[§5.2.6] Translation Lookaside Buffers

The CPU generates *virtual* addresses, which correspond to locations in virtual memory.

In principle, the virtual addresses are translated to physical addresses using a page table.



But this is too slow, so in practice, a *translation lookaside buffer* (TLB) is used.

It is like a special cache that is indexed by page number.

If there is a hit on a page number, then the address of the page in memory (called the *page-frame address*) is immediately obtained.

Therefore, the TLB and the cache must be accessed sequentially.



This adds an extra cycle in case of a hit.

(The page *displacement* is sometimes called the "page offset." But we will call it the displacement to avoid confusion with the block offset," which we just call "offset.")

How can we avoid wasting this time?



Let's look at what happens when a memory address is accessed.

What are the steps in cache access?

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.

We always need to read lines into the sense amplifiers and then select the word (cf. the direct-mapped cache diagram in Lecture 4).

Now, if we know the index *before* address translation takes place, <u>we</u> <u>can perform steps</u> _____ while address translation is occurring.

There is a tradeoff between speed and power efficiency.

- For power efficiency, which order should should steps 1 through 4 be performed in?
- For maximum speed, which of steps 1 through 4 can be performed in parallel?

Let's take a look at address translation.



In this example, what is the page size (in bytes)?

How much physical memory is there?

Our goal is to allow the cache to be indexed before address translation completes.

In order to do that, we need to have the index field be *entirely contained* within the page displacement.

So, if the displacement is *d* bits wide, the width of the index is *j* bits, and the offset is *k* bits, we must have $j + k \le d$.



Cache hit time reduces from two cycles to one!

... because the cache can now be *indexed* in parallel with TLB (although the tag match uses output from the TLB).

But there are some constraints...

• Suppose our cache is direct mapped. Then the index field just contains the line number. So, (line number || block offset) must fit inside the page displacement.

What is the largest the cache can be?

If we want to increase the size of the cache, what can we do?

Options:

• For new machines, select page size such that-

page size $\geq \frac{\text{cache size}}{\text{associativity}}$

• If page size is fixed, select associativity so that-

associativity $\geq \frac{\text{cache size}}{\text{page size}}$

Example: MC88110

- Page size = 4KB
- I-cache, D-cache are both: 8KB, 2-way set-associative (4KB = 8KB / 2)

Example: VAX series

- Page size = 512B
- For a 16KB cache, need assoc. = (16KB / 512B) = 32-way set. assoc.!

The textbook gives these three alternatives for cache indexing and tagging. <u>Answer some questions</u> about them.





Virtually Indexed and Tagged



What's the main disadantage of physically indexed and tagged?

What is the organization we have just been discussing (in the last diagram)?

What is the main disadvantage of virtually indexed and tagged?

Virtually Indexed but Physically Tagged



Multilevel cache design

What are distinguishing <u>features of the different cache levels</u> of the four-level design (from 2013) illustrated on p. 135 of the textbook?

	Distinguish- ing feature	Size	Access time	Implement'n techology
L1 cache				
L2 cache				
L3 cache				
L4 cache				
Main mem.				

What are some advantages of a centralized cache?

What are some advantages of a banked structure?

Inclusion in multilevel caches

Answer these questions about inclusion policies.

Which kind(s) of caches move a block from one level to the other?

Which kind(s) of caches propagate up an eviction from the L2 to the L1?

Which kind(s) of caches have to inform the L2 about a write to the L1?

In an inclusive cache, can L2 associativity be greater than L1 associativity?

Find and describe the typo in this diagram.



Replacement policies

LRU is a good strategy for cache replacement.

In a set-associative cache, LRU is reasonably cheap to implement. Why?

With the LRU algorithm, the lines can be arranged in an *LRU stack*, in order of recency of reference. Suppose a string of references is—

and there are 4 lines. Then the LRU stacks after each reference are—

а	b	С	d	а	b	е	а	b	С	d	е
	а	b	С	d	а	b	е	а	b	С	d
		а	b	С	d	а	b	е	а	b	С
			а	b	С	d	d	d	е	а	b
*	*	*	*			*			*	*	*

Notice that at each step:

- The line that is referenced moves to the top of the LRU stack.
- All lines below that line keep their same position.
- All lines above that line move down by one position.

How many bits per set are required to keep track of LRU status in both of the implementations described in the text?

- Matrix
- Pseudo-LRU







Figure 5.7: Illustration of pseudo-LRU replacement on a 4-way set associative cache.