

Assignment Project Exam Help

<https://eduassistpro.github.io/>

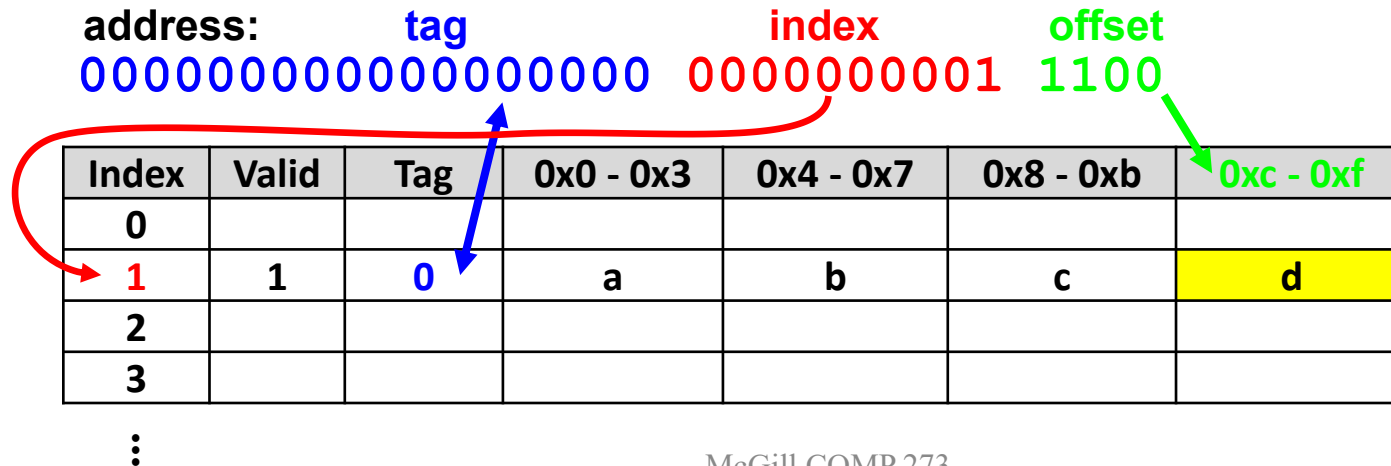
Add WeChat edu_assist_pro

Review

- We would like to have the capacity of disk at the speed of the processor: unfortunately this is not feasible
- So we create a memory hierarchy:
 - each successively loads data from next lower level
 - exploits temporal locality
 - do the common case fast, worry less about the exceptions (design principle of MIPS)
- Locality of reference is a Big Idea

Big Idea Review

- Mechanism for transparent movement of data among levels of a storage hierarchy
 - set of address/value bindings
 - address provides
 - compare deduplicates
 - service hit or miss
 - load new block and binding



Outline

- Block Size Tradeoff
- Types of Cache Misses
- Fully Associative Cache
- N-Way Associative Cache
- Block Replacement Policy
- Multilevel Caches (if time)
- Cache write policy (if time)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Block Size Tradeoff (1/3)

- Benefits of Larger Block Size

- Spatial Locality: if we access a given word, we're likely to access other nearby words soon (Another Big Idea)
- Very applicable with <https://eduassistpro.github.io/> if we execute a given instruction, it's likely that we'll execute a few as well
- Works nicely in sequential array accesses too

Block Size Tradeoff (2/3)

- Drawbacks of Larger Block Size
 - Larger block size means larger miss penalty
 - on a miss, takes longer time to load a new block from next level
 - If block size is too big, then there are too few blocks
 - Result: miss rate goes up
- In general, minimize
Average Access Time
 - = Hit Time + Miss Penalty x Miss Rate

Block Size Tradeoff (3/3)

