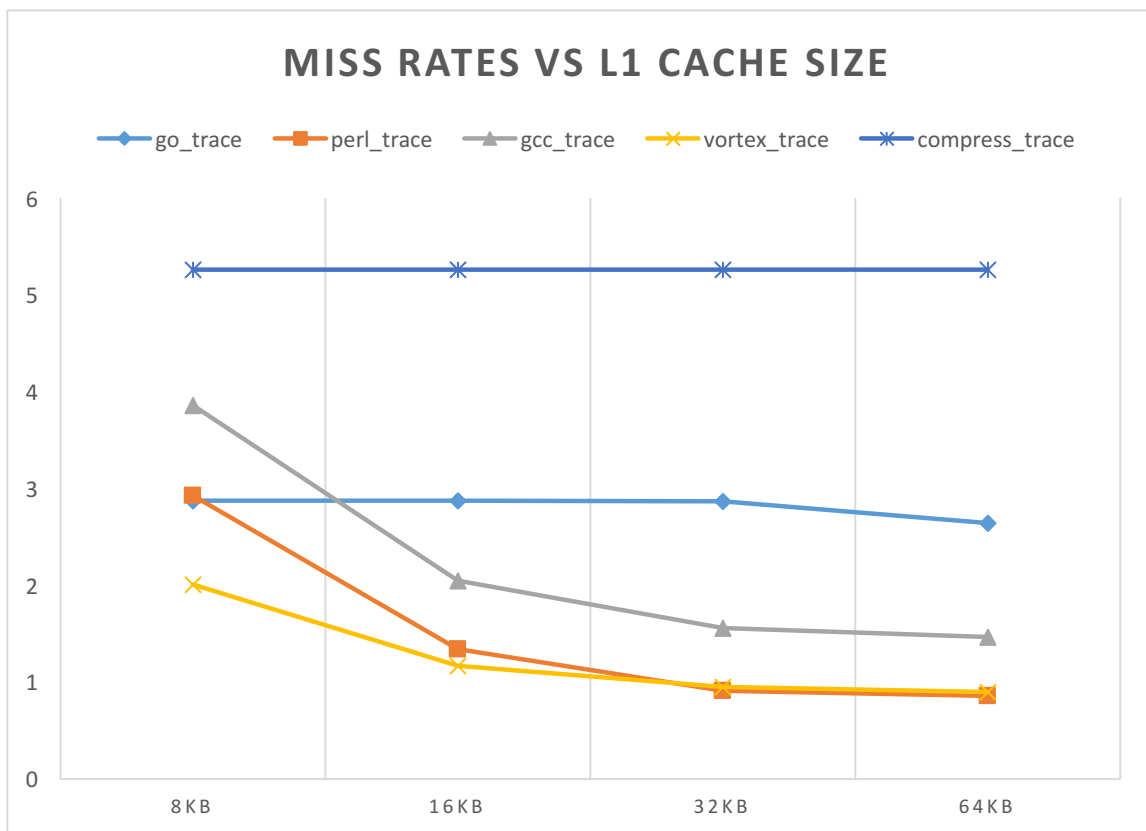
Course number: \_\_\_\_521\_\_\_\_  
(463 or 521 ?)

**Plot #1**

Assume that there is no L2 cache. Assume that LRU replacement policy is used for the L1 cache.

Fix the L1 block size at 64B and the L1 associativity at 4. Vary the cache size between 8KB, 16KB, 32KB and 64KB. Plot the L1 cache miss rates. The plot should have 4 data points per benchmark.

	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
8KB	2.874	2.929	3.859	2.01	5.26
16KB	2.871	1.337	2.045	1.166	5.26
32KB	2.869	0.91	1.558	0.951	5.26
64KB	2.645	0.86	1.466	0.899	5.26



**Observation:** Miss rates decreases for most of the benchmarks as we double the cache size from 8KB to 16KB.

Compress & go benchmark doesn't benefit much for the cache size variations.

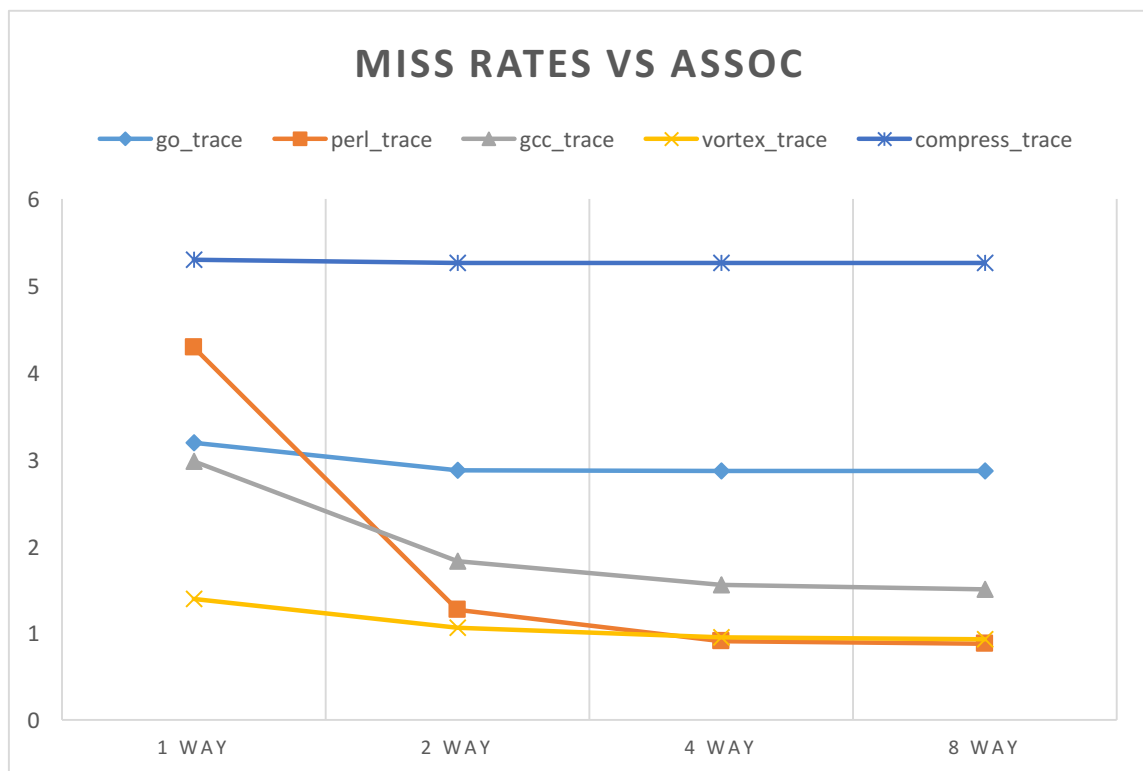
**Discussion:**

- Decrease in miss rates significant in the case of perl & gcc benchmark.
- In the case of go benchmark there is no significant decrease in miss rate.
- The most important reason for the decrease in miss rates is increasing the cache size decreases the capacity misses. Now the cache is large enough that it can hold most the data and evicts.
- This decrease in miss rates has achieves saturation after some size, as there it cannot change the cold miss as well as conflict misses.

**Plot #2:**

Fix the L1 block size at 64B and the L1 cache size at 32KB. Vary the associativity between 1, 2, 4 and 8. Plot the L1 cache miss rates. The plot should have 4 data points per benchmark.

	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
1 way	3.19	4.292	2.977	1.391	5.298
2 way	2.872	1.264	1.831	1.064	5.26
4 way	2.869	0.91	1.558	0.951	5.26
8 way	2.869	0.879	1.5	0.928	5.26



**Observation:** The miss rate goes on decreasing with the increase in the associativity and is applicable for all of the benchmarks

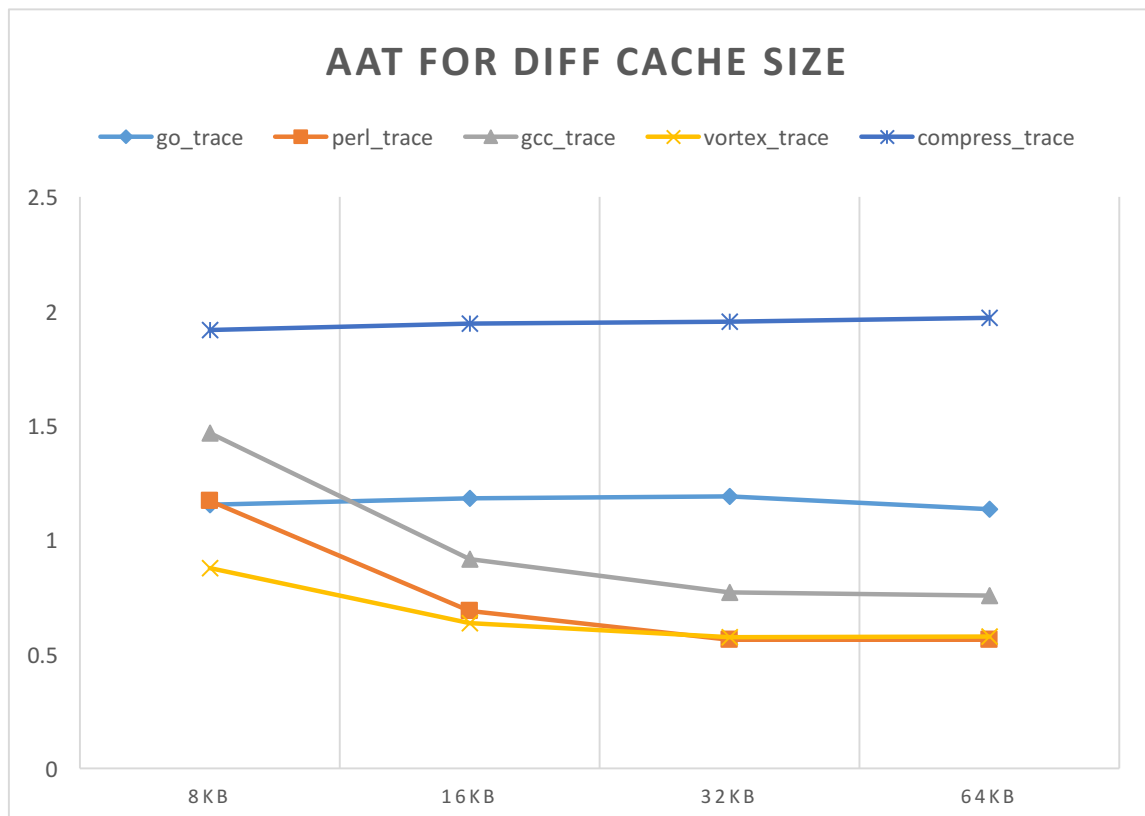
**Discussions:**

- Perl benchmark benefits the most with the increase in assoc from 1 to 2.
- Reason: Increasing the associativity decreases the conflict misses that otherwise happens due to two cachelines mapping to the same set, hence the existing cacheline is to be evicted.
- After 4 way assoc, we can also observe that increasing the assoc further doesn't affect much of the performance because, cold misses & capacity misses are not affected.

To do further analysis on the cache performance, Average access times are computed based on the data used from CACTI simulations.

We plot the Average access time for plot #1.

