

Virtual memory & memory hierarchy

Hung-Wei Tseng

Outline

- Virtual memory
- Architectural support for virtual memory
- Advanced hardware support for virtual memory

Let's dig into this code

```
int main(int argc, char *argv[])
{
    int i,j;
    double **a;
    double sum=0, average;
    int dim=32768;
    if(argc < 2)
    {
        fprintf(stderr, "Usage: %s dimension\n",argv[0]);
        exit(1);
    }
    dim = atoi(argv[1]);
    a = (double **)malloc(sizeof(double *)*dim);
    for(i = 0 ; i < dim; i++)
        a[i] = (double *)malloc(sizeof(double)*dim);
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            a[i][j] = rand();
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            sum+=a[i][j];
    average = sum/(dim*dim);
    fprintf(stderr,"average: %lf\n",average);
    for(i = 0 ; i < dim; i++)
        free(a[i]);
    free(a);
    return 0;
}
```

Let's dig into this code

```
#define _GNU_SOURCE
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>

double a;

int main(int argc, char *argv[])
{
    int i, number_of_total_processes=4;
    number_of_total_processes = atoi(argv[1]);
    // Create processes
    for(i = 0; i< number_of_total_processes-1 && fork(); i++);
    // Generate rand seed
    srand((int)time(NULL)+(int)getpid());
    a = rand();
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n",getpid(), a, &a);
    sleep(10);
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n",getpid(), a, &a);
    return 0;
}
```

Virtual Memory



Program A

0f00bb27	00c2e800
509cbd23	00000008
00005d24	00c2f000
0000bd24	00000008
2ca422a0	00c2f800
130020e4	00000008
00003d24	00c30000
2ca4e2b3	00000008



Program B

0f00bb27	00c2e800
509cbd23	00000008
00005d24	00c2f000
0000bd24	00000008
2ca422a0	00c2f800
130020e4	00000008
00003d24	00c30000
2ca4e2b3	00000008

Processor

Virtual memory

PC for A

PC for A

load

0x0
0f00bb27
509cbd23
00005d24
0000bd24

PC for B

PC for B

load

load

page

0x0

0f00bb27
509cbd23
00005d24
0000bd24
2ca422a0
130020e4
00003d24
2ca4e2b3

00c2e800
00000008
00c2f000
00000008
00c2f800
00000008
00c30000
00000008

264-1

264-1

This approach is called demand paging + swapping

00c2f800	00c2e800
00000008	00000008
00c30000	00c2f000
00000008	00000008
2ca422a0	00c2e800
130020e4	00000008
00003d24	00c2f000
2ca4e2b3	00000008
0f00bb27	2ca422a0
509cbd23	130020e4
00005d24	00003d24
0000bd24	2ca4e2b3

Physical Memory

0f00bb27
509cbd23
00005d24
0000bd24

Swap

Demo revisited

&a = 0x601090

Process A

Process B

**Process A's
Virtual
Memory Space**

**Process B's
Virtual
Memory Space**

```
#define _GNU_SOURCE
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>
```

```
double a;
```

```
int main(int argc, char *argv[])
```

```
{
```

```
    int i, number_of_total_processes=4;
```

```
    number_of_total_processes = atoi(argv[1]);
```

```
    for(i = 0; i< number_of_total_processes-1 && fork(); i++);
```

```
    srand((int)time(NULL)+(int)getpid());
```

```
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n",getpid(), a, &a);
```

```
    sleep(10);
```

```
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n",getpid(), a, &a);
```

```
    return 0;
```

```
}
```

Virtual memory

- An **abstraction** of memory space available for programs/software/programmer
- Programs execute using virtual memory address
- The operating system and hardware work together to handle the mapping between virtual memory addresses and real/physical memory addresses
- Virtual memory organizes memory locations into "**pages**"

Why Virtual memory?

- Allowing multiple applications to share physical main memory
 - Memory protection/isolation among programs/processes is automatically achieved
- Allowing applications to work even the installed physical memory or available physical memory is smaller than the working set of the application
 - Programmer does not need to worry about the physical memory capacity of different machines — make compiled program compatible
 - Multiple programs can work concurrently even though their total memory demand is larger than the installed physical memory

Processor
Core

Registers

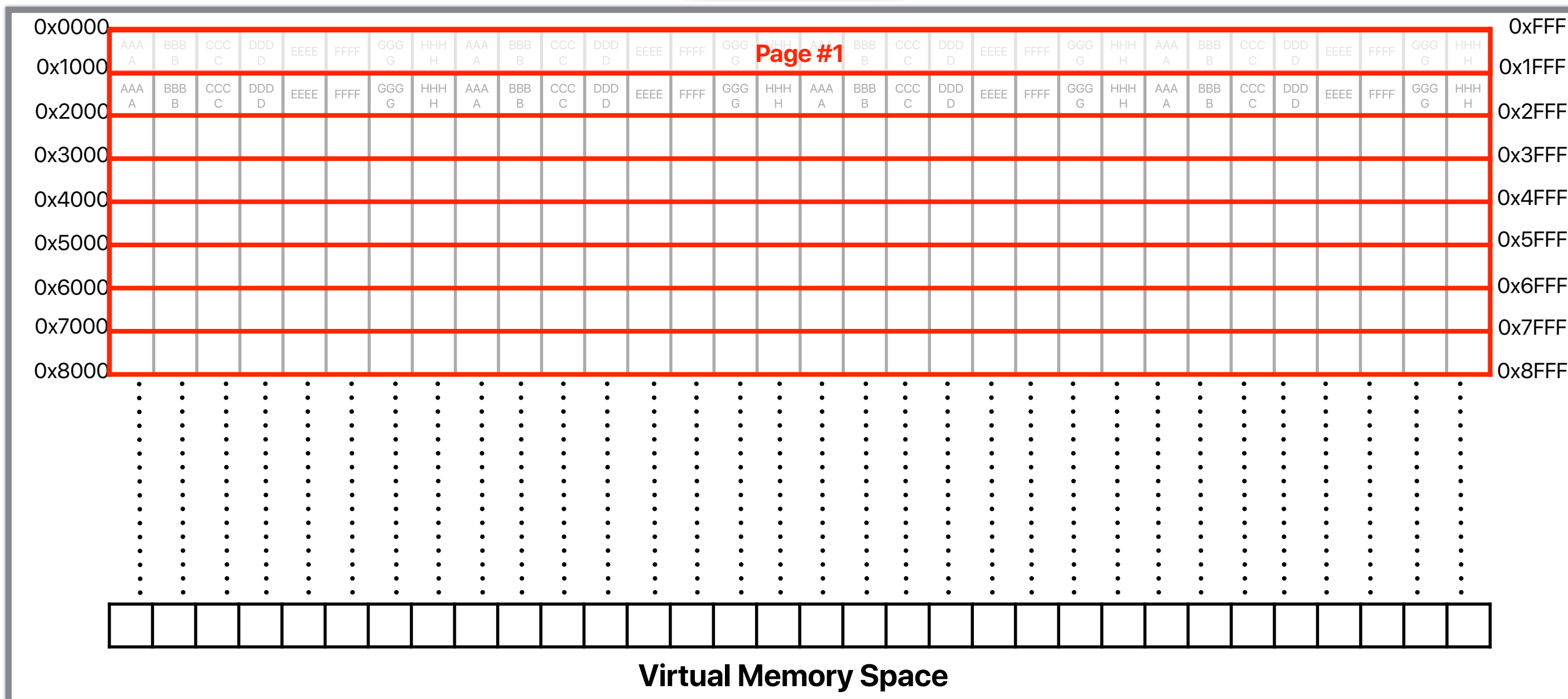
The virtual memory abstraction

load 0x0009

Page table

Main memory
(DRAM)