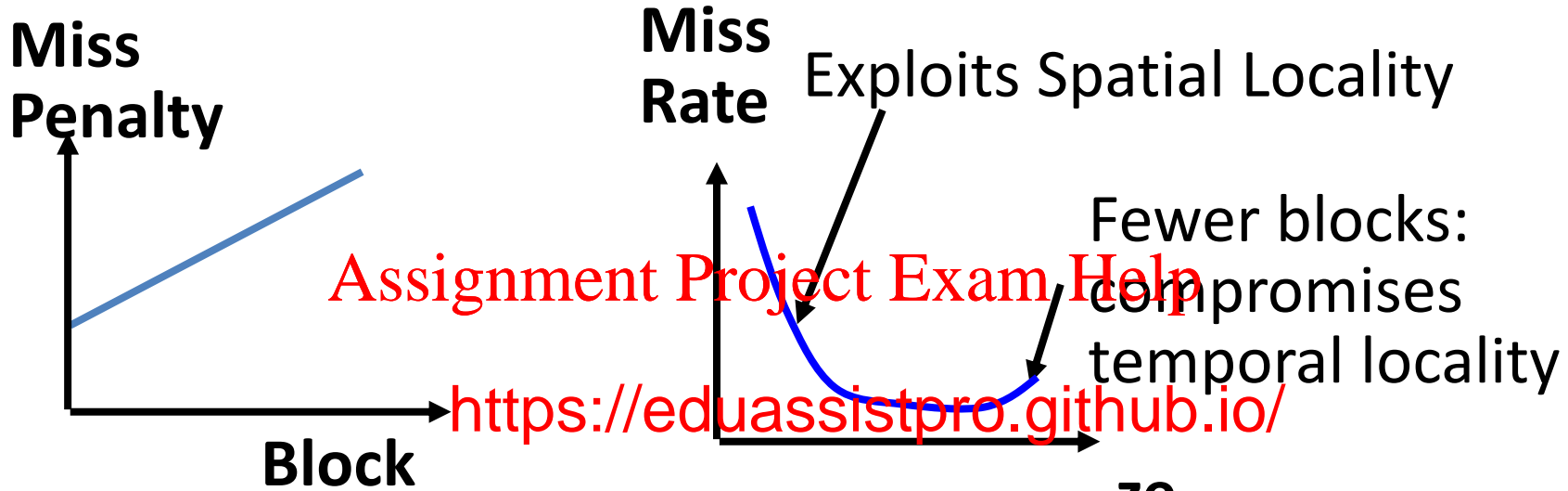
- Hit Time = time to find and retrieve data from current level cache
- Miss Penalty = average time to find data on a current level miss (includes the penalty for searching in successive levels of memory hierarchy)
- Hit Rate = % of requests that are found in current level cache
- Miss Rate = $1 - \text{Hit Rate}$

Assignment Project Exam Help

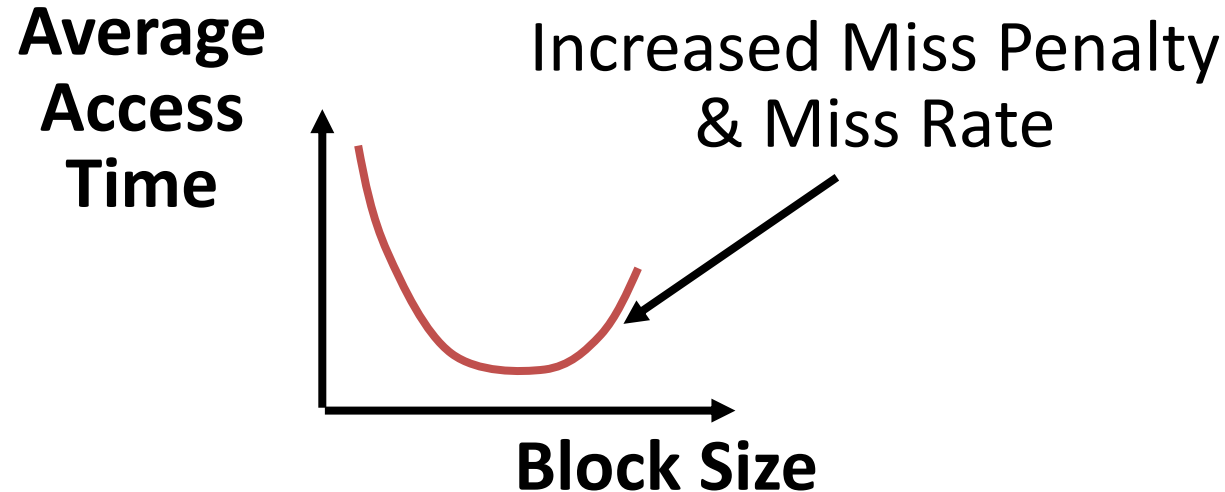
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Block Size Tradeoff Conclusions