	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
8KB	1.130853	1.148453	1.446053	0.854373	1.894373
16KB	1.152656	0.661776	0.888336	0.607056	1.917136
32KB	1.18933	0.56245	0.76981	0.57557	1.95445
64KB	1.165881	0.594681	0.788601	0.607161	2.002681

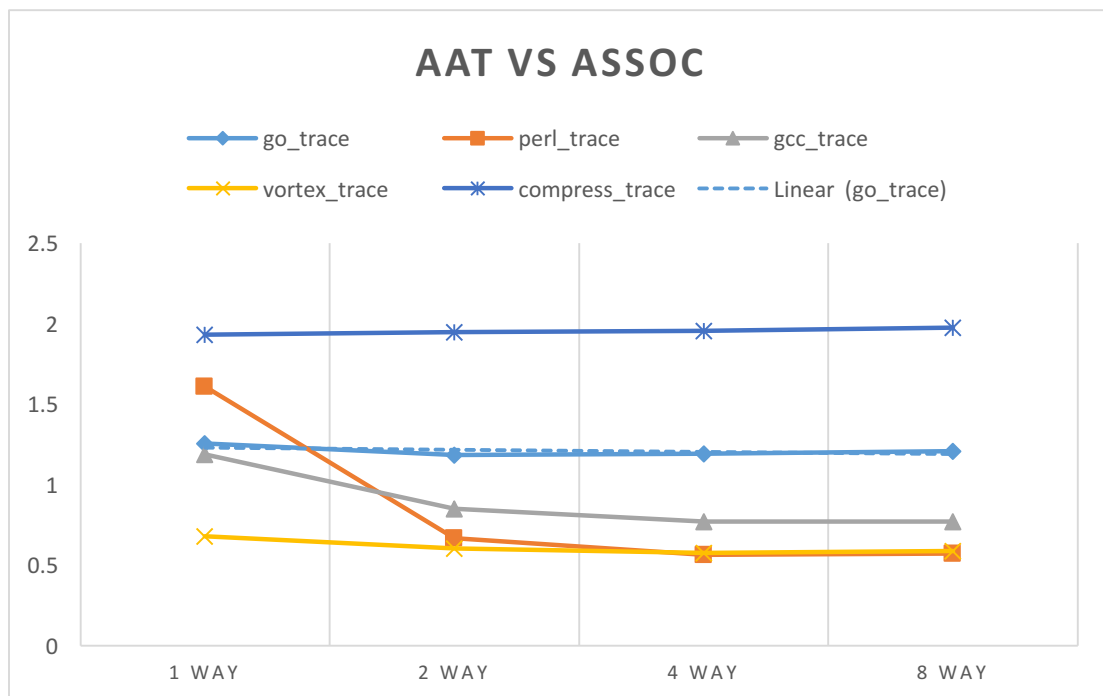


**Observation:**

- The key thing to notice is that as we increase the cache size, the access time goes on reducing approximately till the cache size of 32KB.
- With the increase of cache size from 32KB to 64KB, the access time increases a little bit.
- This trend is not followed by compress and go benchmarks.
- In both of the case, the average access time increases by a small amount.

We also computed the average access time required with the variation of associativity. Given below is the data-

	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
1 way	1.254153	1.606793	1.185993	0.678473	1.928713
2 way	1.181486	0.666926	0.848366	0.602926	1.945646
4 way	1.18933	0.56245	0.76981	0.57557	1.95445
8 way	1.206591	0.569791	0.768511	0.585471	1.971711

**Observation:**

- The key thing to notice is that as we increase the cache size, the access time goes on reducing approximately till the cache size of 4-way.
- With the increase of cache size from 4-way to 8-way, the access time increases a little bit.

- In case of Compress benchmark, the average access time increases with the increase in assoc.

**Implications:**

- Since we know that there was a miss rate decrease for every increase in the cache size, but that doesn't imply that the cache average access time decreases. (Compress benchmark).

**Area consideration:**

	1	2	4	8
8	0.053293238	0.083756164	0.068434156	0.102584886
16	0.096748995	0.130107044	0.105941693	0.130444675
32	0.210543576	0.205554649	0.236647681	0.242170635
64	0.330469394	0.350242085	0.30228937	0.360317611

The area increases as we increase the size of the cache as well as if we increase the assoc of the cache. Taking this under account & considering the fact that smaller area (die area) results lower cost, we should consider cache structure such that it has as small as necessary for greater performance.

**Recommendation:**

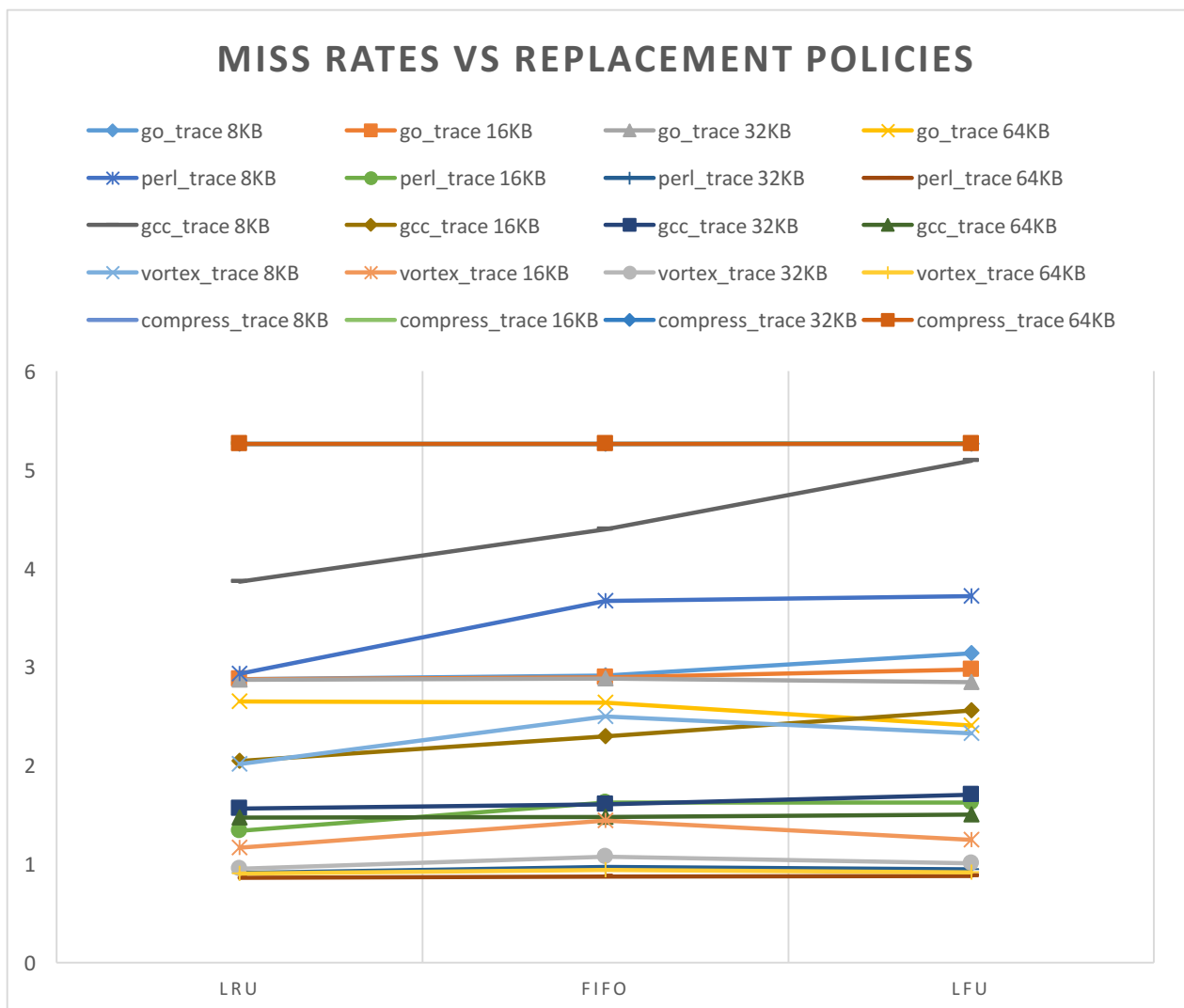
- Considering the area values given in CACTI table & the plots #1 & #2, we will recommend 32KB cache size with 2-way assoc.
- If we are running & considering cache architecture for processor based only on compress benchmark, then we should recommend 8KB 1-way cache as miss rates/AAT remains almost same for all the different configuration. And the fact that the area required is less by a significant amount (approx. 0.05).
- If we are running & considering cache architecture for processor based only on go benchmark, then we would recommend 8KB 4-way cache as miss rates/AAT remains almost same for all the different configuration. And the fact that the area required is less by a significant amount (approx. 0.07).
- If we are running & considering cache architecture for processor based only on perl, vortex, gcc benchmark, then we would recommend 32KB 4-way cache as miss rates/AAT remains almost same for all the different configuration.

**Plot #3**

Fix the L1 block size at 64B and the L1 associativity at 4. Vary cache size between 8KB, 16KB, 32KB and 64KB. For each configuration vary the replacement policy between LRU, FIFO and LFU. Plot the L1 Cache miss rates for each replacement policy. Note

that for each benchmark this plot will have 3 data series (one per replacement policy) and each data series will have 4 data points (one per cache size). That is, 12 data points per benchmark.

	go_trace			perl_trace			gcc_trace			vortex_trace			compress_trace		
	LRU	FIF	LFU	LRU	FIF	LFU	LRU	FIF	LFU	LRU	FIF	LFU	LR	FIF	LFU
8KB	2.87	2.90	3.13	2.92	3.66	3.71	3.85	4.39	5.08		2.49	2.32			5.26
	4	8	6	9	4	4	9	4	7	2.01	4	4	5.26	5.26	5
16K	2.87	2.88	2.97	1.33	1.62	1.62	2.04	2.29	2.55	1.16	1.43	1.24			5.26
B	1	9	3	7	3	4	5	3	6	6	8	1	5.26	5.26	4
32K	2.86	2.87	2.84		0.96	0.94	1.55	1.60	1.70	0.95	1.07	1.00			
B	9	8	2	0.91	8	4	8	1	2	1	5	8	5.26	5.26	5.26
64K	2.64	2.63	2.40		0.86	0.87	1.46	1.47	1.49	0.89	0.93	0.91			
B	5	7	4	0.86	9	9	6	6	7	9	5	3	5.26	5.26	5.26



**Observation:**

- The miss rate is maximum for LFU policy, for most of the benchmarks.
- LRU behave sufficiently good for a large number of benchmarks and cache configuration.
- Hardware implementation of LFU is relatively complex.
- For smaller cache size (typically 8KB), miss rates for LFU are higher than that of LRU & FIFO
- For larger cache size (typically 32-63KB), miss rates are comparable for all of the replacement policies.
- In case of go benchmark, 64KB L1 has lower miss rates for LFU than LRU/FIFO.