Page #1

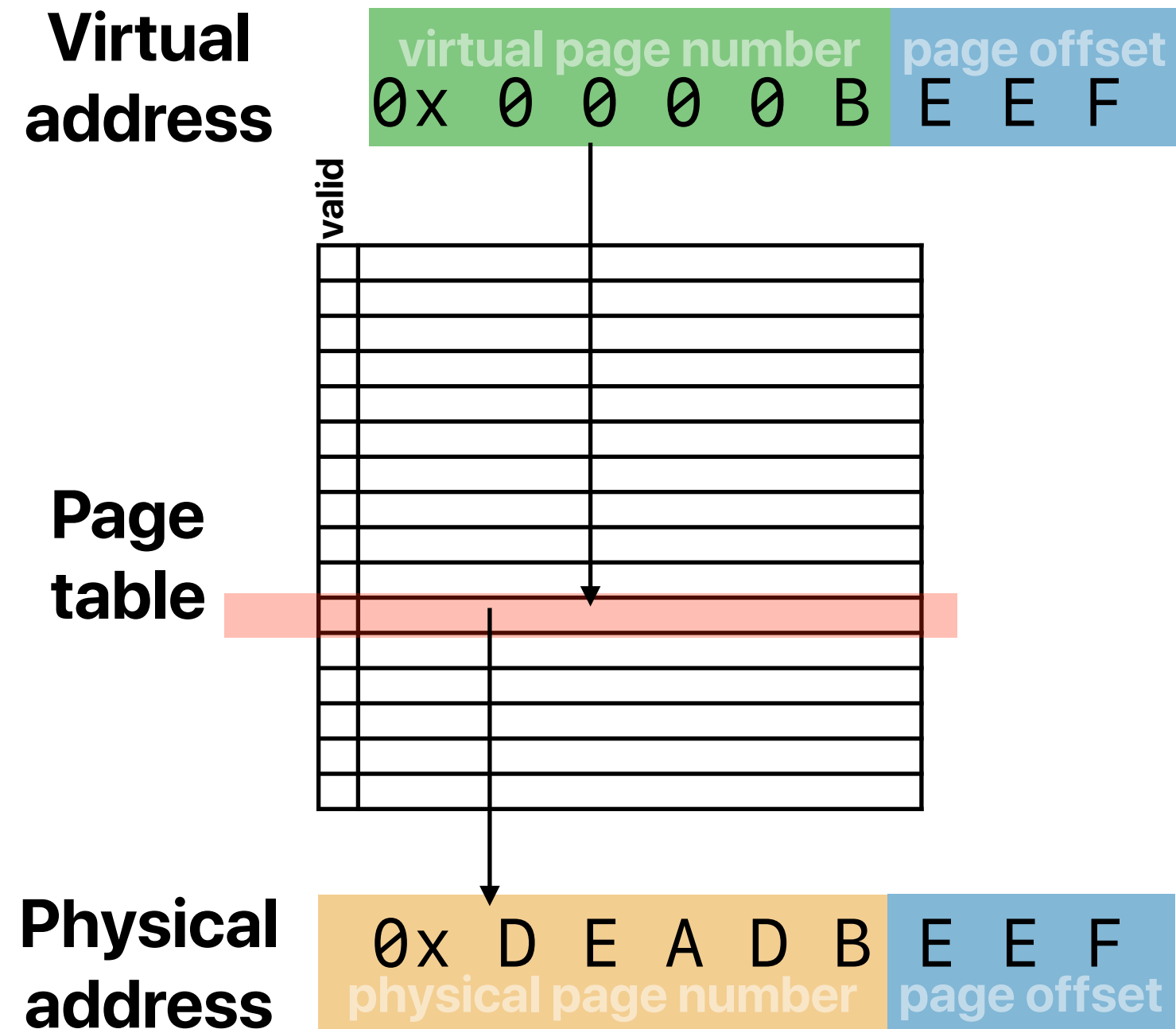


Demand paging

- Treating physical main memory as a “cache” of virtual memory
- The block size is the “page size”
- The page table is the “tag array”
- It’s a “fully-associate” cache — a virtual page can go anywhere in the physical main memory

Address translation

- Processor receives virtual addresses from the running code, main memory uses physical memory addresses
- Virtual address space is organized into "pages"
- The system references the **page table** to translate addresses
 - Each process has its own page table
 - The page table content is maintained by OS



Conventional page table

0x0

0xFFFFFFFFFFFFFFFF

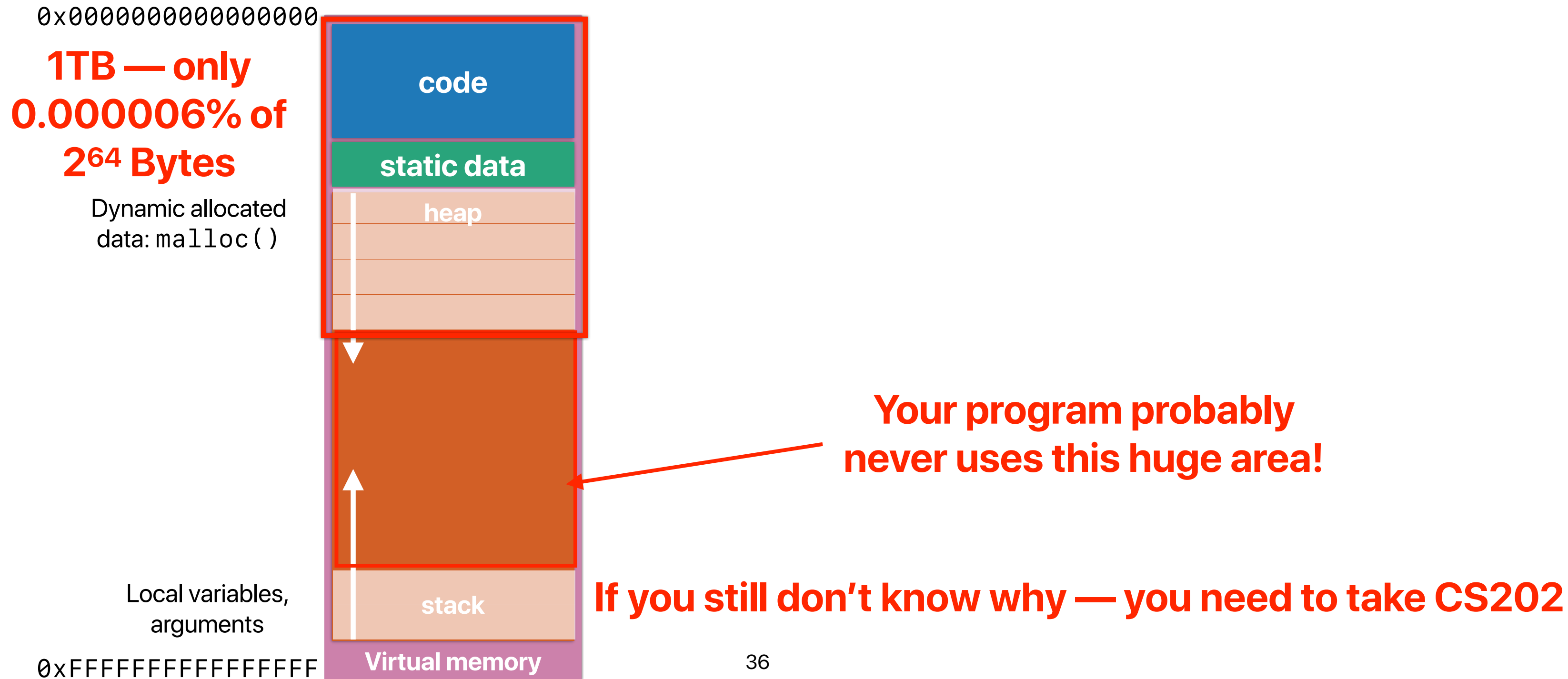
Virtual Address Space

- must be consecutive in the physical memory
- need a big segment! — difficult to find a spot
- simply too big to fit in memory if address space is large!

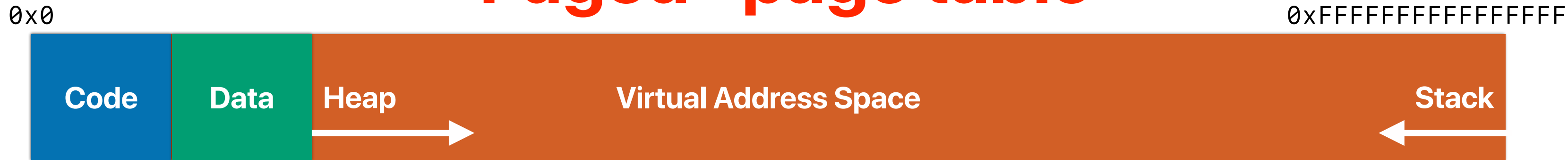
$\frac{2^{64} B}{2^{12} B}$ page table entries/leaf nodes



Do we really need a large table?