Add WeChat edu_assist_pro



Types of Cache Misses (1/2)

- Compulsory Misses

- occur when a program is first started
- cache does not contain program's data yet, so misses are bound to occur
- can't be avoided easily, so won't focus on this course

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Types of Cache Misses (2/2)

- Conflict Misses

- miss that occurs because two distinct memory addresses map to the same cache location
- two blocks (which have the same location) can keep overwriting each other
- big problem in direct-mapped caches
- how do we lessen the effect of these?

Dealing with Conflict Misses

- Solution 1: Make the cache size bigger
 - relatively expensive
- Solution 2: Multiple Index? it in the same Cache

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Fully Associative Cache (1/3)

- Memory address fields:
 - Tag: same as before
 - Offset: same as befo
 - Index: non-existent
- What does this mean?
 - any block can go anywhere in the cache
 - must compare with all tags in entire cache to see if data is there

Assignment Project Exam Help

<https://eduassistpro.github.io/>

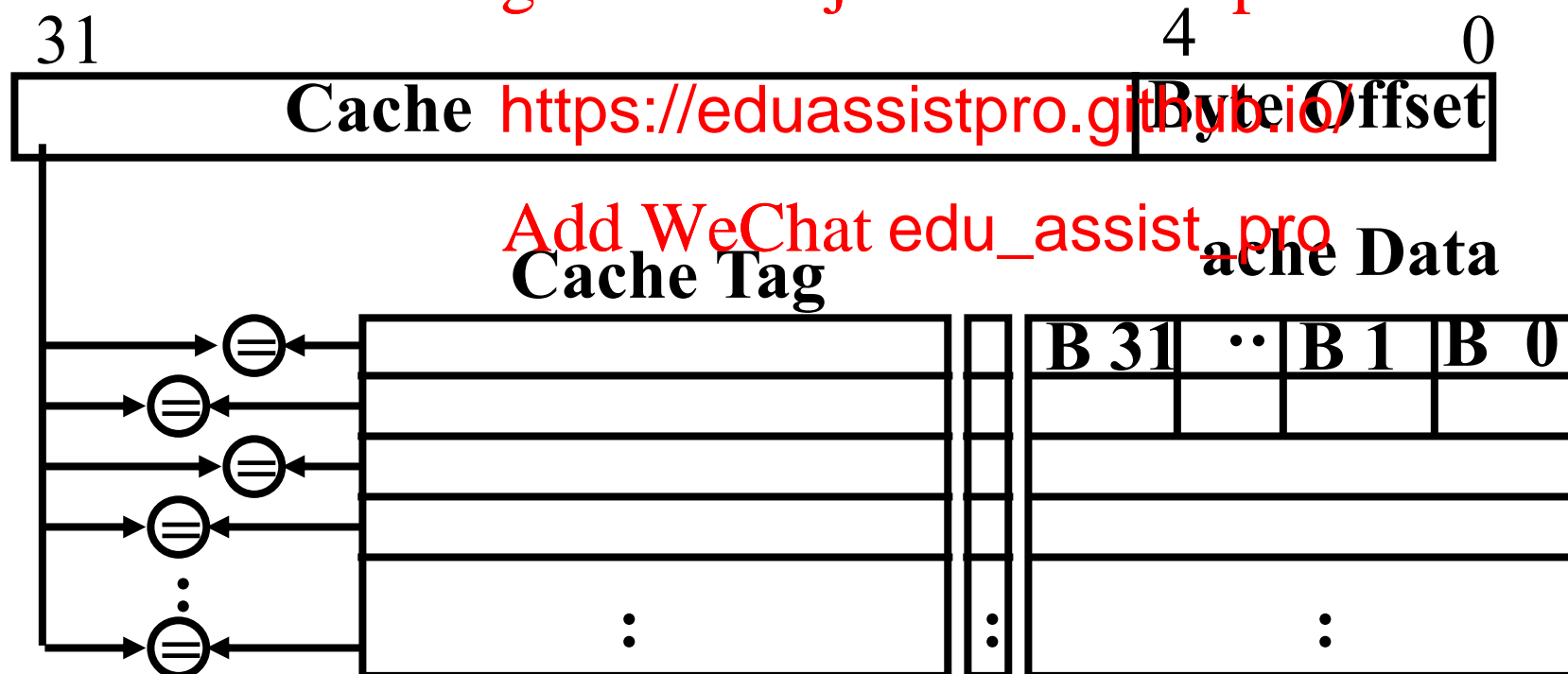
Add WeChat edu_assist_pro

Fully Associative Cache (2/3)