**Recommendation:** We should choose LRU as our replacement policy as it reduces the miss rates.

**Plot #4**

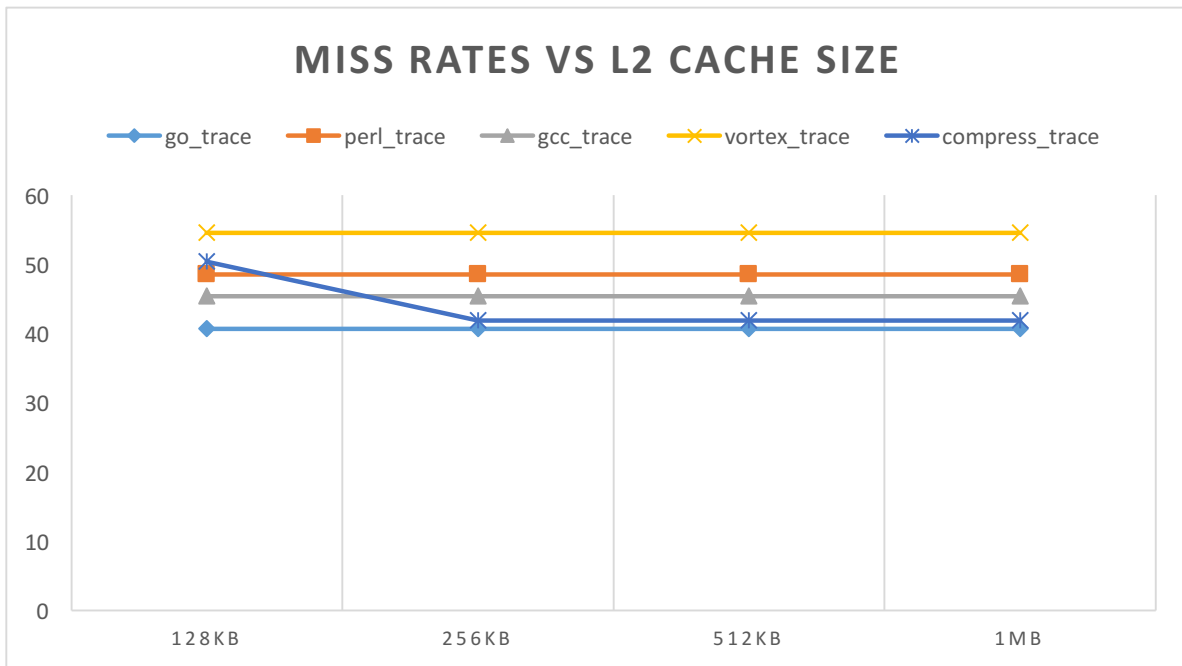
Fix block size at 64B, L1 cache size at 16KB, L1 associativity at 4 and the replacement policy to be LRU. Assume that the L2 Cache is always non-inclusive.

Fix L2 associativity at 8. Vary the L2 cache size between 128KB, 256KB, 512KB and 1MB. Plot the L2 cache miss rates. The plot should have 4 data points per benchmark.

Following is the data points based on the experiments:

	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
128KB	40.7289	48.6056	45.446	54.562	50.4195
256KB	40.7289	48.6056	45.446	54.562	41.885
512KB	40.7289	48.6056	45.446	54.562	41.885
1MB	40.7289	48.6056	45.446	54.562	41.885

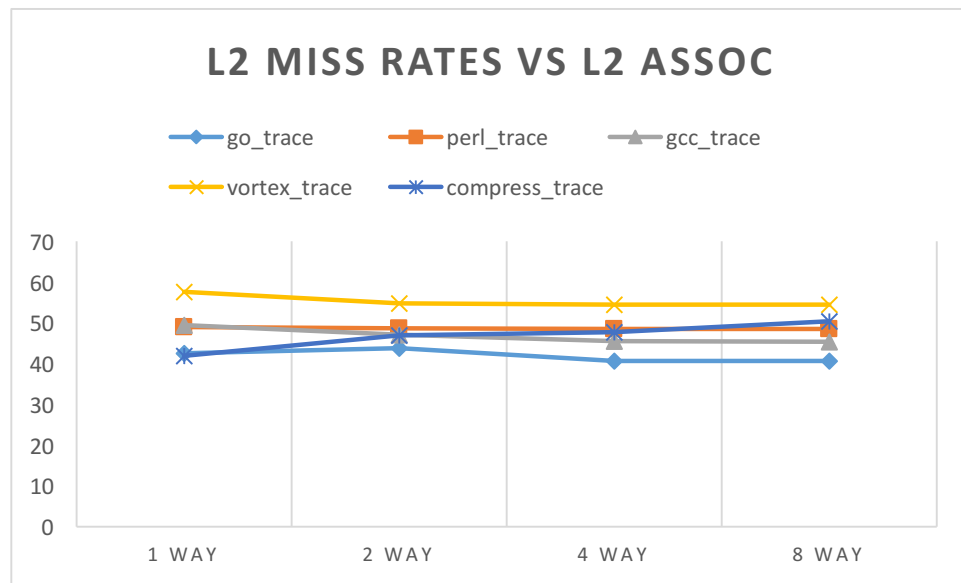


**Observation:**

- Miss rates usually remains the same for most of the benchmark.
- In case of compress benchmark, miss rate gets lowered with the increase in the cache size.

**Plot #5**

	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
1 way	42.511	49.0609	49.4836	57.6642	41.9058
2 way	43.7725	48.7763	47.1674	54.8054	47.0119
4 way	40.7289	48.6056	45.4773	54.562	47.768
8 way	40.7289	48.6056	45.446	54.562	50.4195

**Observation:**

- Miss rates usually decreases as we increase the assoc for most of the benchmarks.
- In case of compress benchmark, miss rate increases with increase in the L2 assoc.

**Average access times analysis:**

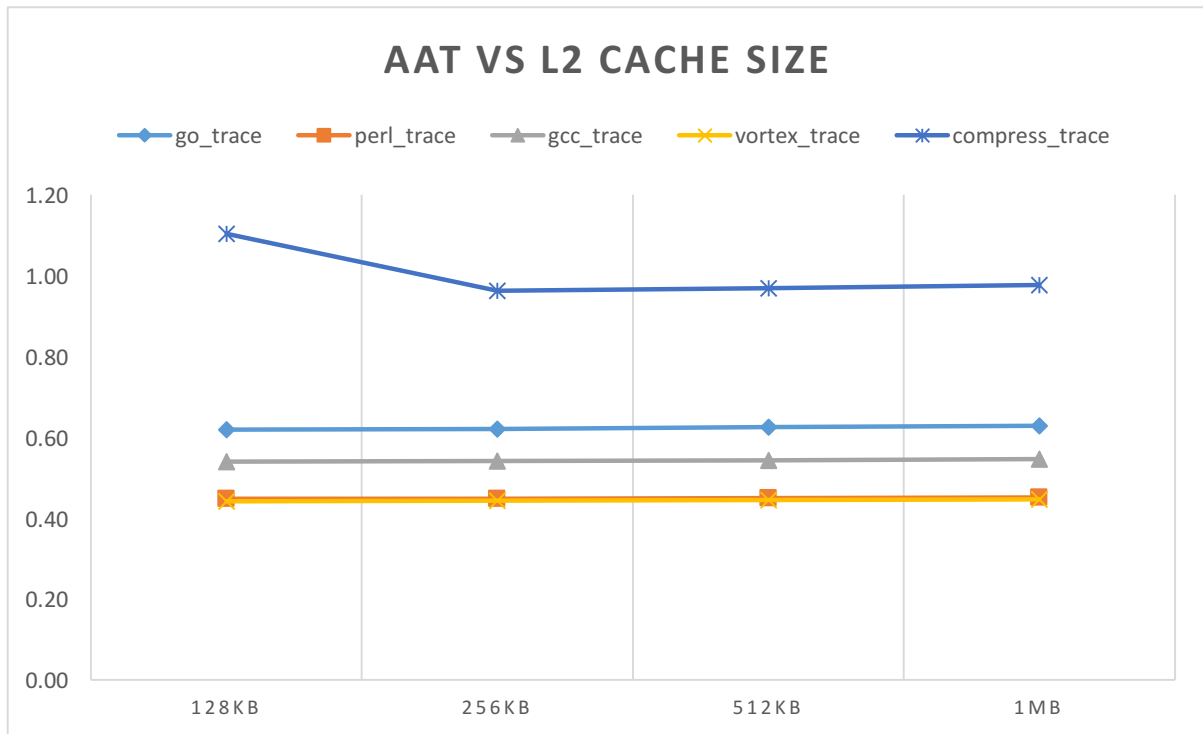
To do further analysis on the cache performance, Average access times are computed based on the data used from CACTI simulations.

$$\text{AATL2} = \text{HitTimeL1} + (\text{MissRateL1} * \text{HitTimeL2}) + (\text{MissRateL1} * \text{MissRateL2} * \text{MissPenalty})$$

*All of the static values are taken from the CACTI simulator sheet.*

The following is the computed AAT for different configuration taken for plot #4.

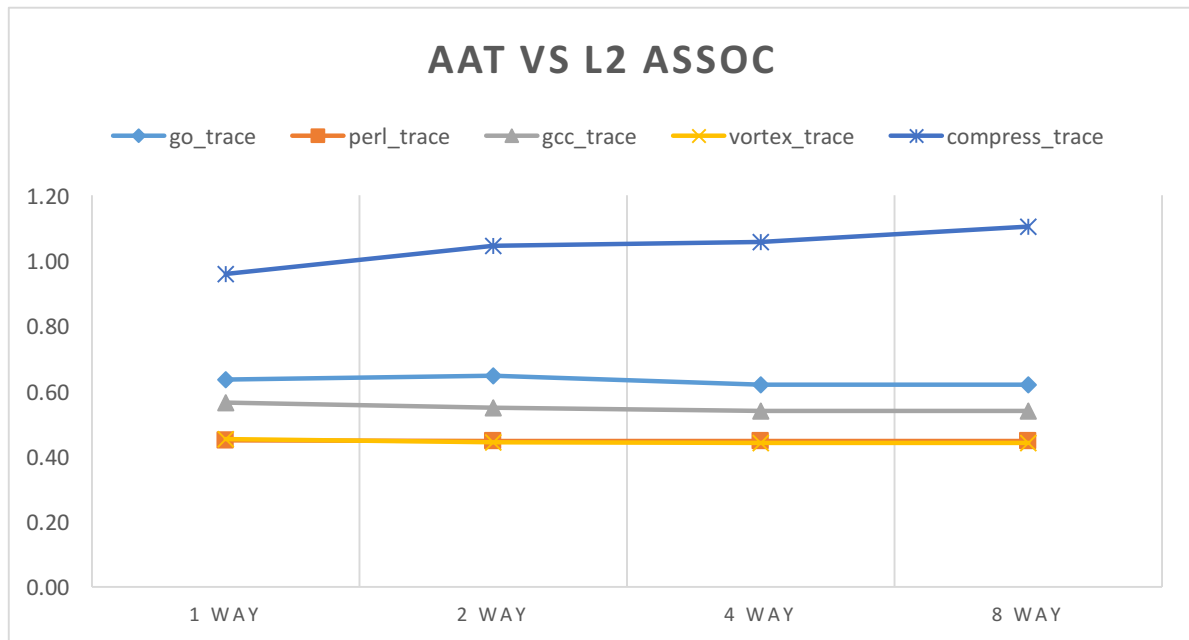
	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
128KB	0.62	0.45	0.54	0.44	1.10
256KB	0.62	0.45	0.54	0.44	0.96
512KB	0.62	0.45	0.54	0.44	0.97
1MB	0.63	0.45	0.55	0.45	0.98

**Observations:**

- As we vary the size of L2 cache, there is no significant difference in the average access time for most of the benchmarks.
- In case of compress benchmark, we see a large drop in the AAT.