"Paged" page table

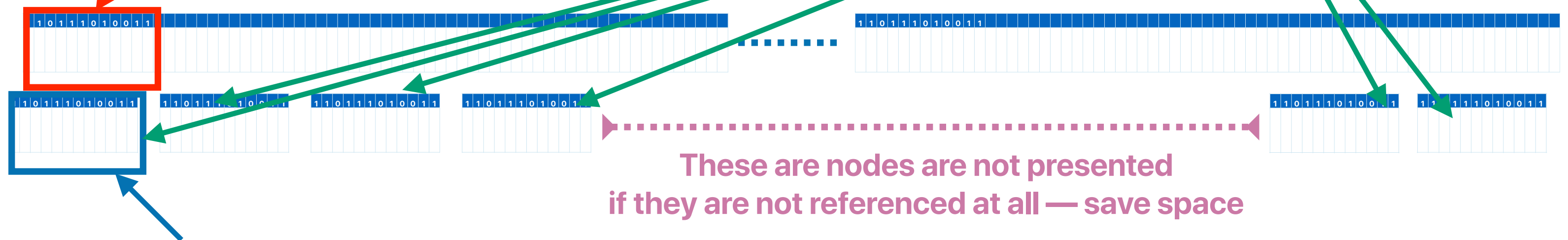


Break up entries into pages!
Each of these occupies exactly a page

$$\frac{2^{12} B}{2^3 B} = 2^9 \text{ PTEs per node}$$

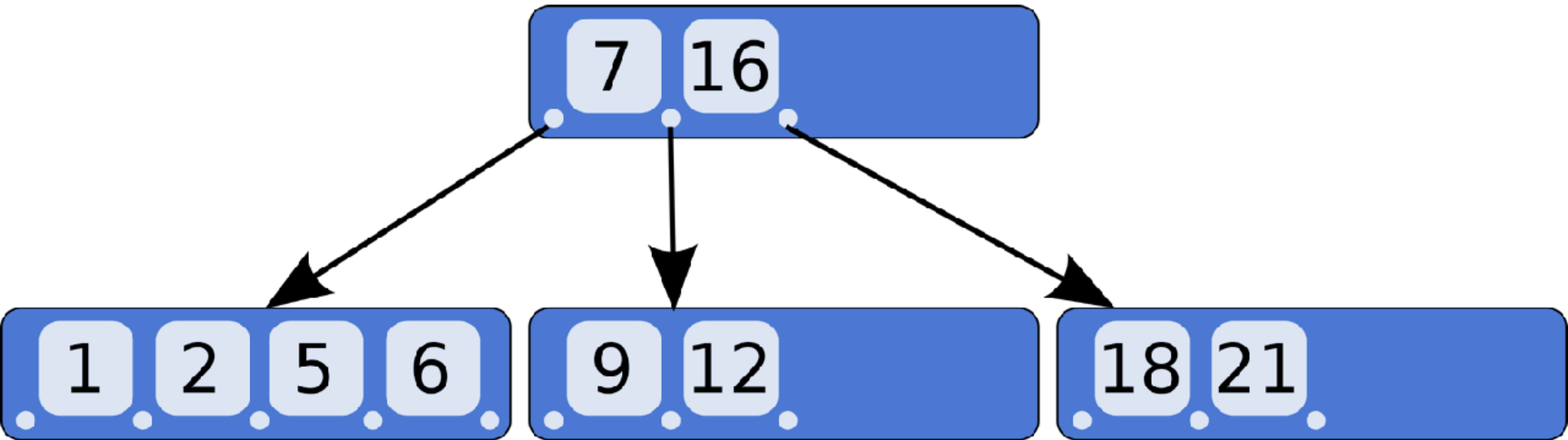
Otherwise, you always need to find more than one consecutive pages — difficult!

Question:
These nodes are spread out,
how to locate them in the memory?



Allocate page table entry nodes "on demand"

B-tree

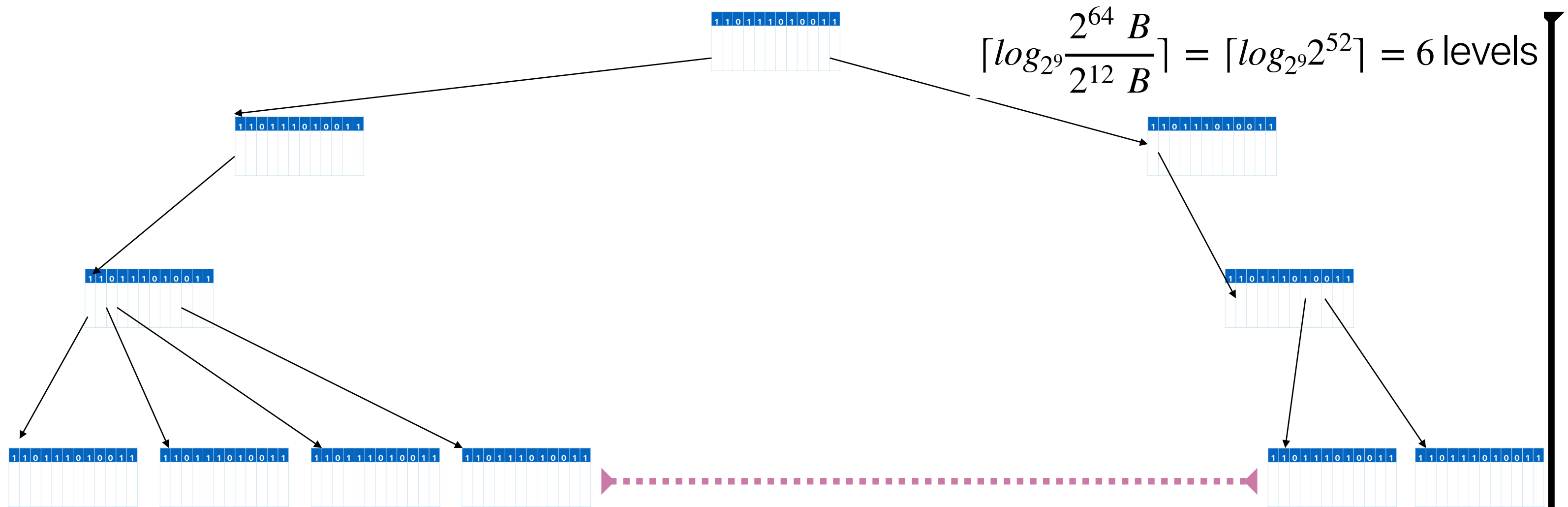
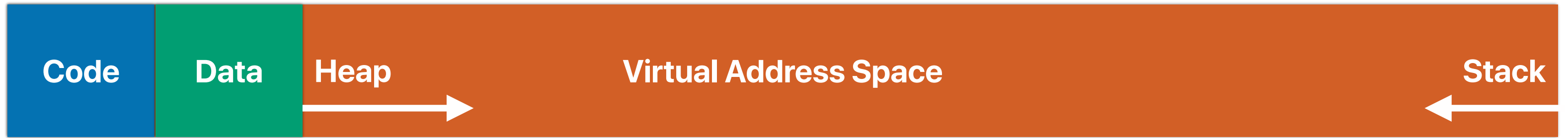