- Fully Associative Cache (e.g., 32 B block)

– compare tags in parallel

Assignment Project Exam Help



Fully Associative Cache (3/3)

- Benefit of Fully Assoc Cache
 - No Conflict Misses (since data can go anywhere)
- Drawbacks of Fully
 - Need hardware comparators for each entry:
 - If we have a 64KB of data in cache with 16B per entry, we need 16K comparators: very expensive
- Small fully associative cache may be feasible

Third Type of Cache Miss

- Capacity Misses

- miss that occurs because the cache has a limited size
- miss that would not have occurred if the size of the cache was larger
- sketchy definition, see <https://eduassistpro.github.io/>

- This is the primary type of miss for private caches.

N-Way Set Associative Cache (1/4)

- Memory address fields:
 - Tag: same as before
 - Offset: same as before
 - Index: points us to the correct set (a set in this case)
- So what's the difference?
 - each set contains multiple blocks
 - once we've found correct set, must compare with all tags in that set to find our data

N-Way Set Associative Cache (2/4)

- Summary:
 - cache is direct-mapped with respect to sets
 - each set is fully asso
 - If we have T blocks to have an T/N direct-mapped cache, where at each ind a fully associative N block cache. Each has its own valid bit and data.

N-Way Set Associative Cache (3/4)

- Given memory address:
 - Find correct set using Index value.
 - Compare Tag with all ~~terminated~~ set.
 - If a match occurs, it's ~~ss.~~
 - Finally, use the offset field as usual to retrieve desired data within the desired block.

N-Way Set Associative Cache (4/4)

- What's so great about this?
 - even a 2-way set associative cache avoids a lot of conflict misses
 - hardware cost isn't too high (no comparators)
- In fact, for a cache
 - it's Direct-Mapped if it's 1-way set associative (1 block per set)
 - it's Fully Associative if it's M-way set associative (M blocks per set)
 - so these two are just special cases of the more general set associative design

Block Replacement Policy (1/2)

- Direct-Mapped Cache: index completely specifies which position a block can go in on a miss
- N-Way Set Assoc (N) a set, but block can occupy any position <https://eduassistpro.github.io/>
- Fully Associative: block can be written to any position (there is no index) [Add WeChat edu_assist_pro](#)
- Question: if we have the choice, where should we write an incoming block?

Block Replacement Policy (2/2)

- Solution!
- If there are any locations with valid bit off (empty), then usually write the new block into the first one.
- If all possible locations have valid bit on, we must use a replacement policy by which line which block gets “cached out” on a miss.

Block Replacement Policy: LRU

- LRU (Least Recently Used)
 - Idea: cache out block which has been accessed (read or write) least recently
 - Pro: temporal locality <https://eduassistpro.github.io/> implies likely future use: in fact, this is a very effective policy
 - Con: with 2-way set assoc, easy to keep track (one LRU bit); with 4-way or greater, requires complicated hardware and much time to keep track of this

Block Replacement Example

- We have a 2-way set associative cache with a four word *total* capacity and one word blocks. We perform the following word accesses (ignore byt

0, 2, 0, 1, 4, 0, 2, <https://eduassistpro.github.io/>

How many hits and how many [Add WeChat edu_assist_pro](#) there for the LRU block replacement policy?

Block Replacement Example: LRU

- Addresses 0, 2, 0, 1, 4, 0, ...