Similarly, we compute the AAT by for different cache configuration taken under consideration in plot #5

	go_trace	perl_trace	gcc_trace	vortex_trace	compress_trace
1 way	0.64	0.45	0.57	0.45	0.96
2 way	0.65	0.45	0.55	0.44	1.04
4 way	0.62	0.45	0.54	0.44	1.06
8 way	0.62	0.45	0.54	0.44	1.10

**Observations:**

- As we vary the assoc of L2 cache, there is no significant difference in the average access time for most of the benchmarks.
- In case of compress benchmark, we see a gradual but significant increase in the average access time.

**Area considerations:**

The area increases as we increase the size of the cache as well as if we increase the assoc of the cache. Taking this under account & considering the fact that smaller area (die area) results lower cost, we should consider cache structure such that it has as small as necessary for greater performance. Following is the plot that specifies area for different cache configuration.

*These area values will help in choosing/recommending cache structure.*

	128KB	256KB	512KB	1MB
1	0.657687839	1.277806565	2.487864981	3.747960026
2	0.694245498	1.562127169	2.225827274	4.349252234
4	0.667017966	1.141294802	2.177361671	4.673162925
8	0.559933334	1.293540536	2.640142073	4.874201405

**Recommendation:**

- If we are running & considering cache architecture for processor, then we should recommend 128KB 8-way cache as miss rates/AAT remains almost same for all

the different configuration. And the fact that the area required is less by a significant amount (approx. 0.055).

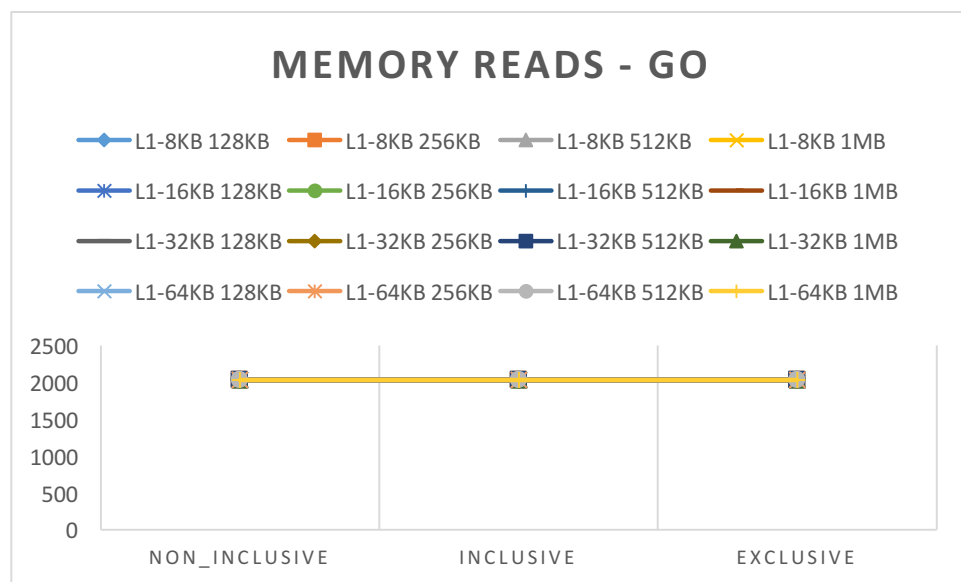
- Miss rate is of 48% for perl benchmark.
- AAT is of 0.45.

### Plot #6

Plot the number of cache misses resulting in a memory read for each configuration for each inclusion policy (inclusive, exclusive and non-inclusive). Note that for each benchmark this plot will have 3 data series (one per inclusion policy) and each data series will have 16 points (one per configuration).

We study the read misses for L2 cache size

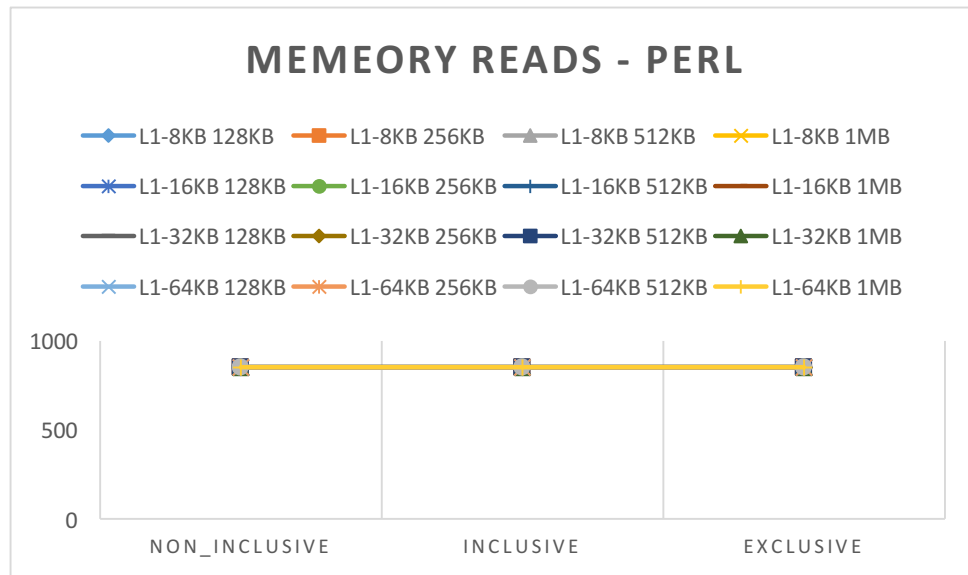
### Plot #6(a) – Go Benchmark



### Observation:

- The memory reads remain the same for all configuration.

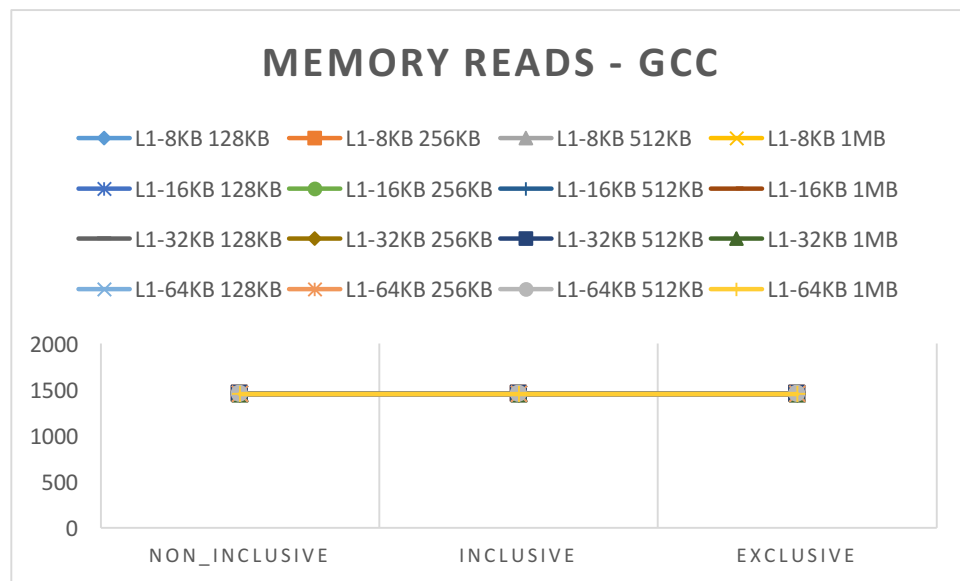
### Plot #6(a) – Perl Benchmark



**Observation:**

- The memory reads remain the same for all configuration.

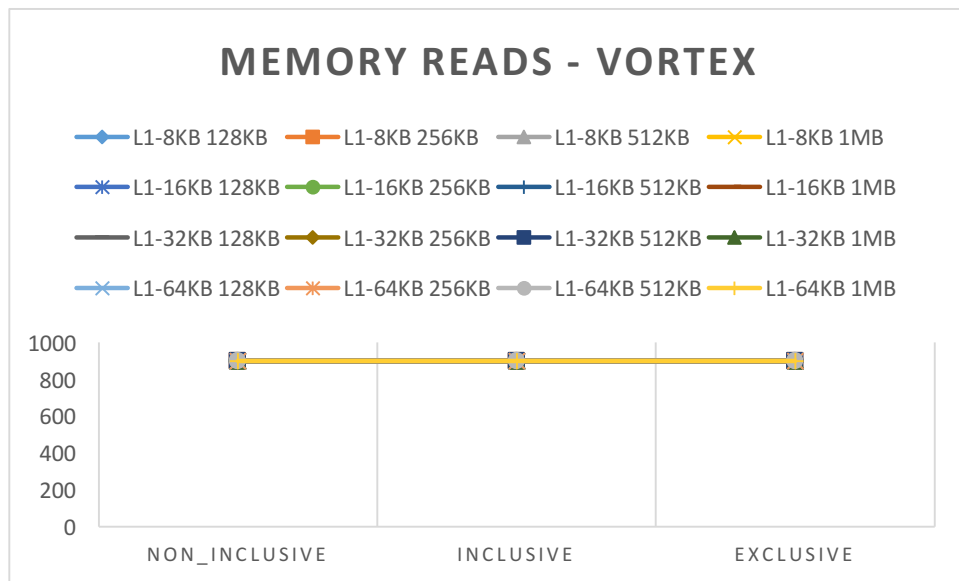
**Plot #6(a) – GCC Benchmark**



**Observation:**

- The memory reads remain the same for all configuration.