<https://en.wikipedia.org/wiki/B-tree#/media/File:B-tree.svg>

Hierarchical Page Table

0x0

0xFFFFFFFFFFFFFFFF

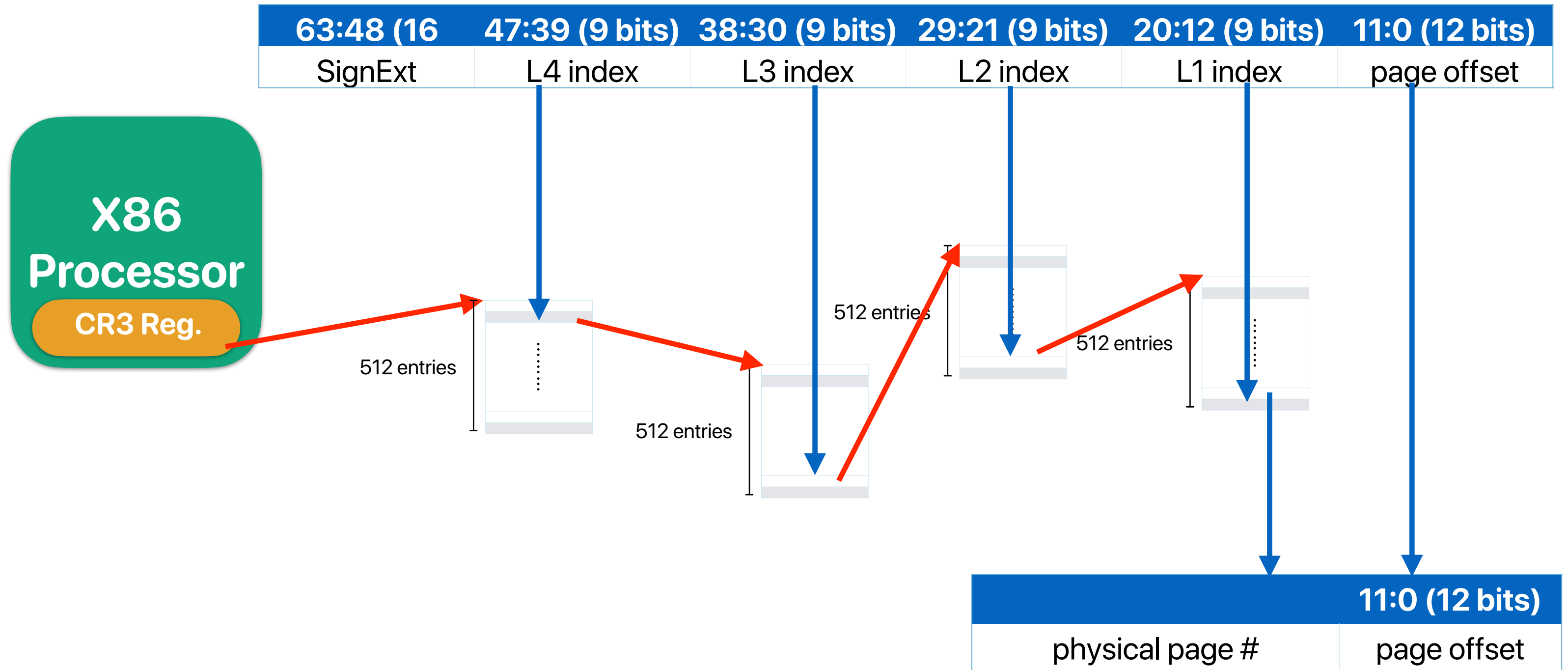


$$\lceil \log_2 \frac{2^{64} B}{2^{12} B} \rceil = \lceil \log_2 2^{52} \rceil = 6 \text{ levels}$$

$\frac{2^{64} B}{2^{12} B}$ page table entries/leaf nodes (worst case)

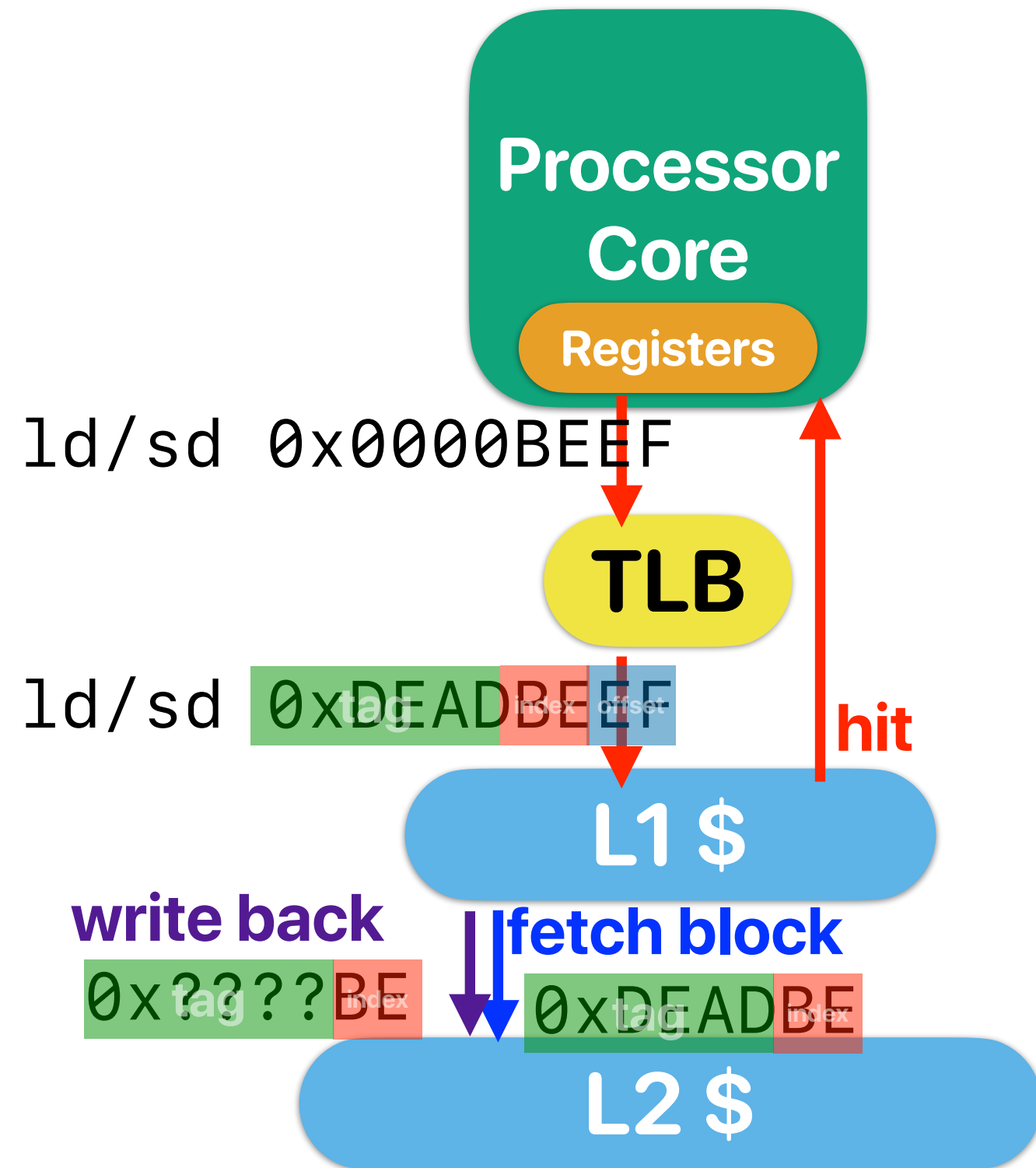
These are nodes are not presented as they are not referenced at all.

Address translation in x86-64



Avoiding the address translation overhead

TLB: Translation Look-aside Buffer



- TLB — a small SRAM stores frequently used page table entries
- Good — A lot faster than having everything going to the DRAM
- Bad — Still on the critical path