0: miss, bring into set 0 (loc 0)

0 remainder 2 is 0, so set 0

2: miss, bring into set 0 (loc 1)

2 remainder 2 is 0, so set 0

Assignment Project Exam Help

<https://eduassistpro.github.io/>

1: miss, bring into set 1

1 remainder 2 is 1, so set 1

4: miss, bring into set 0 (loc 1, replace 2)

4 remainder 2 is 0, so set 0

0: hit

	loc 0	loc 1
set 0	0	iru
set 1		
set 0	iru 0	2
set 1		
set 0	0	iru 2
set 1		
set 0	0	iru 2
set 1	1	iru
set 0	iru 0	4
set 1	1	iru
set 0	0	iru 4
set 1	1	iru

Ways to reduce miss rate

- Larger cache
 - limited by cost and technology
 - hit time of first level cache < cycle time
- More places in the cache of memory - associativity
 - fully-associative
 - any block any line
 - k-way set associated
 - k places for each block
 - direct map: $k=1$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Big Idea

- How do we choose between options of associativity, block size, replacement policy?
- Design against a performance metric
 - Minimize: $Average\ Access\ Time = Hit\ Time + Miss\ Penalty \times Miss\ Ratio$
 - influenced by technology and problem

Example

- Assume
 - Hit Time = 1 cycle
 - Miss rate = 5%
 - Miss penalty = 20 cycles
- Average memory access time = 2 cycles

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

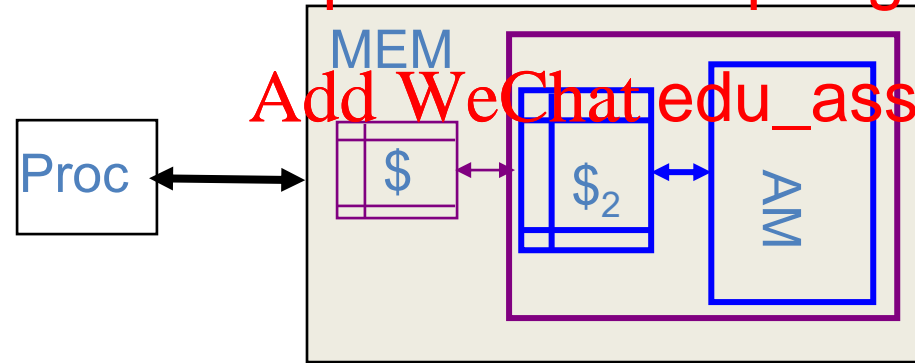
Improving Miss Penalty

- When caches first became popular, Miss Penalty ~ 10 processor clock cycles
- Today: 1000 MHz Processor (1 ns per clock cycle) and 100 ns to go to DRAM
 \Rightarrow 100 process

Assignment Project Exam Help

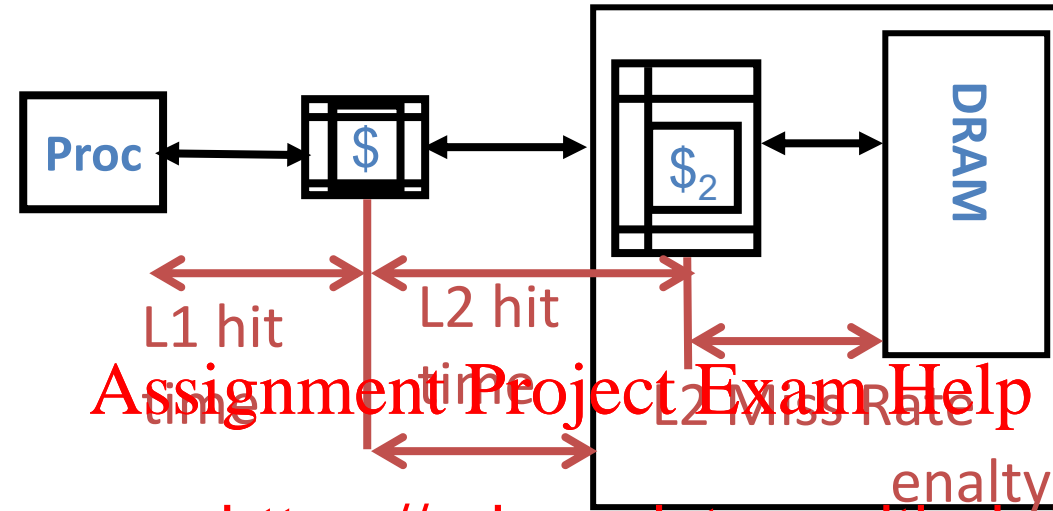
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Solution: another cache between memory and the processor cache: Second Level (L2) Cache