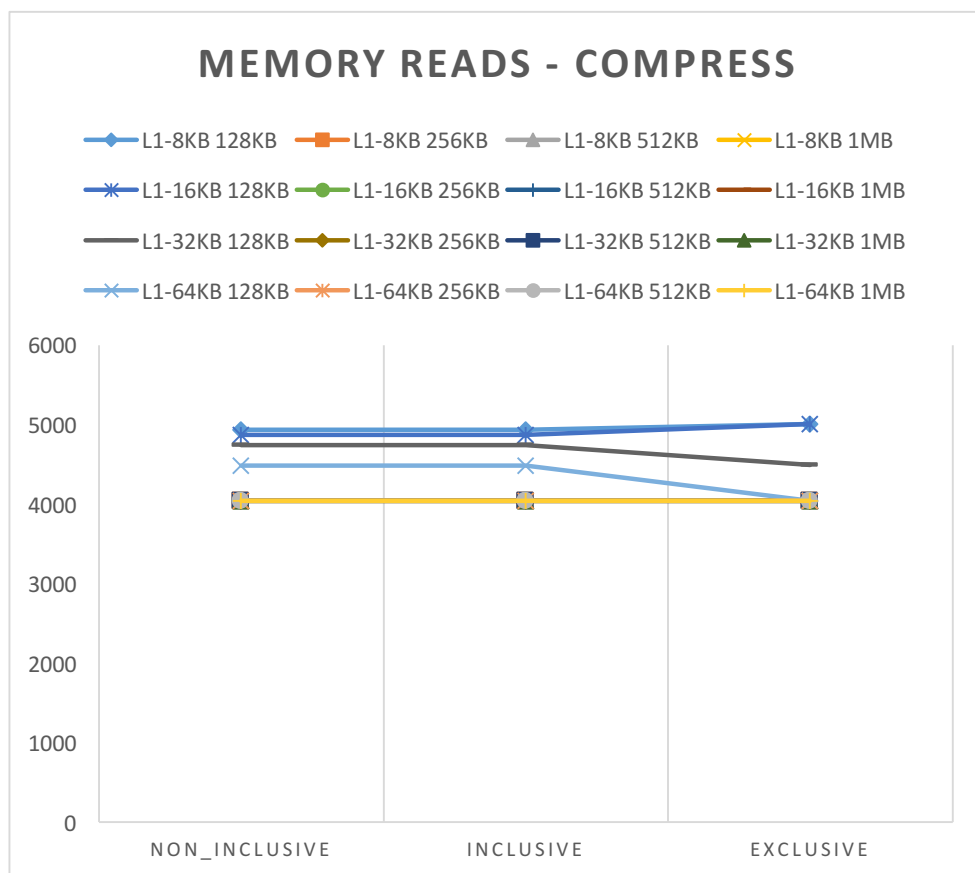
**Plot #6(a) – Vortex Benchmark**



**Observation:**

- The memory reads remain the same for all configuration.

**Plot #6(a) – Compress Benchmark**



**Observation:**

- In case of larger L1 (typically 32/64KB), and smaller L2 (typically 128KB), the memory reads are less in the case of Exclusive caches.
- In case of smaller L1 (typically of 8/16KB) and smaller L2 (128KB), the memory reads are more in case of Exclusive cache as compared to Inclusive or Non-Inclusive.
- The Memory reads remain the same in case of Inclusive & Non-Inclusive cache.

**Implications:**

- In case of Exclusive, since L1 content cannot be present in L2, L2 can have large amount of unique data.
- If the size of L1 is large and L2 is small, then in case of Non-Inclusive & Inclusive, much of the L2 will be taken by the data present in L1. This may result in eviction of some useful data from L2. (Useful meaning the data that may be referenced later).
- This results in higher number of memory reads for Non-Inclusive & Inclusive cache.

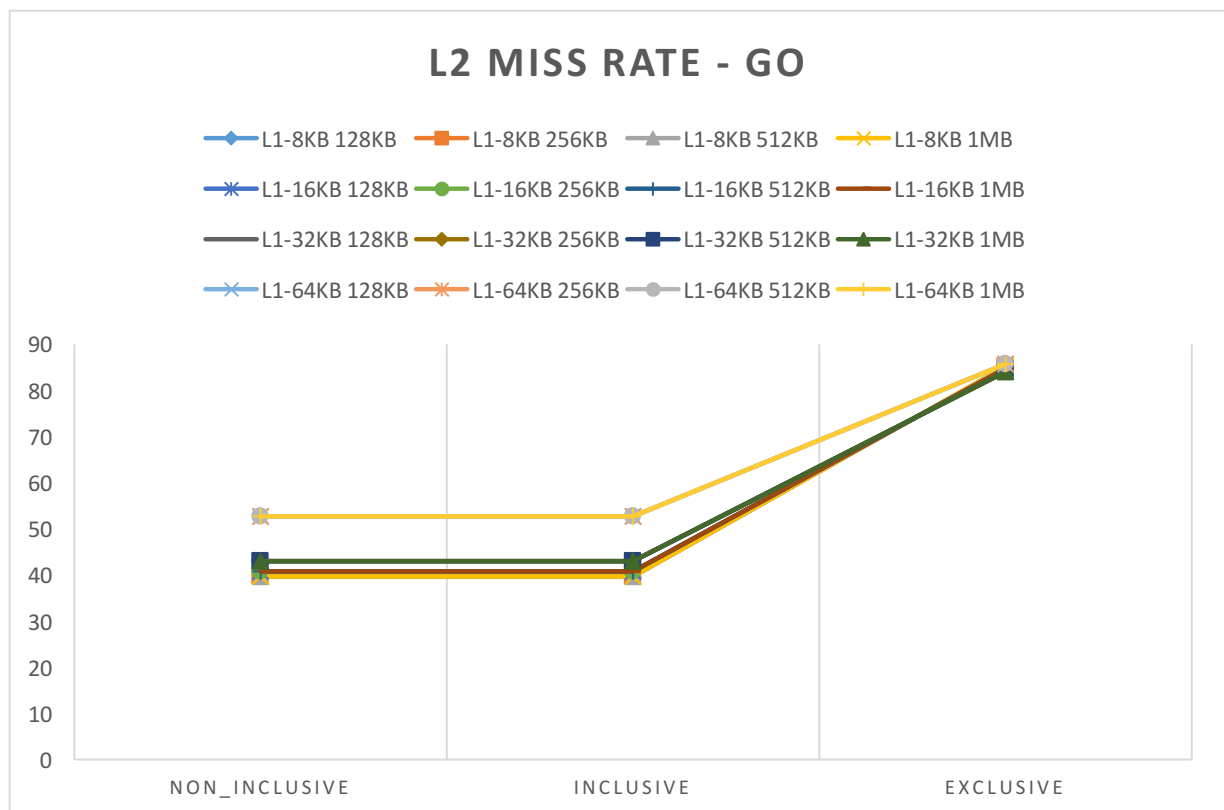


**Plot #7**

Assume that the AAT for Exclusive caches can be calculated similarly to the AAT calculation for non-inclusive caches. Even if blocks brought in from memory are directly installed in the L1 cache (no insertion time at L2 cache), we still encounter the hit latency of L2 cache as a result of checking the L2 cache before reading from memory. Plot the AAT when using non-inclusive, inclusive and exclusive L2 caches for each configuration. The number of data points will be similar to Plot 6.

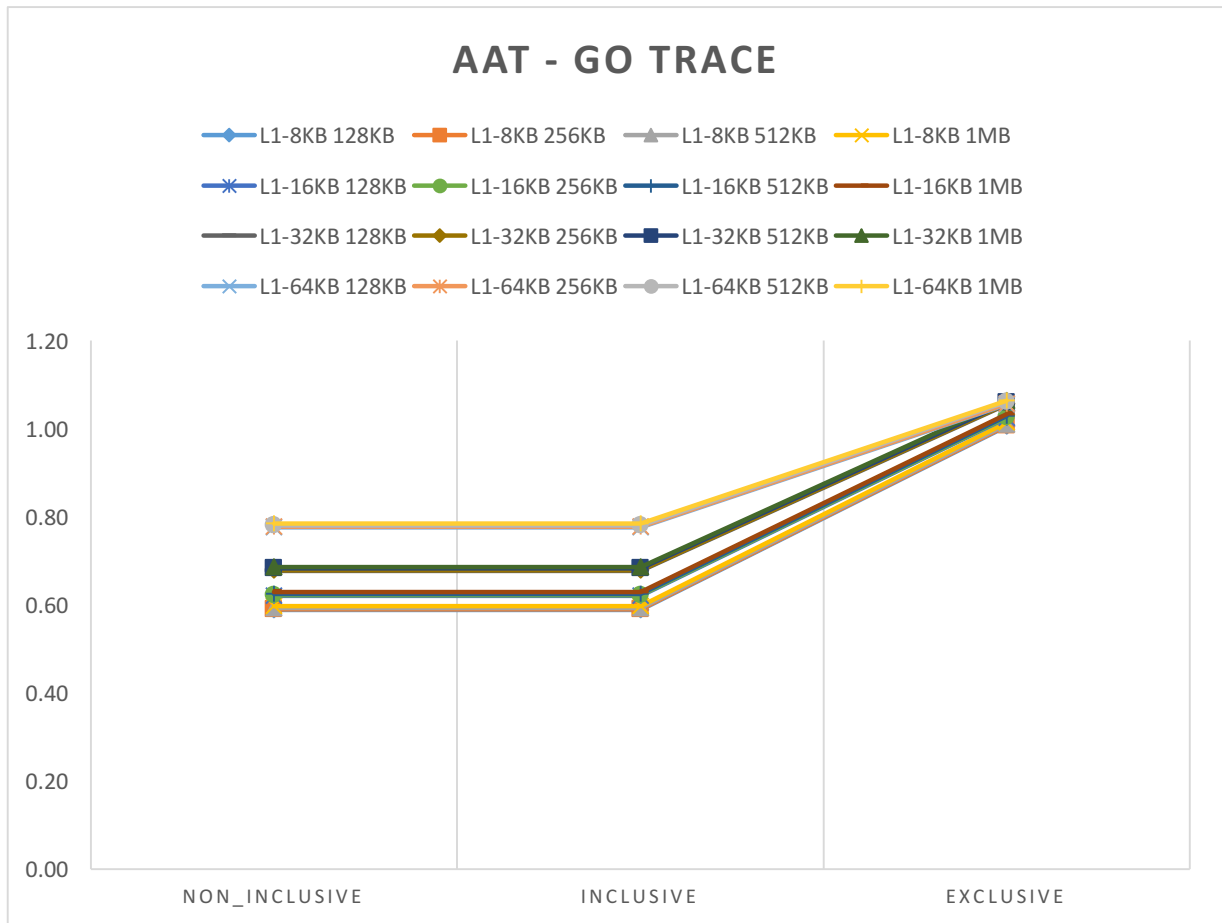
For this experiment we consider different configuration given in the spec sheet. The miss rates for these configuration is tabulated as follows.

*For each of the benchmark, plot #1 would refer to the miss rates for different configuration whereas, plot #2 would refer to the AAT required. For each of the benchmark we would make sufficient observation and finally we would recommend a cache configuration that will be better in terms of configuration. Note: we would not consider area required for the cache.*

**Plot #7(a) – Go benchmark**

**Observation:**

- Miss rate remains the same for the case of Non-inclusive and Inclusive cache configuration.
- Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches.
- As we increase the size of L1, the L2 miss rates increases by a significant factor.