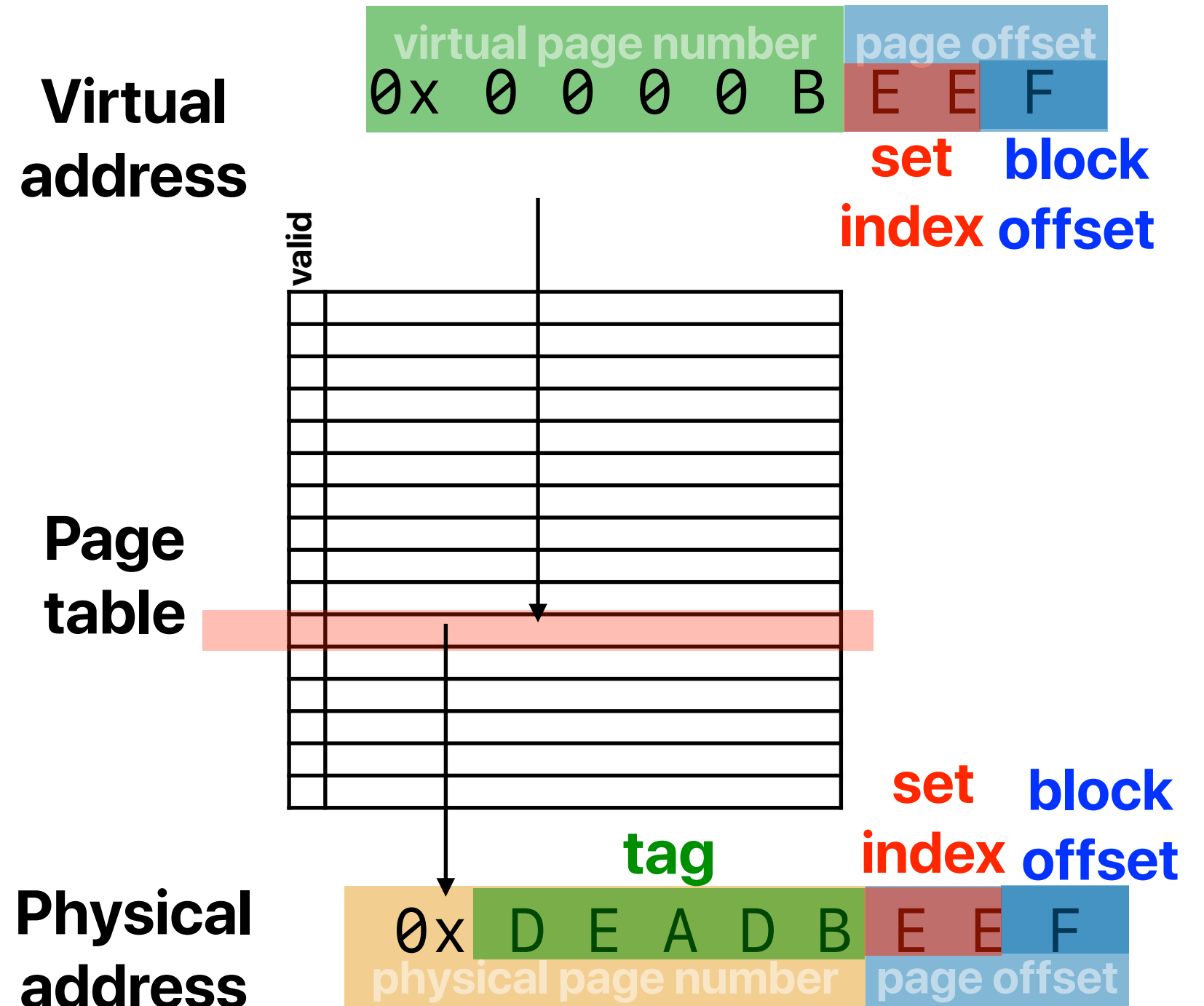
TLB + Virtual cache

- L1 \$ accepts virtual address — you don't need to translate
- Good — you can access both TLB and L1-\$ at the same time and physical address is only needed if L1-\$ misses
- Bad — it doesn't work in practice
 - Many applications have the same virtual address but should be pointing different **physical addresses**
 - An application can have "aliasing virtual addresses" pointing to the same **physical address**



Virtually indexed, physically tagged cache

- Can we find physical address directly in the virtual address — Not everything — but the page offset isn't changing!
- Can we indexing the cache using the "partial physical address"?
 - Yes — Just make set index + block set to be exactly the page offset



Virtually indexed, physically tagged cache

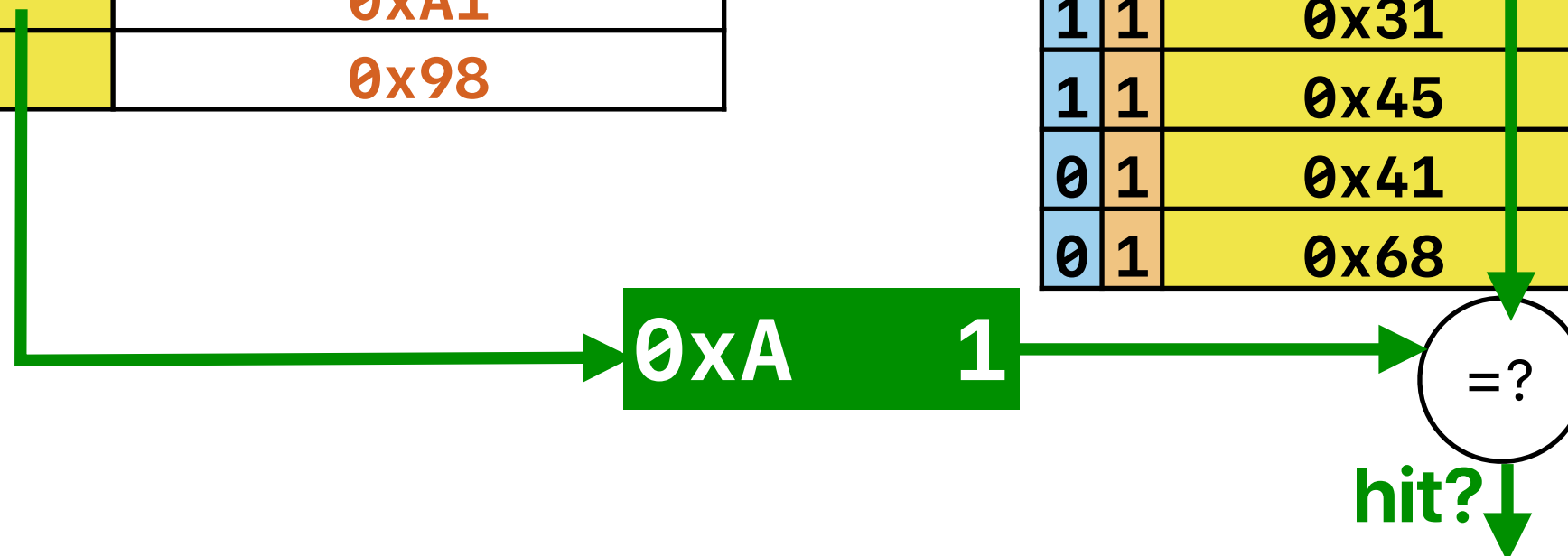
memory address:

0x0	8	2	4
		set	block

memory address: 0b0000100000100100

V	virtual page #	physical page #
1	0x29	0x45
1	0xDE	0x68
1	0x10	0xA1
0	0x8A	0x98

V D		tag	data
1	1	0x00	AABBCCDDEEGGFFHH
1	1	0x10	IIJJKKLLMMNNOOPP
1	0	0xA1	QQRRSSTTUUVVWWXX
0	1	0x10	YYZZAABBCCDDEEFF
1	1	0x31	AABBCCDDEEGGFFHH
1	1	0x45	IIJJKKLLMMNNOOPP
0	1	0x41	QQRRSSTTUUVVWWXX
0	1	0x68	YYZZAABBCCDDEEFF



Virtually indexed, physically tagged cache

- If page size is 4KB —

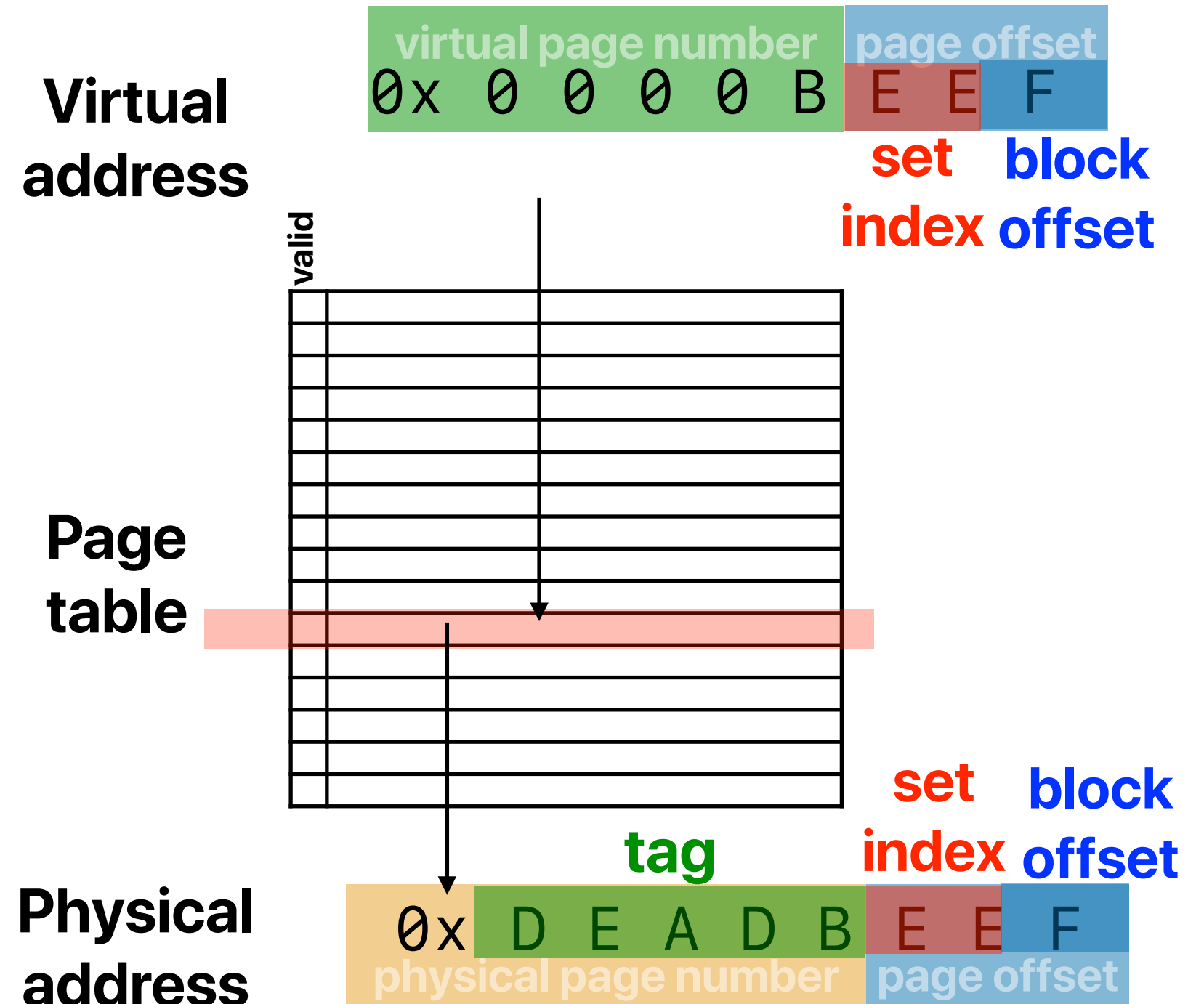
$$lg(B) + lg(S) = lg(4096) = 12$$

$$C = ABS$$

$$C = A \times 2^{12}$$

if $A = 1$

$$C = 4KB$$



Translation Caching: Skip, Don't Walk (the Page Table)

Thomas W. Barr, Alan L. Cox, Scott Rixner

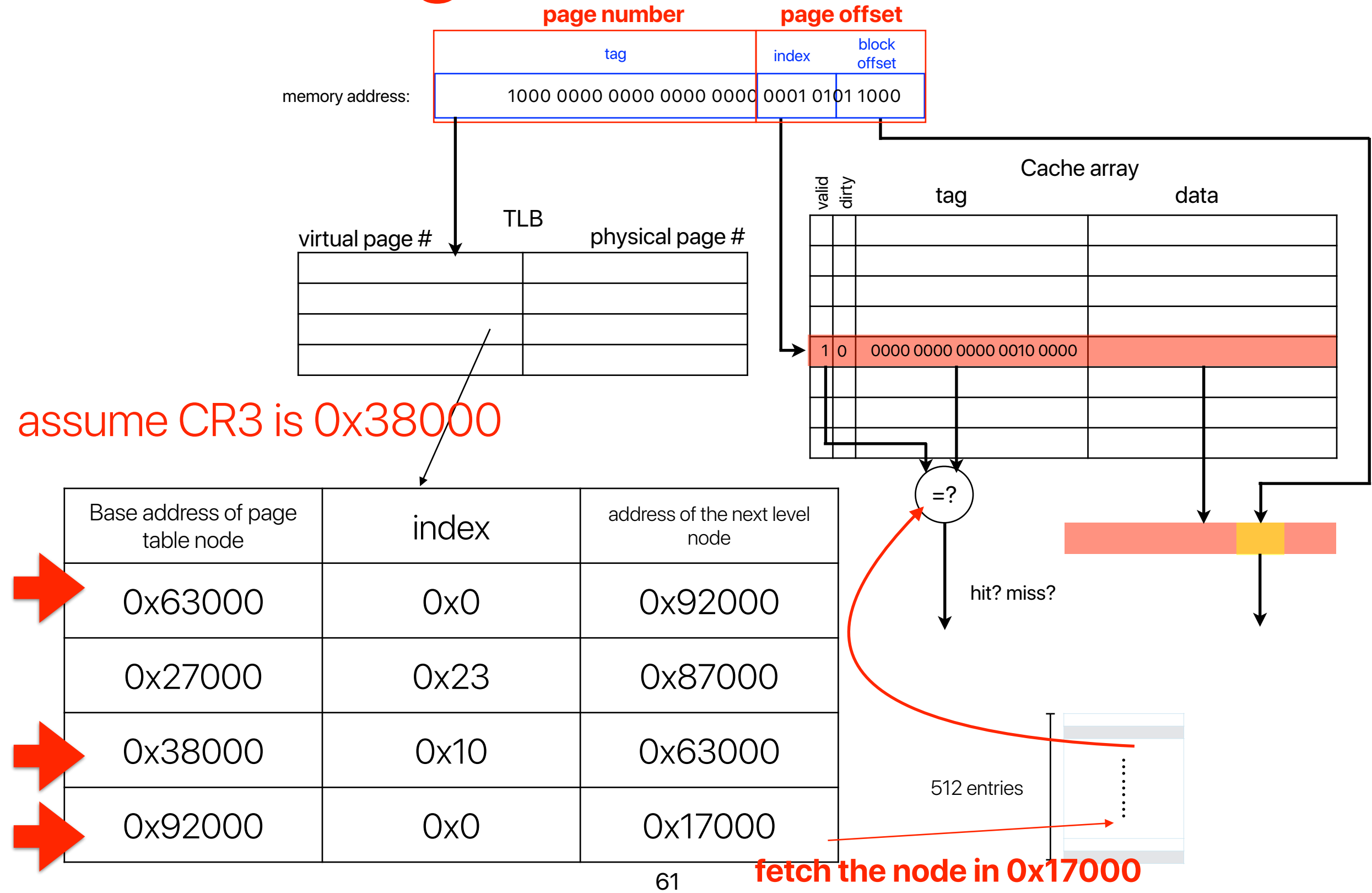
Why should we care about this paper?

- TLB miss is expensive
 - You have to walk through multiple nodes in the hierarchical page table
 - Each node is a memory access — 100 ns
- Modern processors use memory management units (MMUs)
 - MMUs have caches, but not optimized for the timing critical TLB miss
 - Page table caches
 - Translational caches

What this paper proposed

- Not really proposing anything. More an empirical analysis paper
- Design space exploration to find out the optimal solution

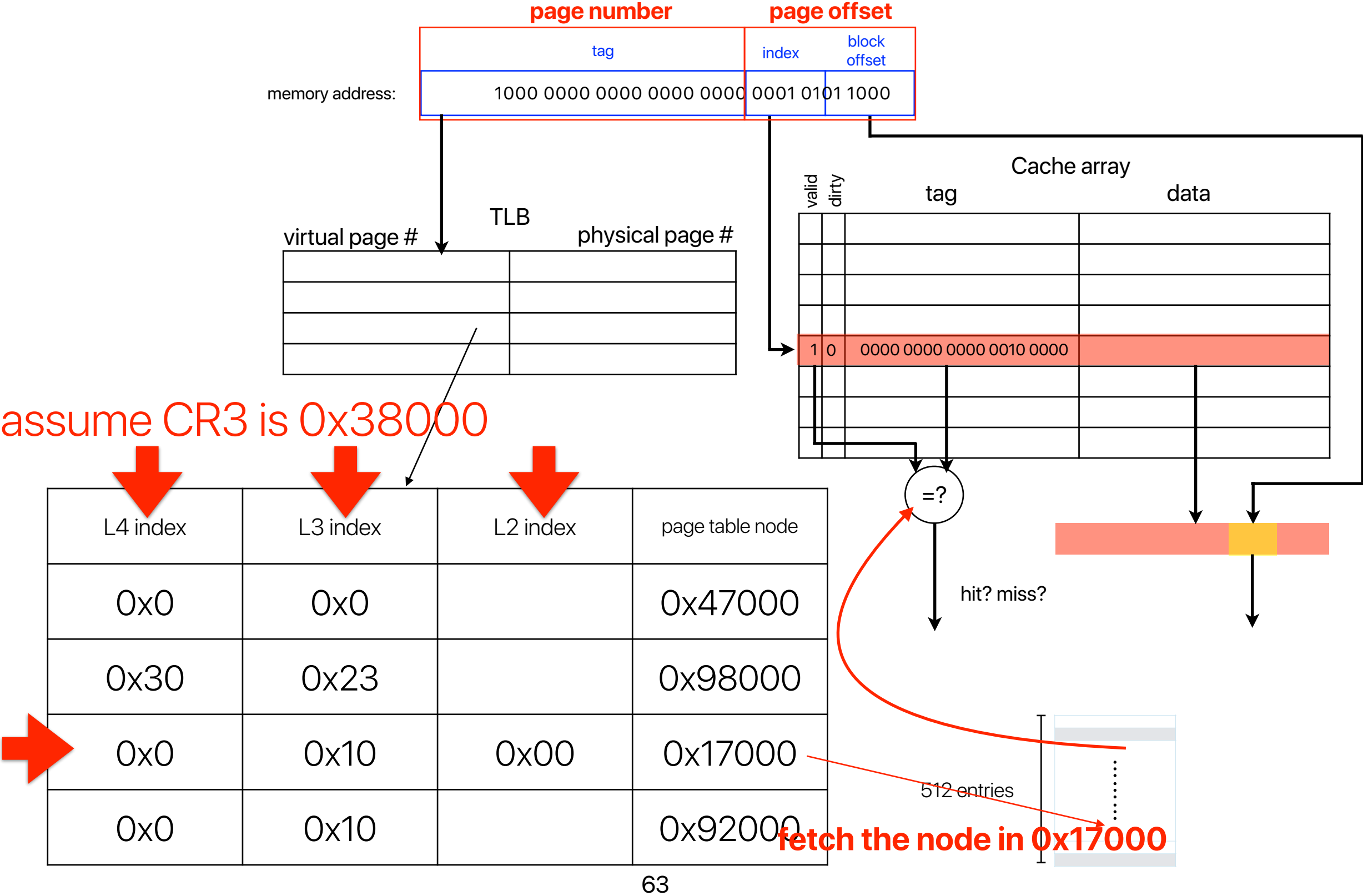
Page table caches



Page table caches

- PTC caches the addresses of “page table nodes”
- PTC uses the physical address of page table nodes as the index
 - Unified page table cache (UPTC)
 - Split page table cache (SPTC)
 - Each page level get a private cache location

Translation cache



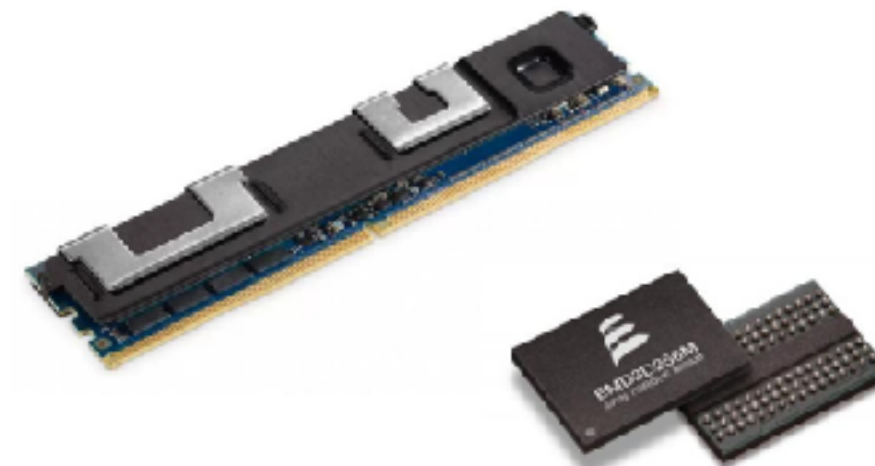
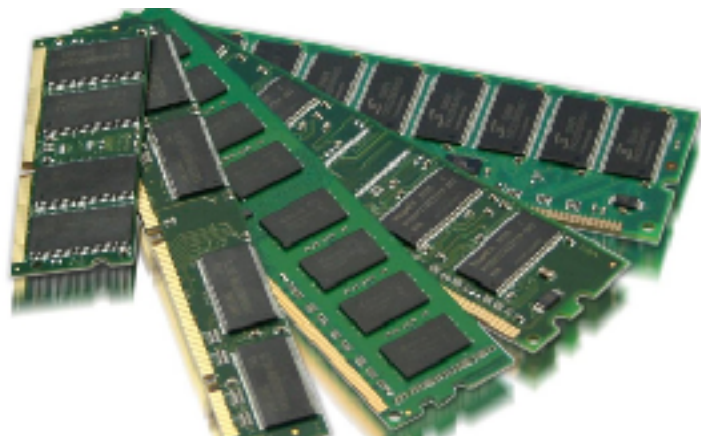
Translation caches

- Indexed by the prefix of the requesting virtual address
 - Split translational cache (STC)
 - Unified translational cache (UTC)
 - Translational-path Cache (TPC)
- Pros:
 - Allowing each level lookup to perform independently, in parallel
- Cons:
 - Less space efficient

Efficient Virtual Memory for Big Memory Servers

**Arkaprava Basu, Jayneel Gandhi, Jichuan Chang, Mark D. Hill and
Michael M. Swift**

Applications with big memory footprints and high-density memory technologies



Why should we care about this paper?

- We care about big memory applications (e.g. memcached)
- Machines with TBs of physical memory are available
- Big-memory workloads pay high costs in paging
 - Up to 51% execution time burned due to TLB miss with 4KB page size
 - 10% of execution time in TLB miss even with 2MB page sizes

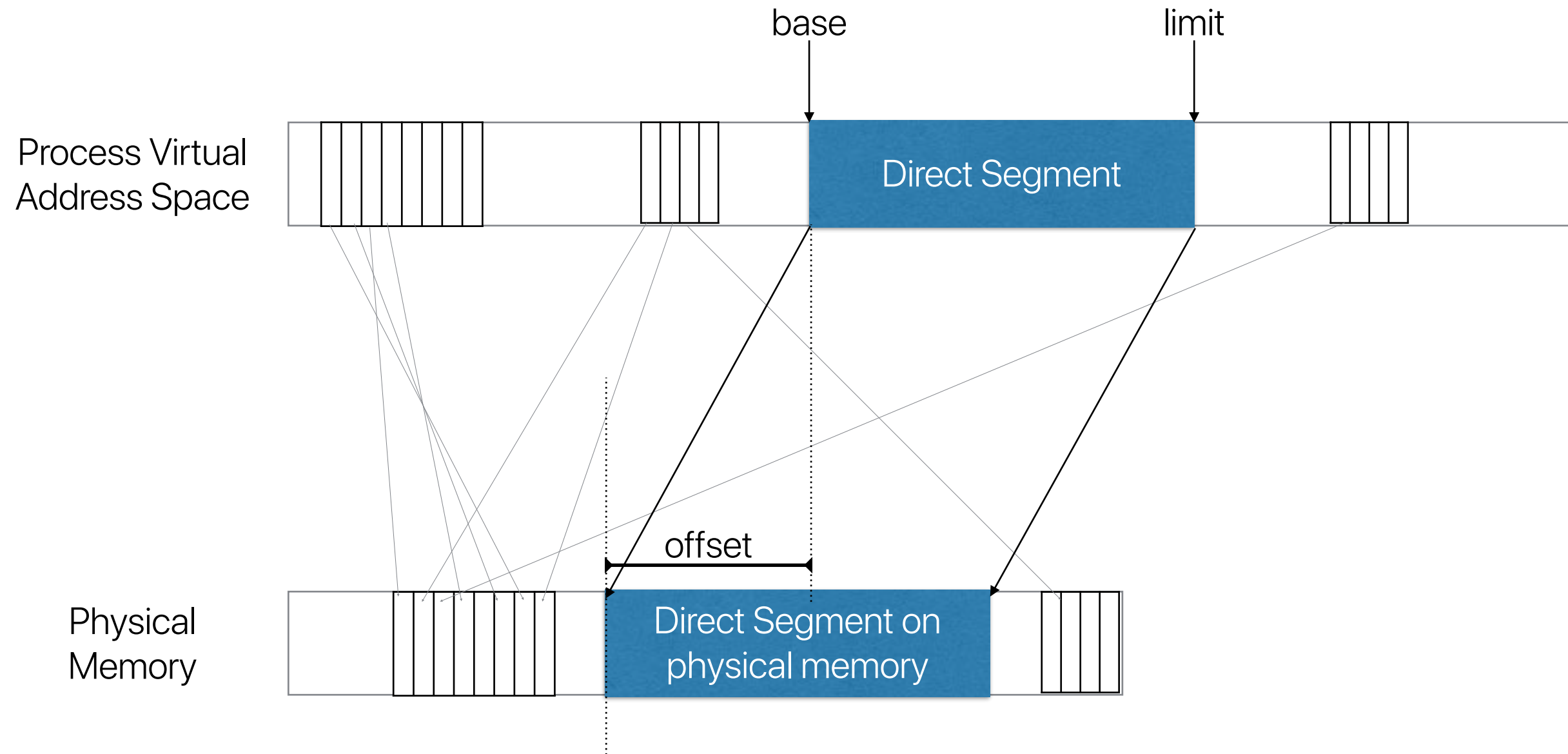
Characteristics of Big Memory Apps

- Few swapping — since the computer has big memory
- Few copy-on-write — since the workload is mostly reads
- Does not require per-page protection
— remind yourself how you usually allocate memory?

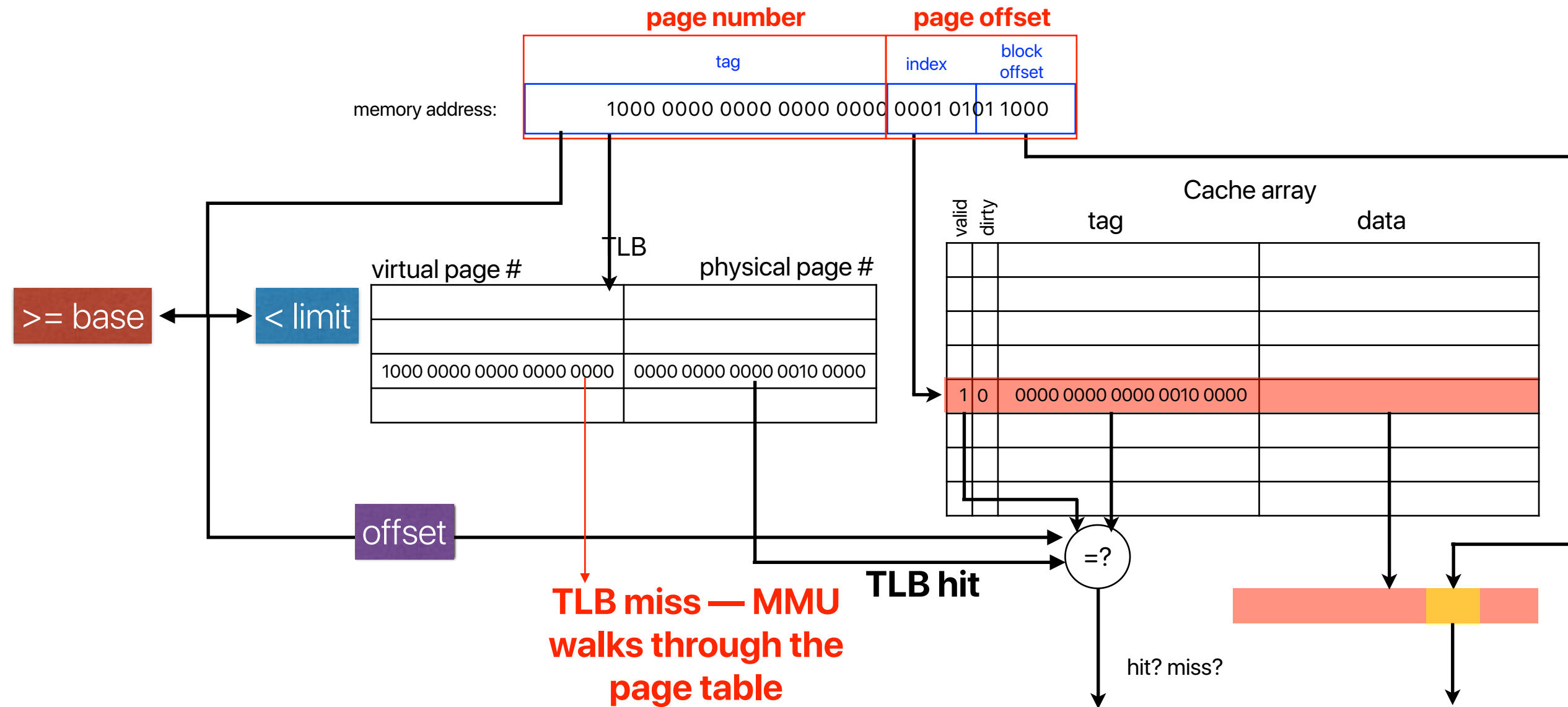
What this paper proposed?

- Mapping part of a process's linear virtual address to a "direct segment" rather than a page
- Direct segment
 - Similar to classic segmentation: adding base, offset, limit registers to each core
 - If the virtual address falls in the range between base and limit, no TLB access is necessary
- Virtual memory outside a direct segment still uses conventional demand paging

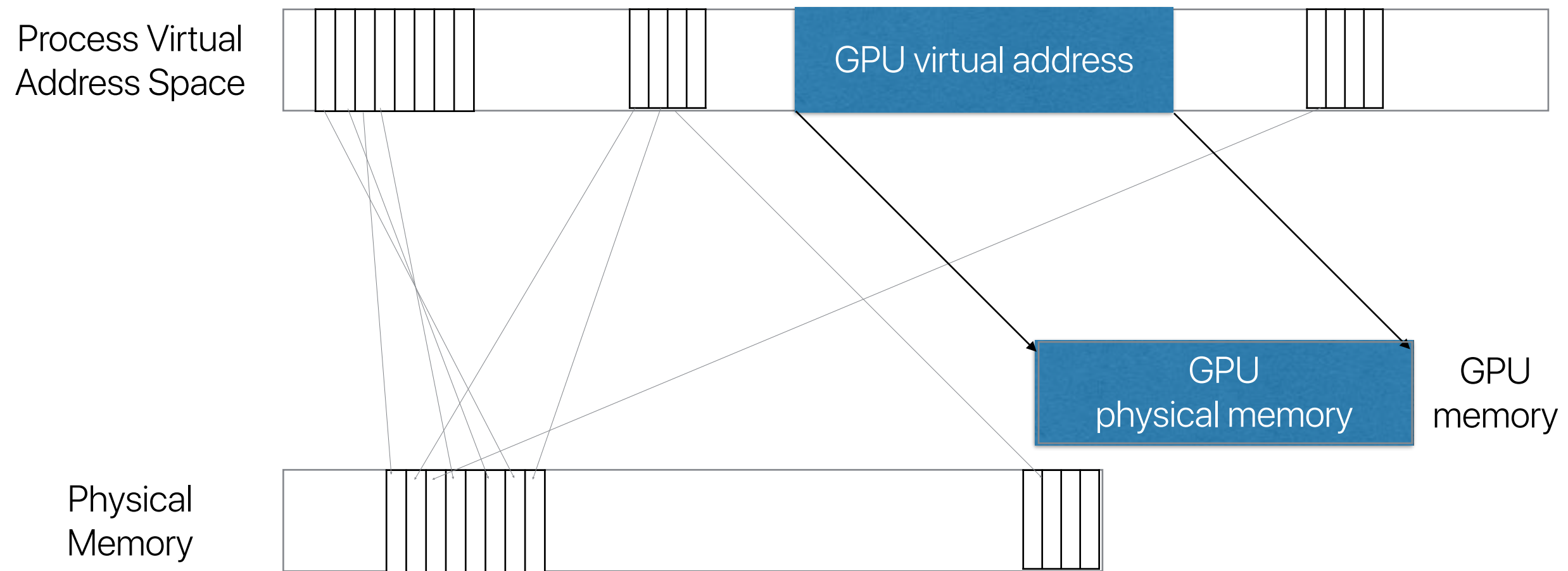
Direct Segment



Architecture overview



Nvidia's unified virtual memory



OS/Software support

- The operating system provides a primary region abstraction
 - 64-bit memory space is big enough to find a contiguous primary region
- The operating system allocates a physical memory region for the primary region using direct segment
 - extending libc function interface — think about your projects
- Optimizing the physical memory allocation
 - memory compaction
 - program architectural registers: base, offset, limit
 - growing/shrinking segments