Analyzing Multi-level cache hierarchy



<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Avg Mem Access Time = L1 Hit Time + L1 Miss Rate * L1 Miss Penalty

L1 Miss Penalty = L2 Hit Time + L2 Miss Rate * L2 Miss Penalty

**Avg Mem Access Time = L1 Hit Time +
L1 Miss Rate * (L2 Hit Time + L2 Miss Rate * L2 Miss Penalty)**

Typical Scale

- L1
 - size: tens of KB
 - hit time: complete in one clock cycle
 - miss rates: 1-5%
- L2
 - size: hundreds of KB
 - hit time: few clock cycles
 - miss rates: 10-20%
- L2 miss rate is fraction of L1 misses that also miss in L2
 - why so high?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Example: without L2 cache

- Assume
 - L1 Hit Time = 1 cycle
 - L1 Miss rate = 5%
 - L1 Miss Penalty = 10
- Average memory access time = 6 cycles

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Example with L2 cache

- Assume
 - L1 Hit Time = 1 cycle
 - L1 Miss rate = 5%
 - L2 Hit Time = 5 cycles
 - L2 Miss rate = 15%
 - L2 Miss Penalty = 100 cycles
- L1 miss penalty = $5 + 0.15 * 100 = 20$
- Average memory access time = $1 + 0.05 * 20$
= 2 cycle

3x faster with L2 cache

What to do on a write hit?

- Write-through

- update the word in cache block and corresponding word in memory

- Write-back

- update word in cache <https://eduassistpro.github.io/>
 - allow memory word to be “stale”
 - *add ‘dirty’ bit to each line indicating memory needs to be updated when block is replaced*
 - *OS flushes cache before I/O !!!*

- Performance trade-offs?

“And in conclusion...” (1/2)

- Caches are NOT mandatory:
 - Processor performs arithmetic
 - Memory stores data
 - Caches simply make d
- Each level of memory hierarchy is a set of next higher level
- Caches speed up due to **temporal locality**: store data used recently
- Block size > 1 word speeds up due to **spatial locality**: store words adjacent to the ones used recently

“And in conclusion...” (2/2)

- Cache design choices:
 - size of cache: speed v. capacity
 - direct-mapped v. associative
 - for N-way set assoc: <https://eduassistpro.github.io/>
 - block replacement policy
 - 2nd level cache?
 - Write through v. write back?
- Use [performance model](#) to pick between choices, depending on programs, technology, budget, ...

A real example

- And additional reading (for fun):

<http://igoro.com/archive/gallery-of-processor-cache-effects/>

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

average time loop 1 = 126858657.625

average time loop 2 = 71715726.125

ratio is 1.7689098957721807

but first loop does 32 times more work!!

Loop 1 gets work done 18 times faster!

```

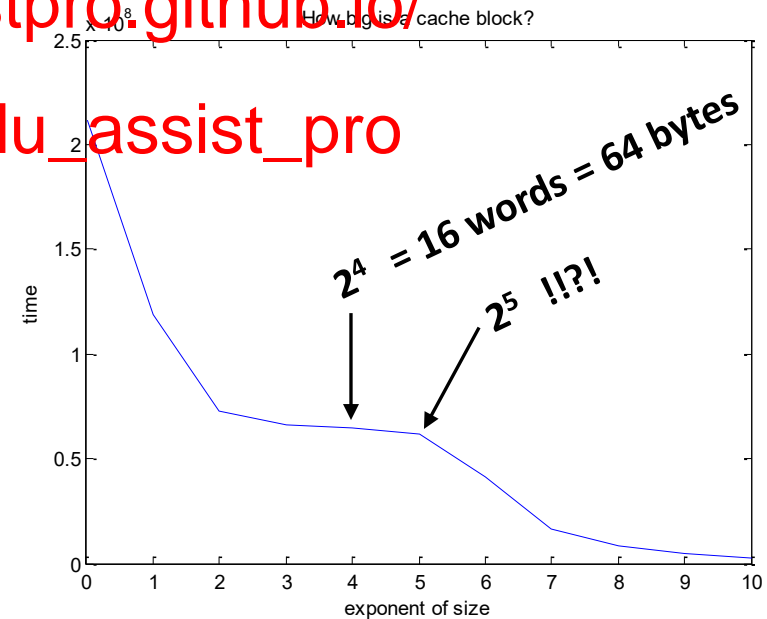
A=[
0 1 240591499
1 2 134307003
2 4 84736089
3 8 74437939
4 16 70215291
5 32 73695400
6 64 52077957
7 128 19758427
8 256 10488407
9 512 6369311
10 1024 3736937
];
plot(A(:,1),A(:,3));
title('How big is a cache block?');
ylabel('time');
xlabel('exponent of size');