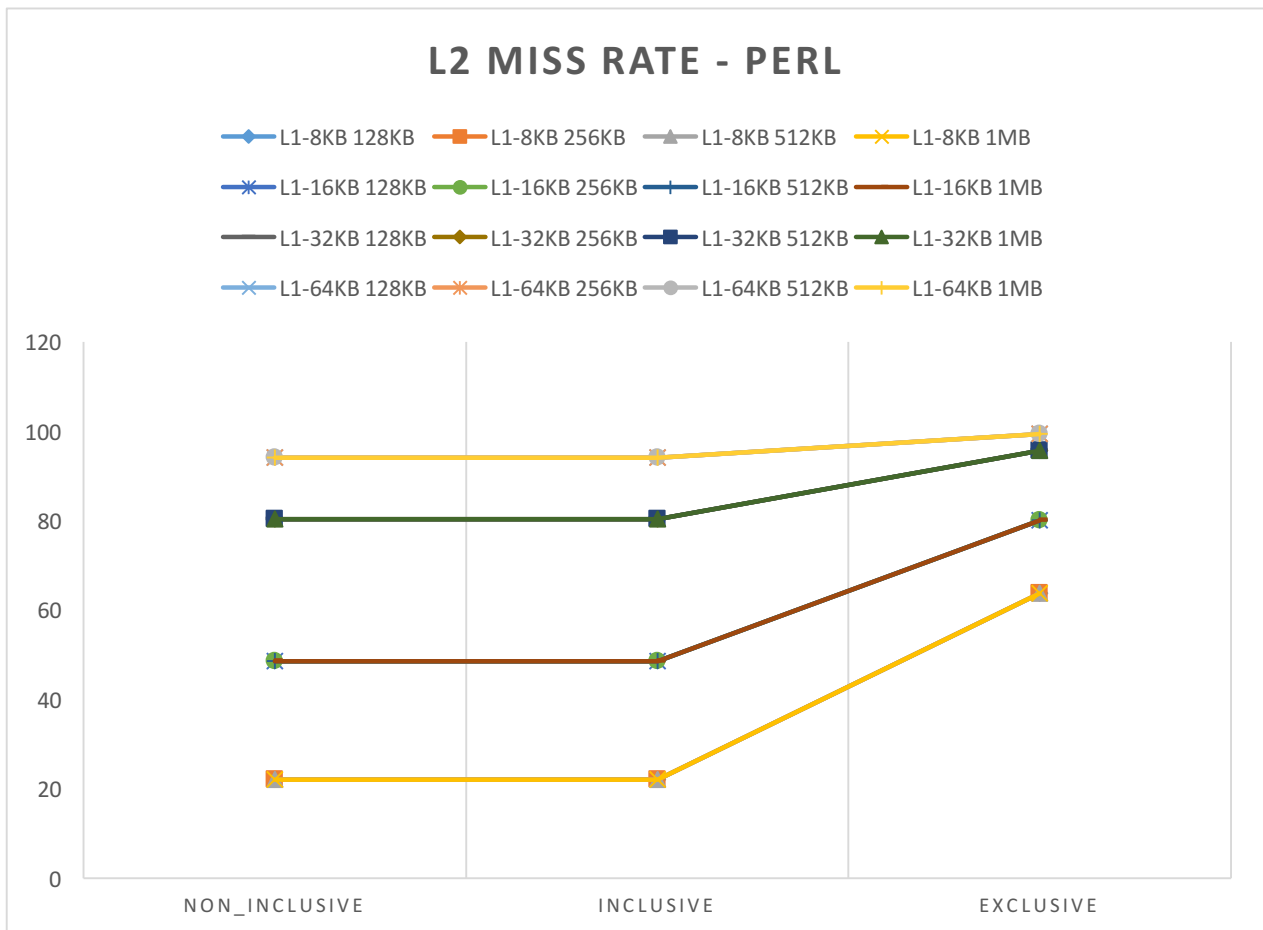
**Observation:**

- AAT remains the same for the case of Non-inclusive and Inclusive cache configuration mainly because both have same miss rates.
- Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches, hence the average access time is also greater.

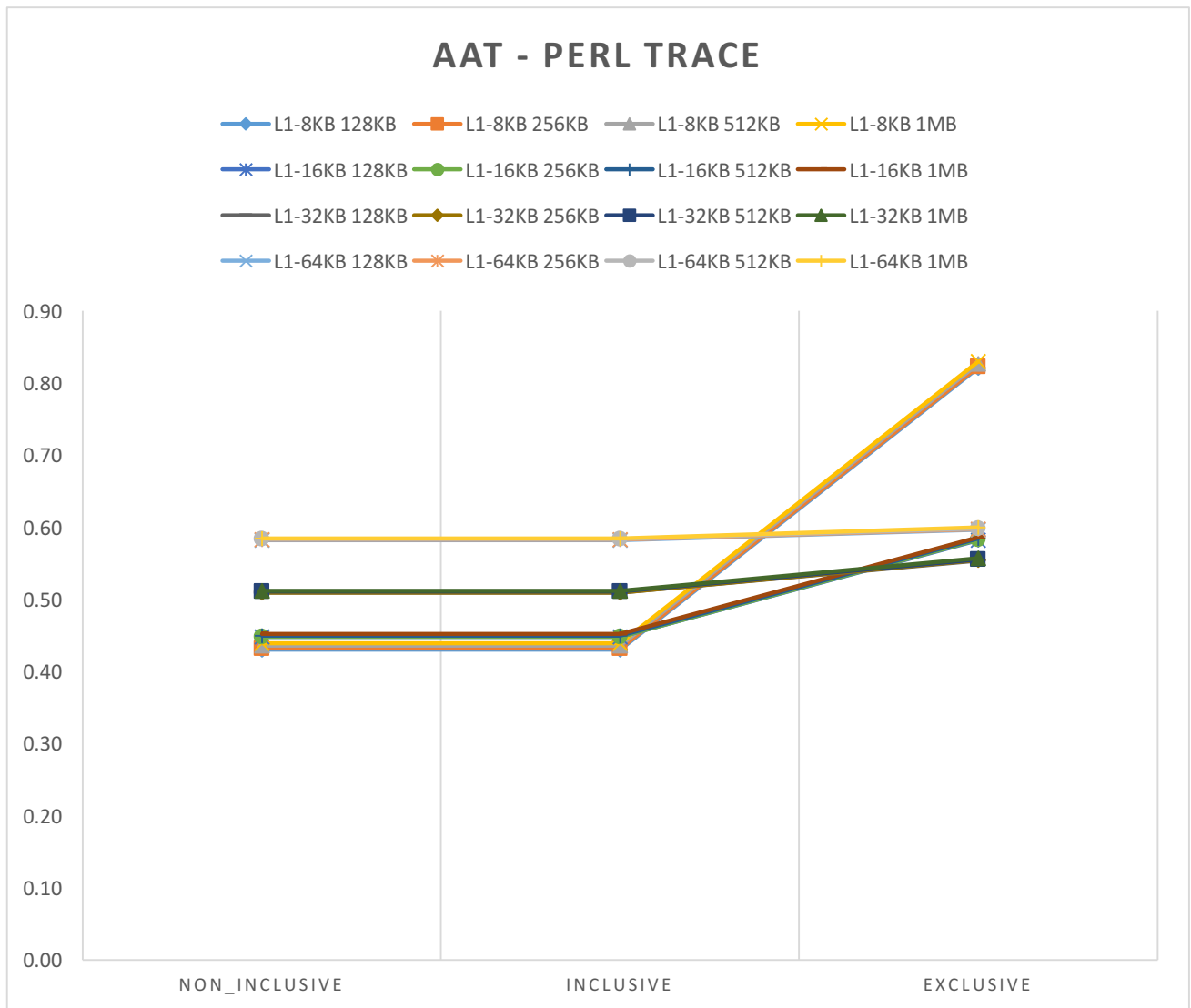
**Plot #7(b) – Perl benchmark**

perl_trac				
e	L1-8KB	L1-16KB	L1-32KB	L1-64KB

	128K	256K	512K		128K	256K	512K		128K	256K	512K		128K	256K	512K	
	B	B	B	1MB	B	B	B	1MB	B	B	B	1MB	B	B	B	1MB
NON_IN C	22.10	22.10	22.10	22.1 0	48.61	48.61	48.61	48.6 1	80.34	80.34	80.34	80.3 4	94.05	94.05	94.05	94.0 5
INCL	22.10	22.10	22.10	22.1 0	48.61	48.61	48.61	48.6 1	80.34	80.34	80.34	80.3 4	94.05	94.05	94.05	94.0 5
EXCL	63.79	63.79	63.79	63.7 9	80.02	80.02	80.02	80.0 2	95.72	95.72	95.72	95.7 2	99.39	99.39	99.39	99.3 9

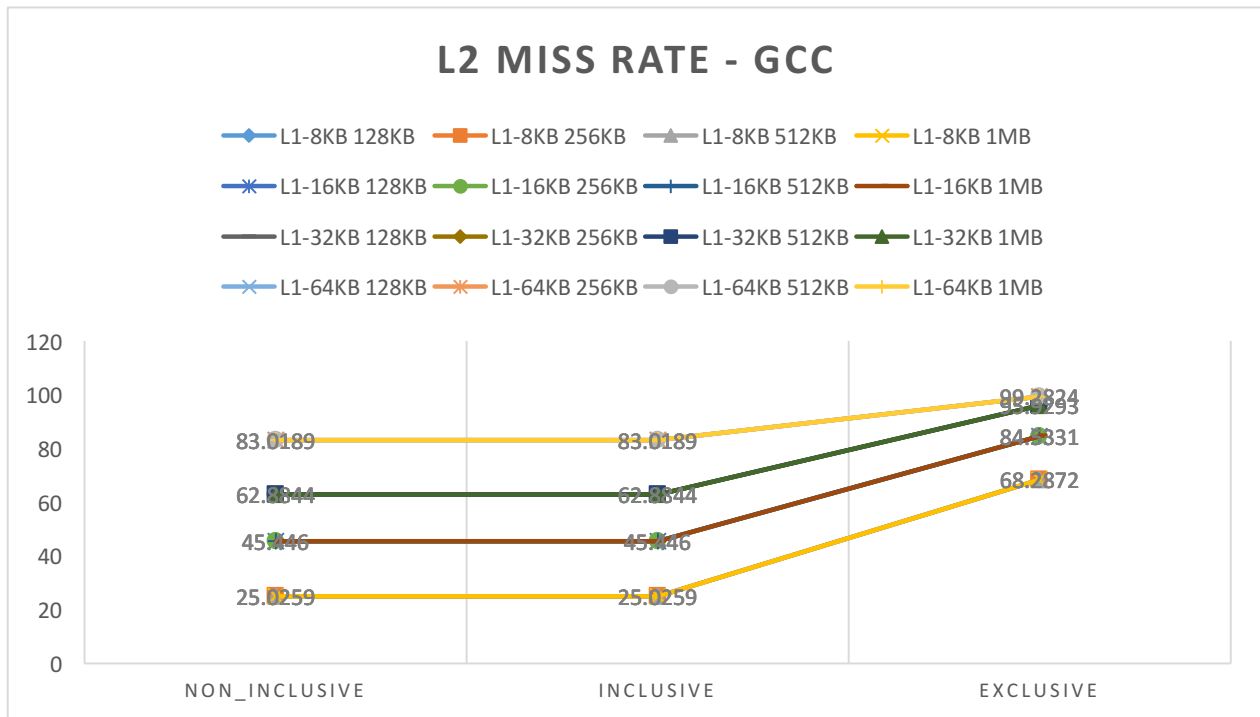
**Observation:**

- Miss rate remains the same for the case of Non-inclusive and Inclusive cache configuration.
- Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches.
- As we increase the size of L1, the L2 miss rates increases by a significant factor.

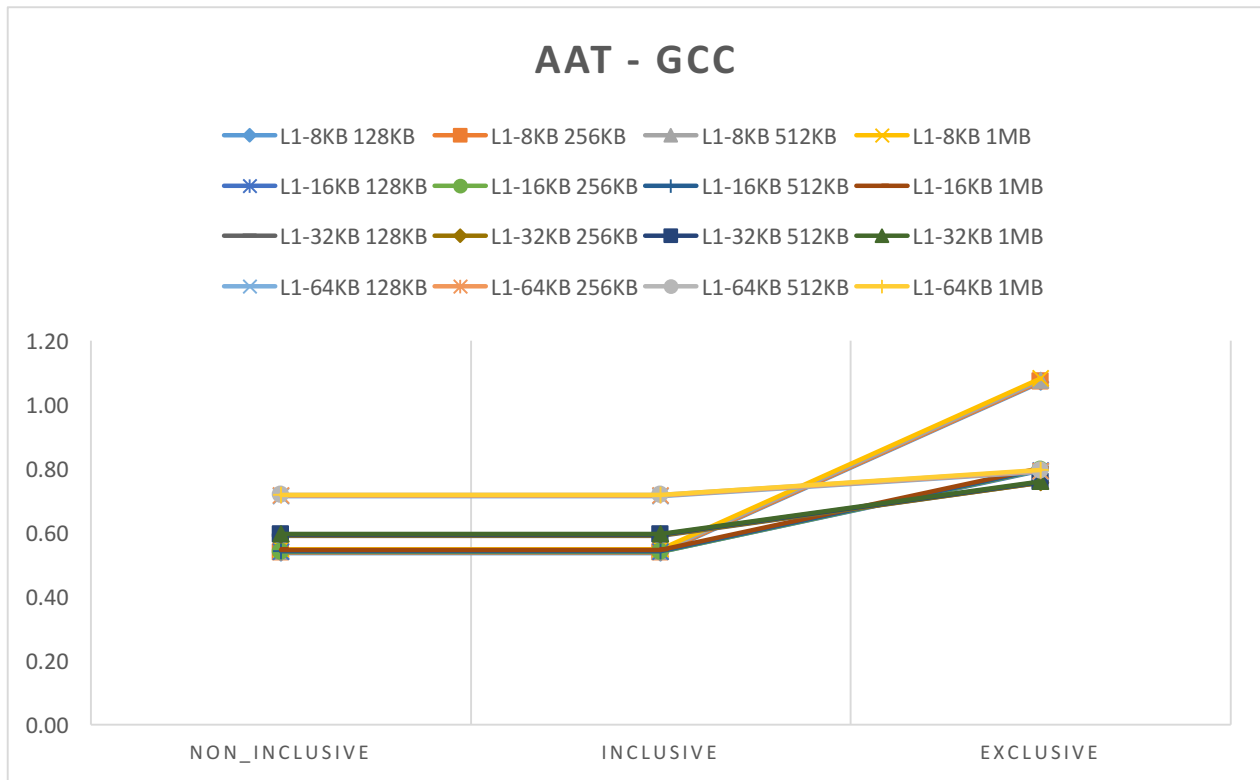


Plot #7(c) – GCC benchmark

gcc_trace	L1-8KB				L1-16KB				L1-32KB				L1-64KB			
	128 KB	256 KB	512 KB	1M B	128 KB	256 KB	512 KB	1M B	128 KB	256 KB	512 KB	1M B	128 KB	256 KB	512 KB	1M B
NON_INC	25.0	25.0	25.0	25.0	45.4	45.4	45.4	45.4	62.8	62.8	62.8	62.8	83.0	83.0	83.0	83.0
LUSIVE	259	259	259	259	46	46	46	46	844	844	844	844	189	189	189	189
INCLUSI	25.0	25.0	25.0	25.0	45.4	45.4	45.4	45.4	62.8	62.8	62.8	62.8	83.0	83.0	83.0	83.0
VE	259	259	259	259	46	46	46	46	844	844	844	844	189	189	189	189
EXCLUSI	68.2	68.2	68.2	68.2	84.5	84.5	84.5	84.5	95.9	95.9	95.9	95.9	99.2	99.2	99.2	99.2
VE	872	872	872	872	331	331	331	331	293	293	293	293	824	824	824	824

**Observation:**

- Miss rate remains the same for the case of Non-inclusive and Inclusive cache configuration.
- Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches.
- In case of 8KB L1 non-inclusive/inclusive cache outperform exclusive caches by a factor of approx. 2
- As we increase the size of L1, the L2 miss rates increases by a significant factor.

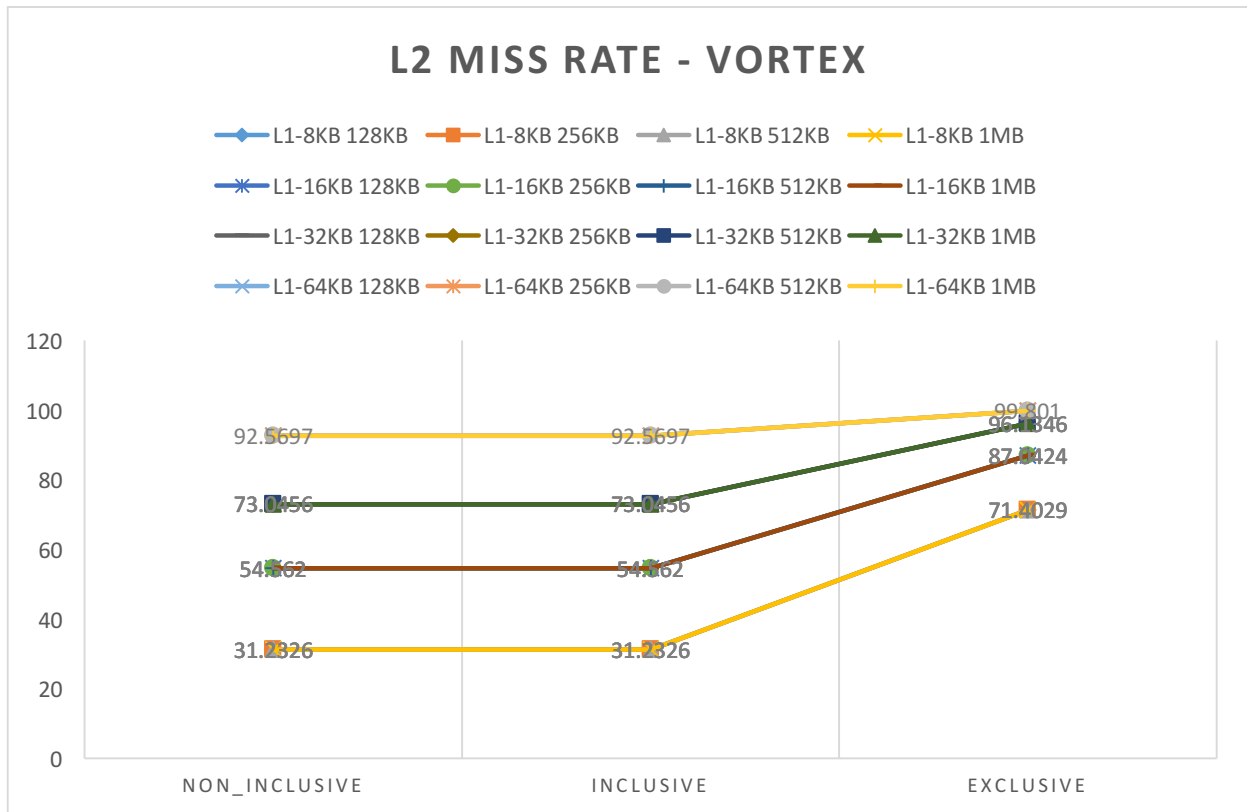
**Observation:**

- AAT remains the same for the case of Non-inclusive and Inclusive cache configuration mainly because both have same miss rates.

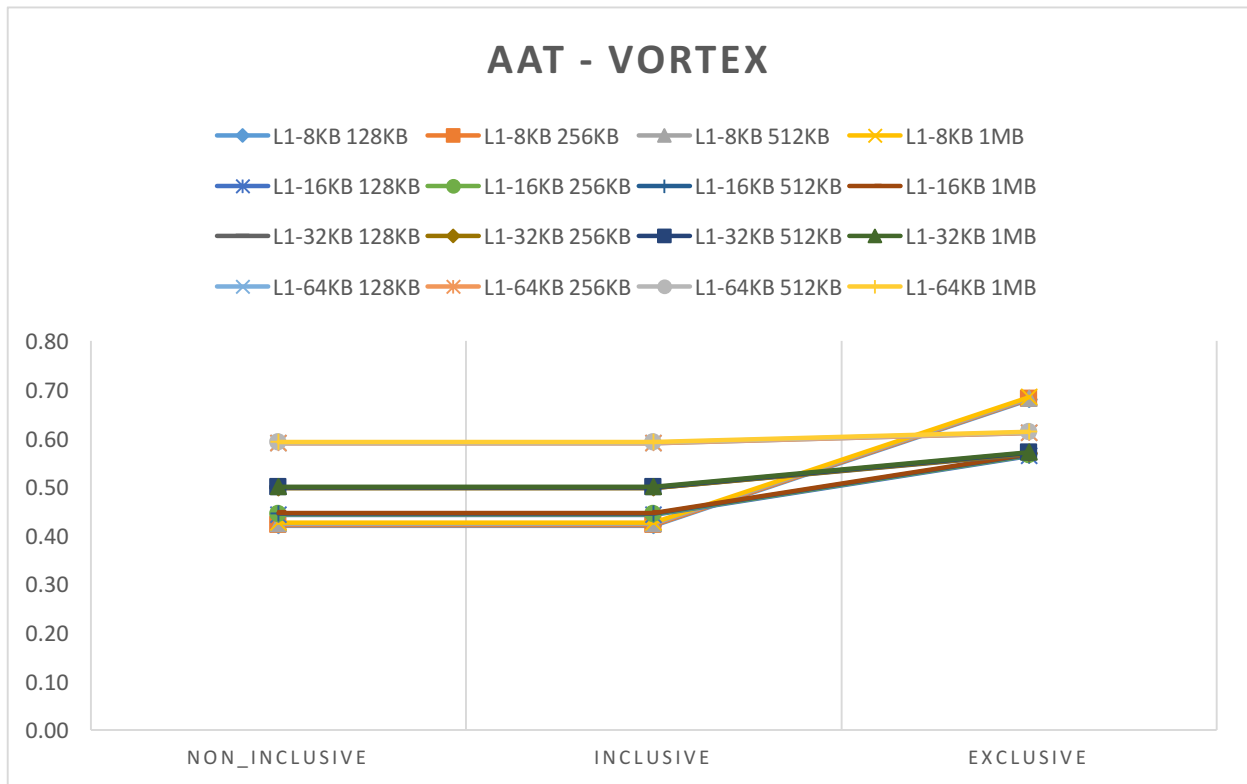
Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches, hence the average access time is also greater.

**Plot #6(d) & #7(d) – Vortex benchmark**

vortex_tra	L1-8KB				L1-16KB				L1-32KB				L1-64KB			
c	128	256	512	1M	128	256	512	1M	128	256	512	1M	128	256	512	1M
	KB	KB	KB	B	KB	KB	KB	B	KB	KB	KB	B	KB	KB	KB	B
NON_INC	31.2	31.2	31.2	31.2	54.5	54.5	54.5	54.5	73.0	73.0	73.0	73.0	92.5	92.5	92.5	92.5
LUSIVE	326	326	326	326	62	62	62	62	456	456	456	456	697	697	697	697
INCLUSI	31.2	31.2	31.2	31.2	54.5	54.5	54.5	54.5	73.0	73.0	73.0	73.0	92.5	92.5	92.5	92.5
VE	326	326	326	326	62	62	62	62	456	456	456	456	697	697	697	697
EXCLUSI	71.4	71.4	71.4	71.4	87.0	87.0	87.0	87.0	96.1	96.1	96.1	96.1	99.8	99.8	99.8	99.8
VE	029	029	029	029	424	424	424	424	346	346	346	346	01	01	01	01

**Observation:**

- Miss rate remains the same for the case of Non-inclusive and Inclusive cache configuration.
- Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches.
- In case of 8KB L1 & 1MB L2, non-inclusive/inclusive cache outperform exclusive caches by a factor of approx. 2
- As we increase the size of L1, the L2 miss rates increases by a significant factor.
- For 64 KB L1, the miss rate of L2 is close to 90%.

**Observation:**

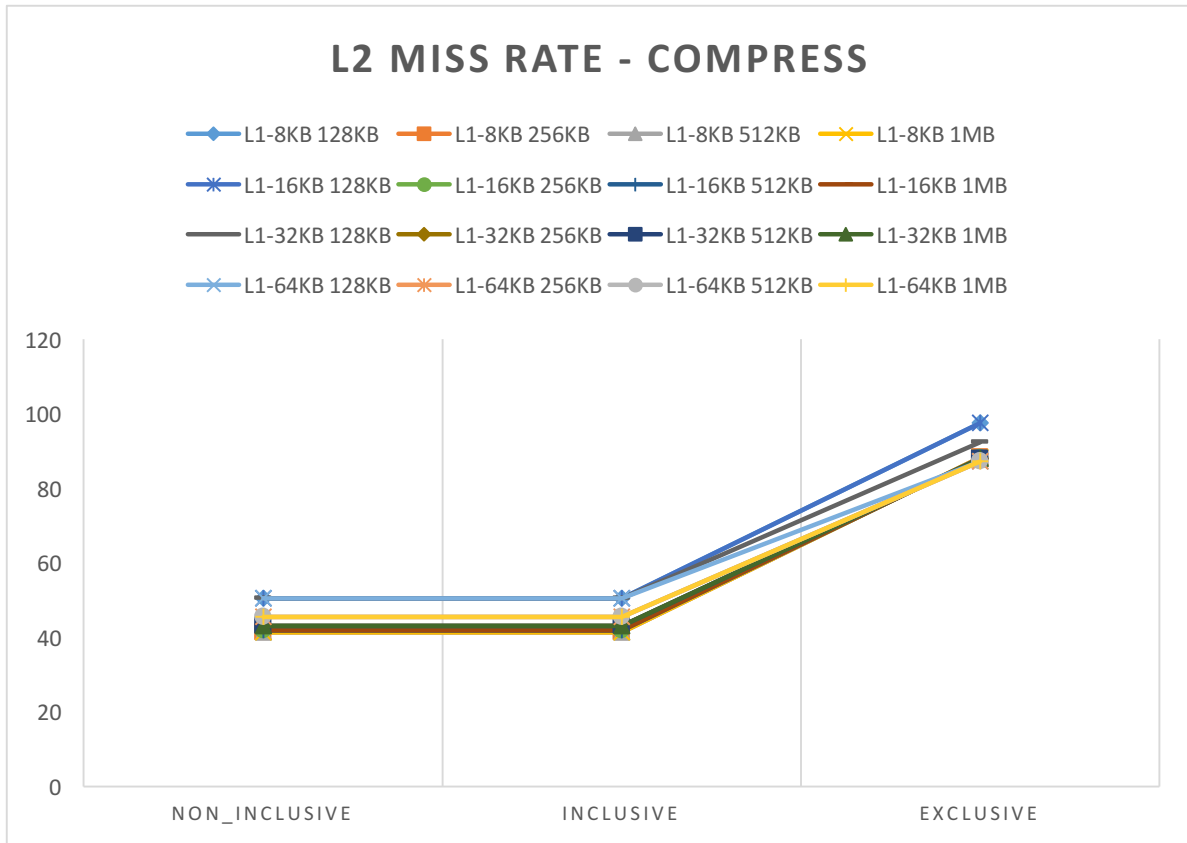
- AAT remains the same for the case of Non-inclusive and Inclusive cache configuration mainly because both have same miss rates.

Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches, hence the average access time is also greater.

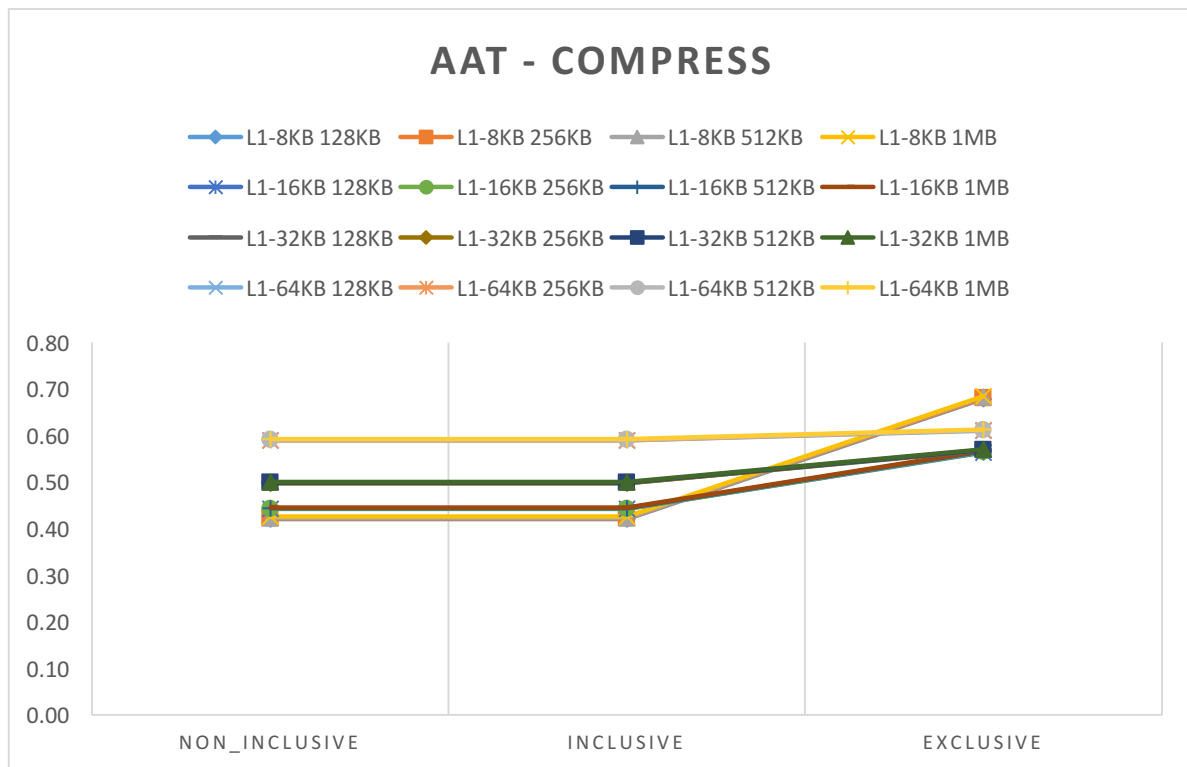
**Plot #6(e) & #7(e) – Compress benchmark**

compress	L1-8KB				L1-16KB				L1-32KB				L1-64KB			
	128 KB	256 KB	512 KB	1M B	128 KB	256 KB	512 KB	1M B	128 KB	256 KB	512 KB	1M B	128 KB	256 KB	512 KB	1M B
NON_INC	50.4	41.3	41.3	41.3	50.4	41.8	41.8	41.8	50.4	43.0	43.0	43.0	50.4	45.5	45.5	45.5
LUSIVE	14	37	37	37	195	85	85	85	309	259	259	259	557	047	047	047
INCLUSI	50.4	41.3	41.3	41.3	50.4	41.8	41.8	41.8	50.4	43.0	43.0	43.0	50.4	45.5	45.5	45.5
VE	14	37	37	37	195	85	85	85	309	259	259	259	557	047	047	047
EXCLUSI	97.5	88.2	88.2	88.2	97.5	88.1	88.1	88.1	92.3	87.8	87.8	87.8	87.1	87.1	87.1	87.1
VE	558	987	987	987	253	528	528	528	461	497	497	497	946	946	946	946



**Observation:**

- Miss rate remains the same for the case of Non-inclusive and Inclusive cache configuration.
- Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches.
- Miss rates becomes almost double in case of Exclusive cache.

**Observation:**

- AAT remains the same for the case of Non-inclusive and Inclusive cache configuration mainly because both have same miss rates.

Exclusive cache configuration results in higher cache miss rate as compared to that of Non-inclusive & Inclusive caches, hence the average access time is also greater.

**Implications:**

- The general observation is that the miss rate is less in the case of non inclusive caches & inclusive caches for wide range of benchmarks.
- Miss rate is same for inclusive as well as non-inclusive cache.
- In case of exclusive caches, the miss rate is higher.
- Average access time required is also greater in case of exclusive caches.
- As we increase the size of L1, miss rates in L1 cache decreases, but miss rate in L2 increase.
- This may be explained by the fact that L1 is able to address most of the request except for the cold misses. These accesses are compulsory miss and will also miss in L2.
- It can be argued that misses in L2 may have significant amount of cold misses in case of 64KB L1.

**Non Inclusive vs Inclusive:**

- Inclusive cache requires back invalidations, which in-turn put pressure on L1

- Non-Inclusive cache on the other hand doesn't require any back-invalidation, hence there is no extra traffic from L2 to L1.
- Inclusive caches will be better in the case where there are multiple processors and each of them have their own local cache (Not shared). If there is cache coherence problem, it just has to look into the lower level of the cache.
- Since we are considering multiprocessor configuration, this deduction is out of scope.
- In case of single processor configuration, as we have seen based on our experiments, the miss rates are same.
- Consider a case where L1 cache is large and L2 cache is small.
- Inclusion property specifies that content of L1 should be always be present in L2. So effectively L2 can hold less unique data. In case of smaller L2 cache, this usually creates a problem. So we might see a larger L2 cache miss rates. We should recommend in using Non-Inclusive cache instead.
- Or in other words: Non-Inclusive cache will outperform Inclusive cache for this case.

**Exclusive vs Non-Inclusive:**

- Consider the case where the L1 cache size is large enough (typically of 64KB) while the size of L2 is small (typically of 128 KB). In the case of Non-Inclusive, any read in L1 would result in bringing the block in L1 as well as L2.
- In case of Exclusive cache, any request in L1 would only result in bringing the block in L1 and not in L2.
- Since, L2 doesn't have the blocks that were present in L1, L2 in turn has large space for unique data.
- Since L2 is small & L2 may contain blocks that were present in L1, any eviction in L1 would result in allotting a space in L2. Now if there is a capacity miss in L2 on the block that was evicted this would result in memory read.
- This capacity miss would less probably occur in case of Exclusive as L2 can hold whole 128 KB of unique data. This means there is more space in L2.
- In summary, we should recommend using Exclusive cache in the case of smaller size of L2 with relatively larger L1.