```

Assignment Project Exam Help

<https://eduassistpro.github.io/>

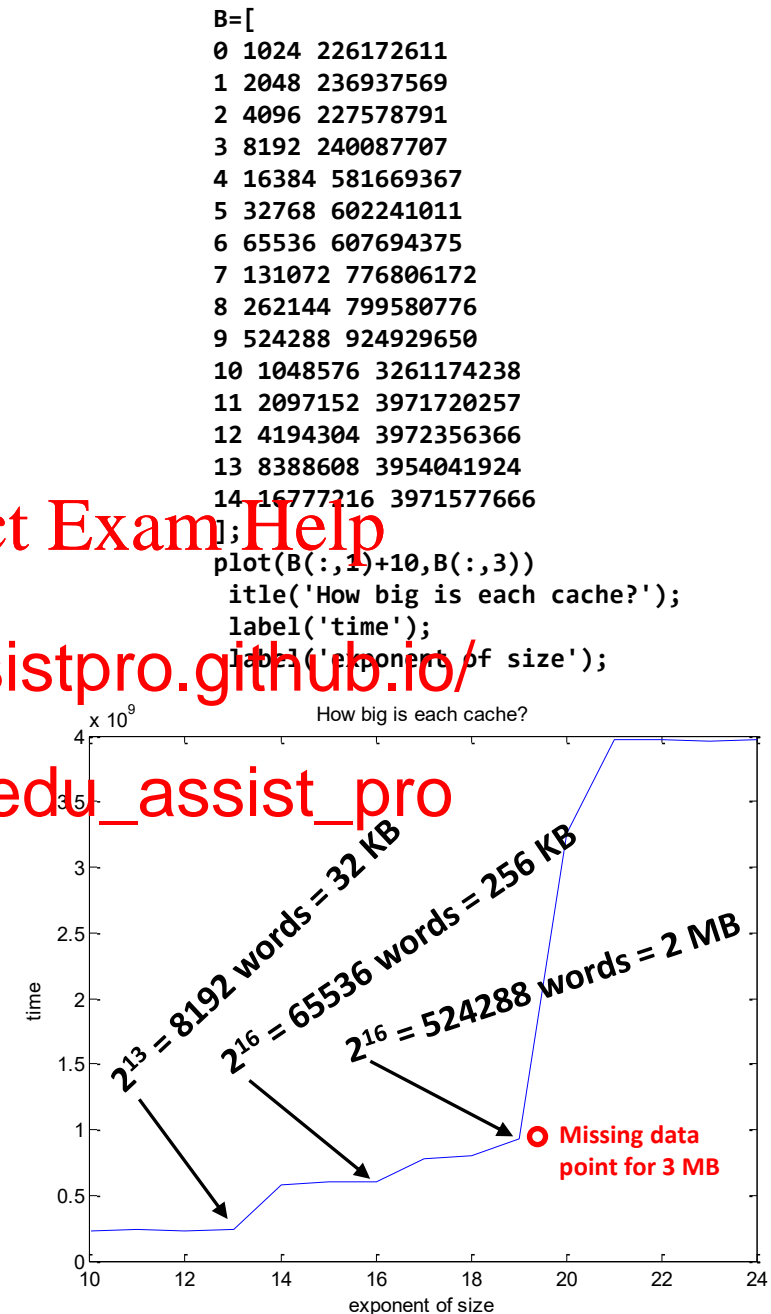
Add WeChat edu_assist_pro



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Review and More Information

- Sections 5.3 - 5.4 of textbook

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro