

cs5965 Advanced OS Implementation

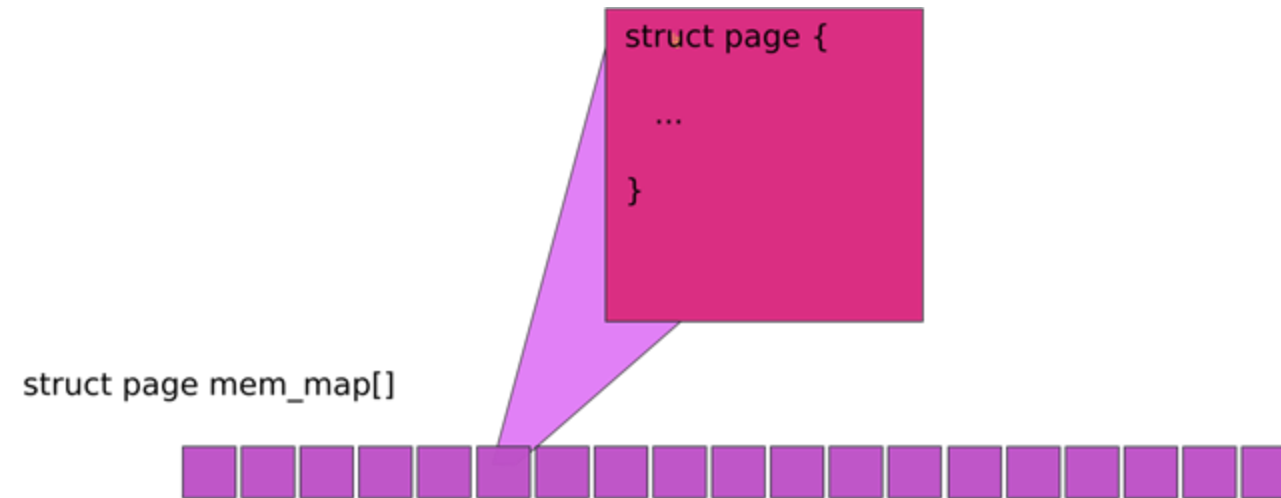
Lecture 05 – Memory management
Anton Burtsev

Physical memory



Physical memory

We need a smaller array to describe physical pages, e.g., `mem_map[]` in Linux



- Memory allocation

Simplest memory allocator

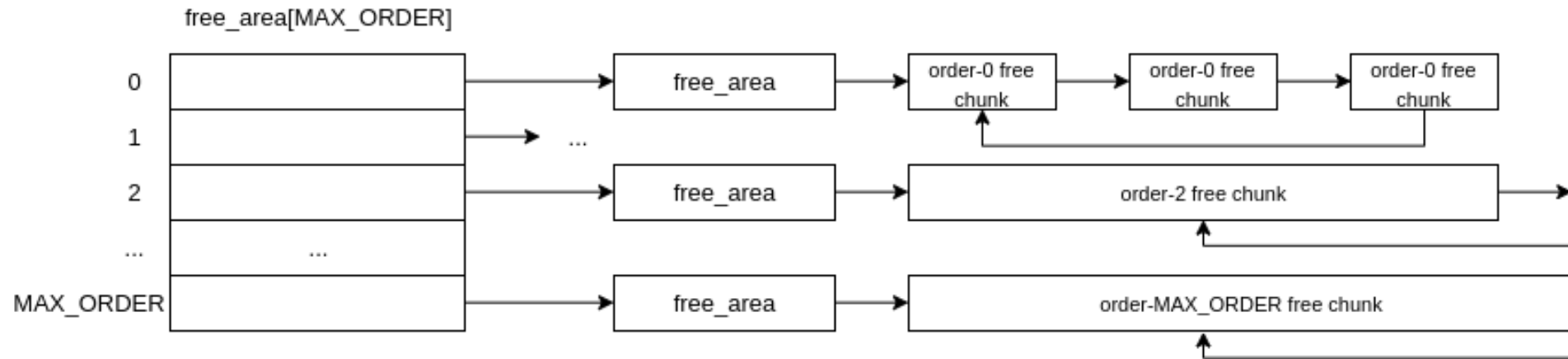
- Bitmap of all pages
- Bootmem allocator in Linux
- Allocation searches for an unused page
- Multiple sub-page allocations can be served from the same page by advancing a pointer
- Works ok, but what is the problem?

Boot memory allocator

- Bitmap of all pages
- Bootmem allocator in Linux
- Allocation searches for an unused page
- Multiple sub-page allocations can be served from the same page by advancing a pointer
- Works ok, but what is the problem?
- Linear scan of the bitmap
 - Too long

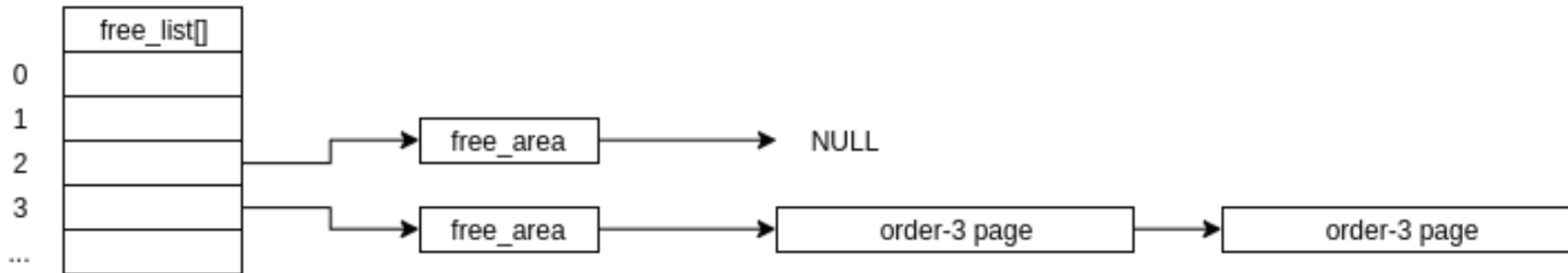
- Buddy:
- Physical Memory Allocator

Buddy memory allocator



Buddy memory allocator

`buddy_allocate_me(0x4000 bytes)` → 0x4 pages → an order-2 page



What's wrong with buddy?

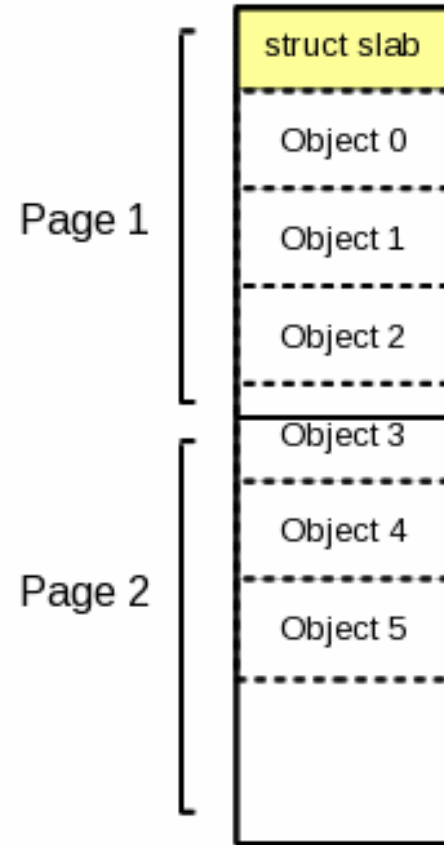
What's wrong with buddy?

- Buddy allocator is ok for large allocations
 - E.g. 1 page or more
- But what about small allocations?
- Buddy uses the whole page for a 4 bytes allocation
 - Wasteful
- Buddy is still slow for short-lived objects

- Slab:
- Allocator for object of a fixed size

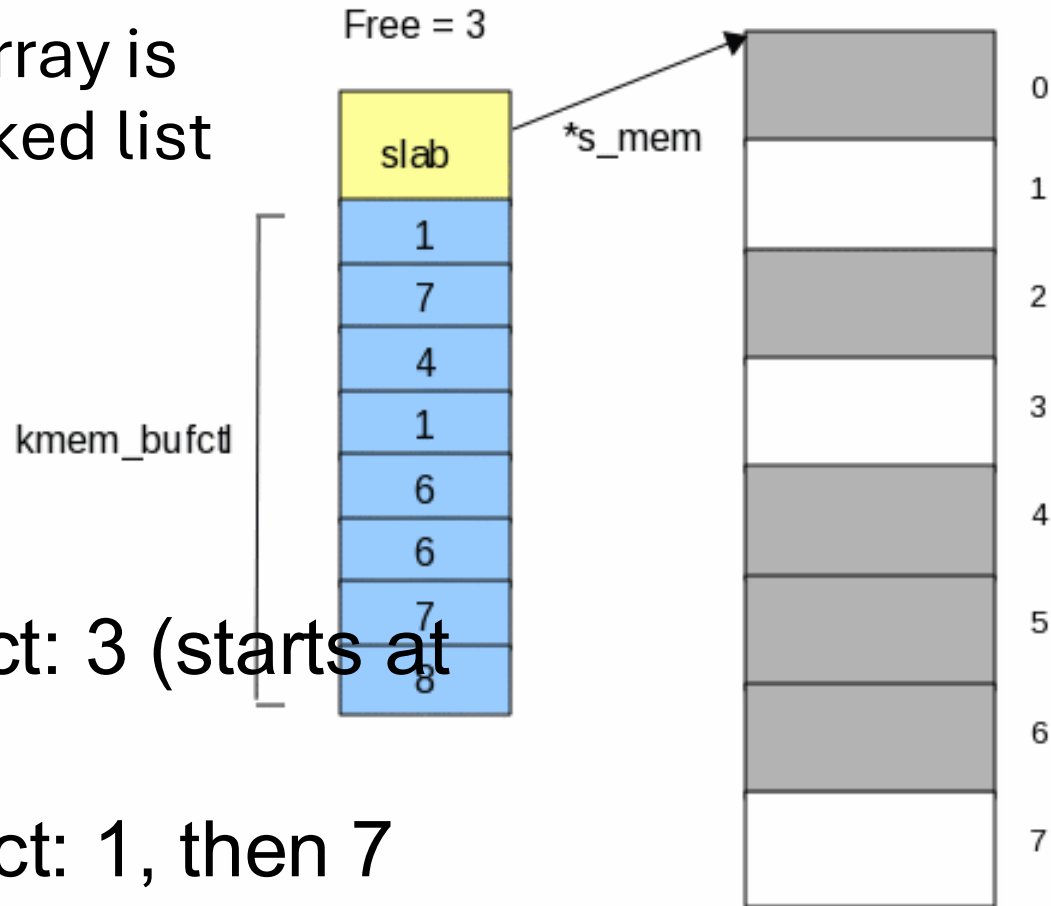
Slab

- A 2 page slab with 6 objects



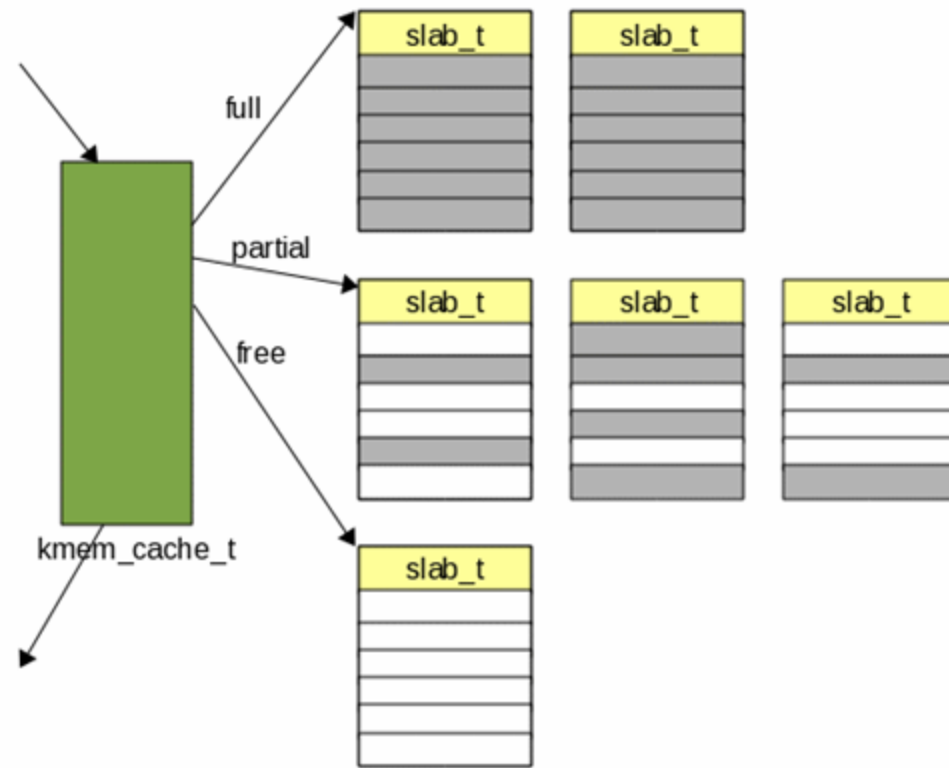
Keeping track of free objects

- kmem_bufctl array is effectively a linked list



- First free object: 3 (starts at 0)
- Next free object: 1, then 7

A cache is formed out of slabs



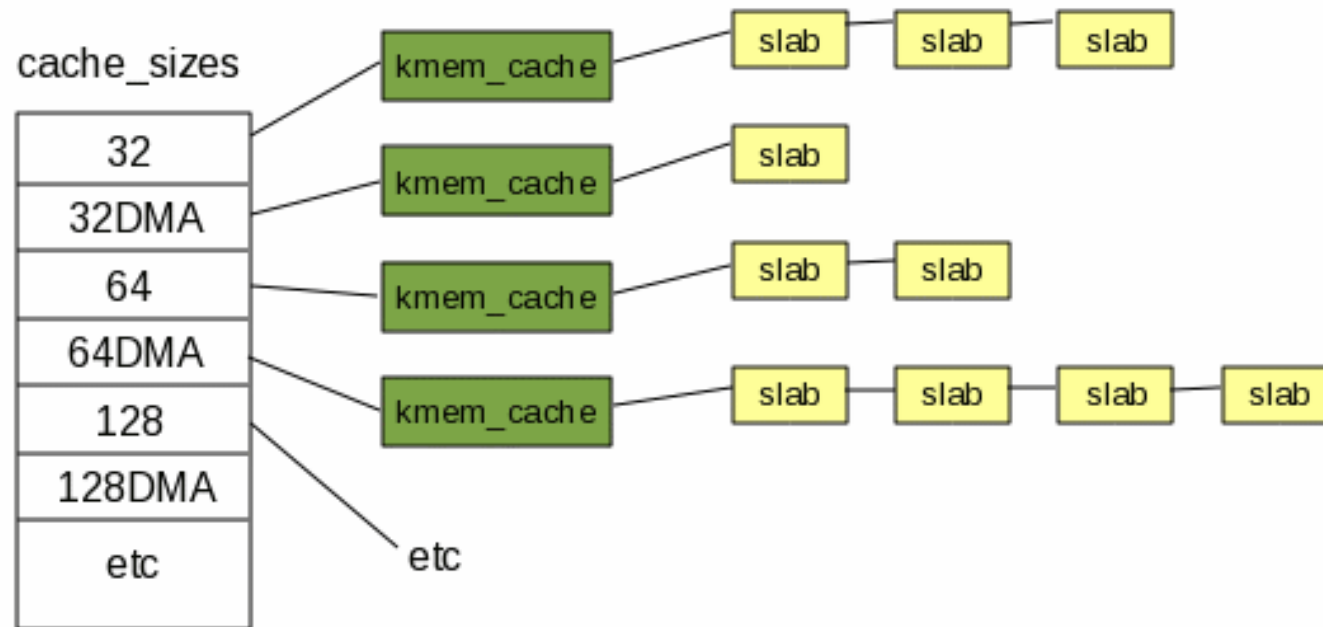
Slab is fine, but what's wrong?

Slab is fine, but what's wrong?

- We can only allocate objects of one size

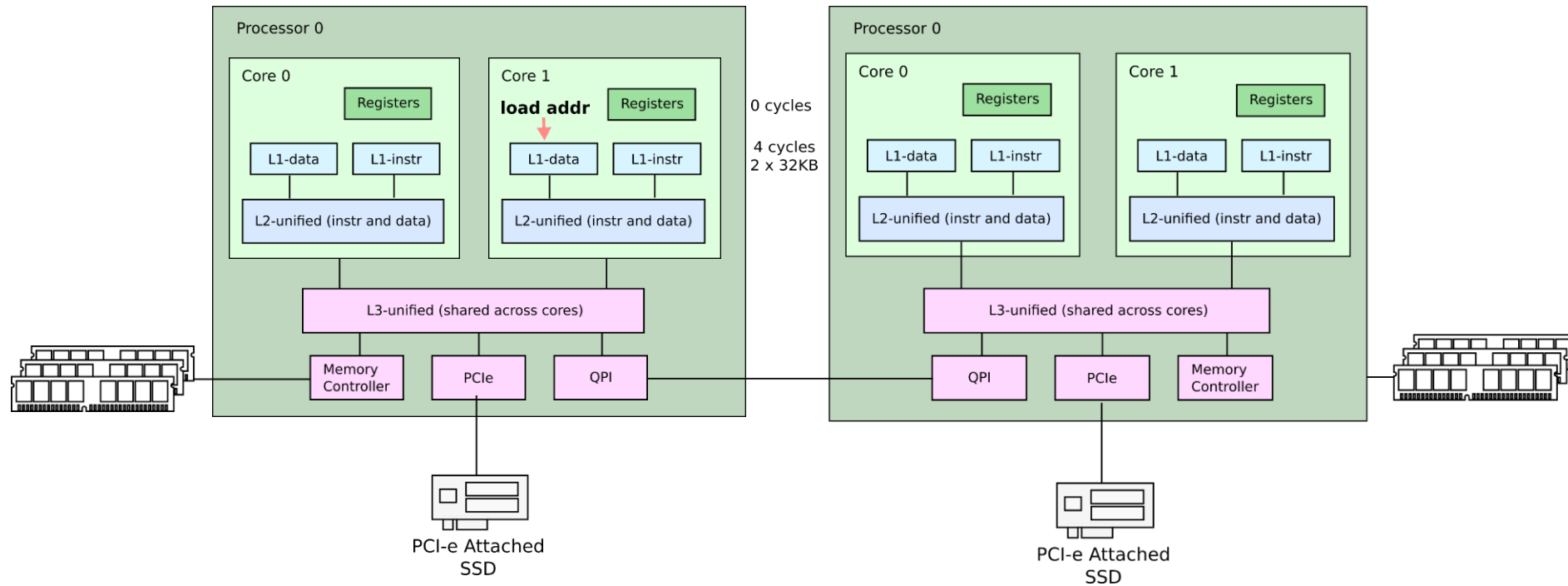
Kmalloc(): variable size objects

- A table of caches
- Size: 32, 64, 128, etc.

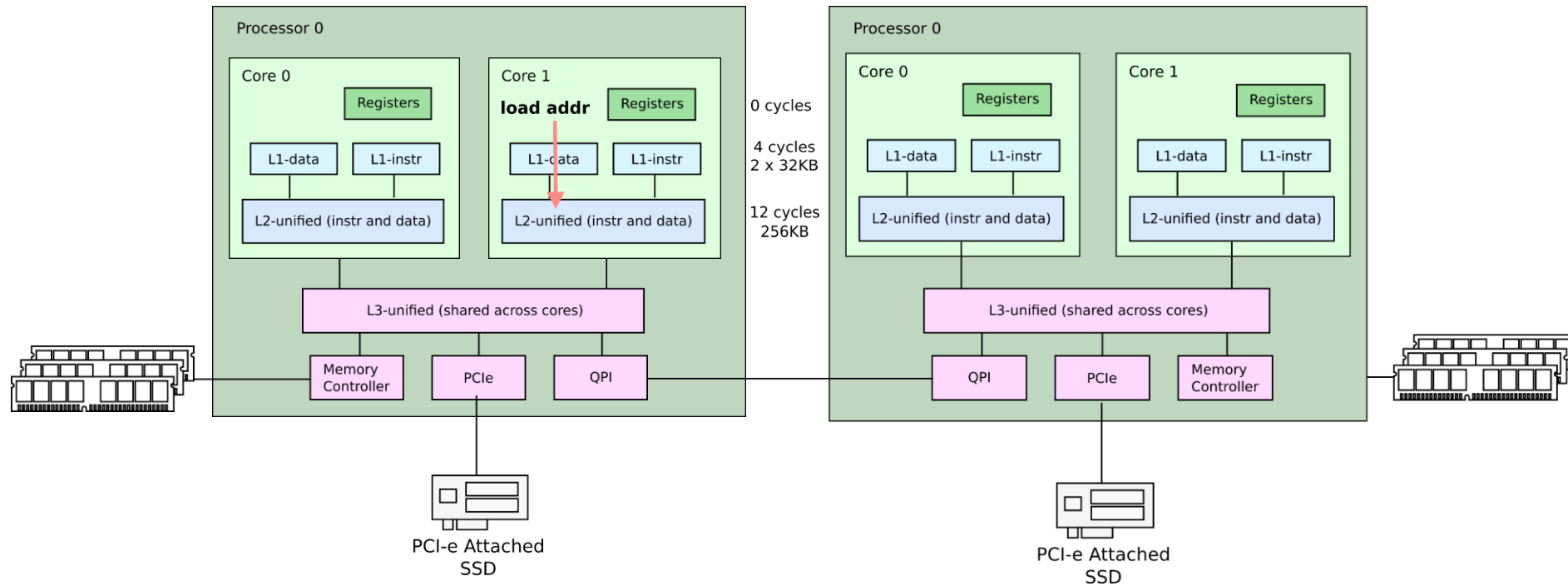


Non-Uniform Memory Access

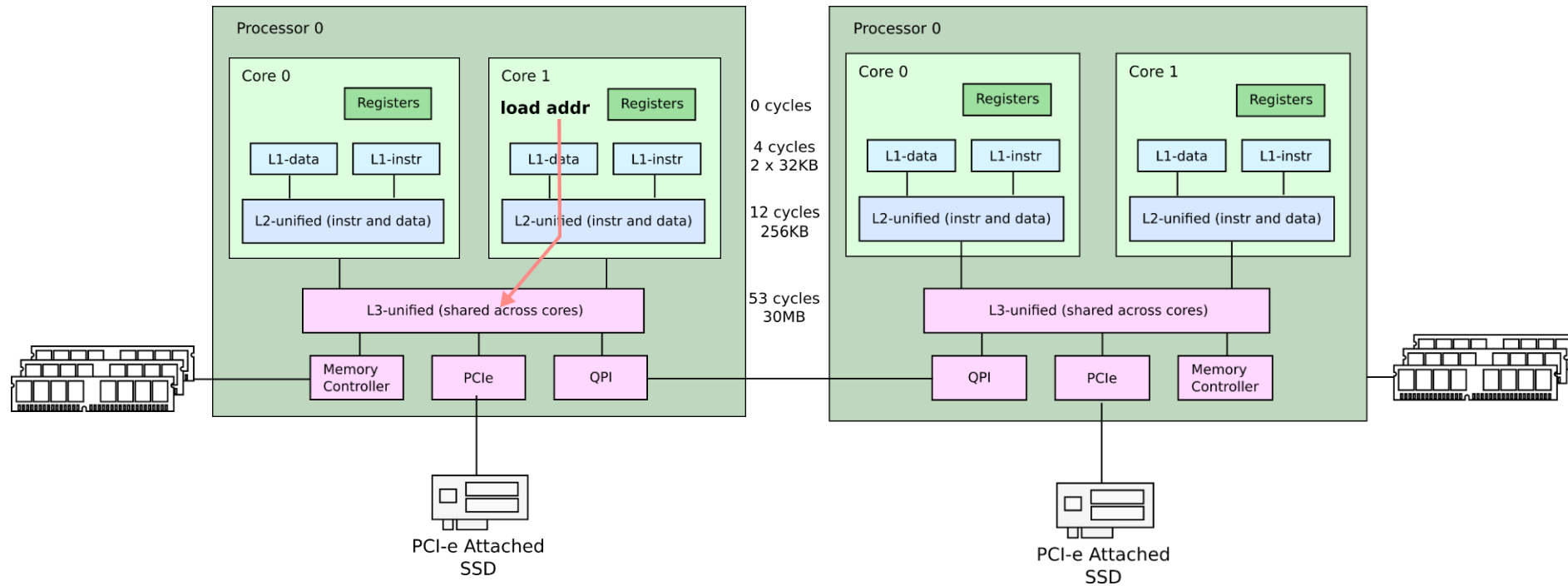
Latencies: load from local L1



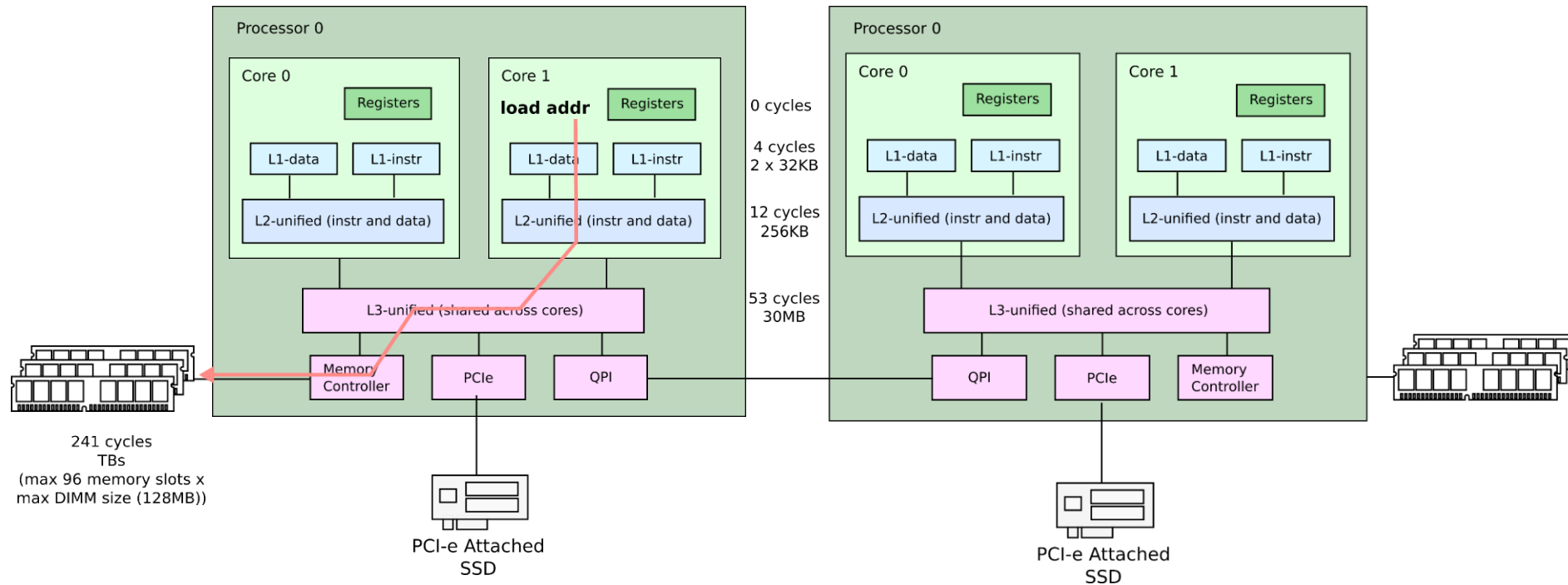
Latencies: load from local L2



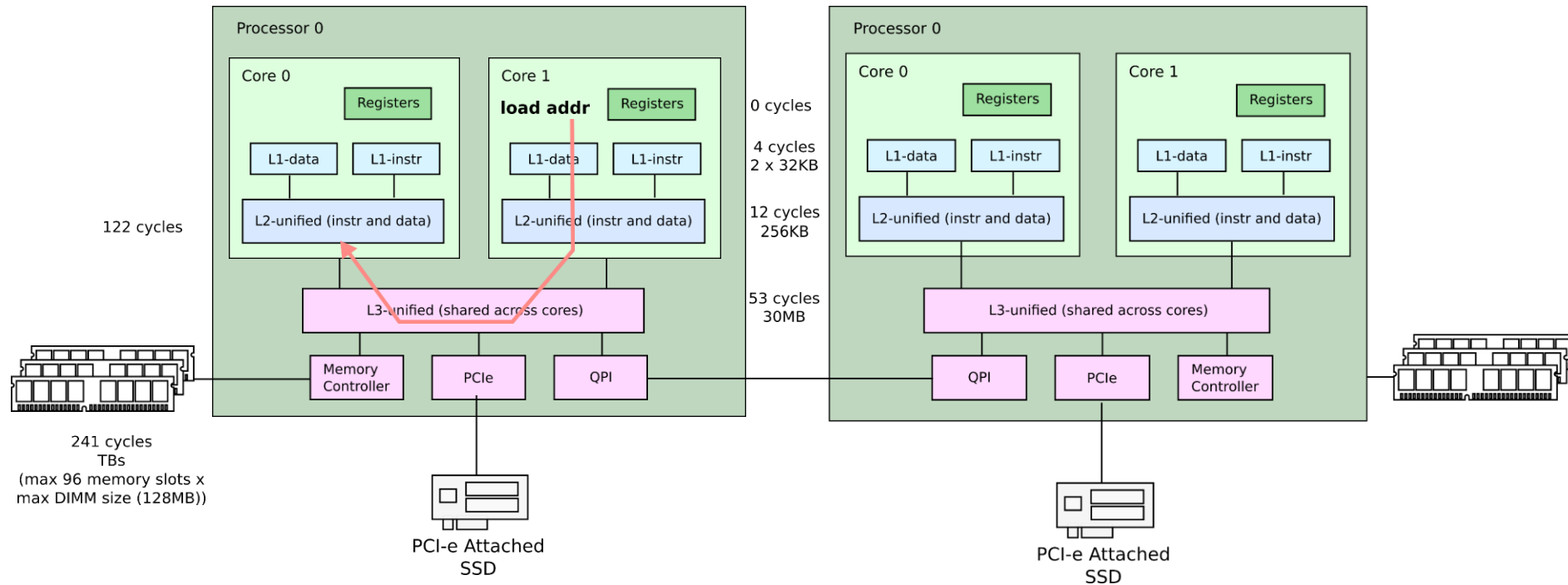
Latencies: load from local L3



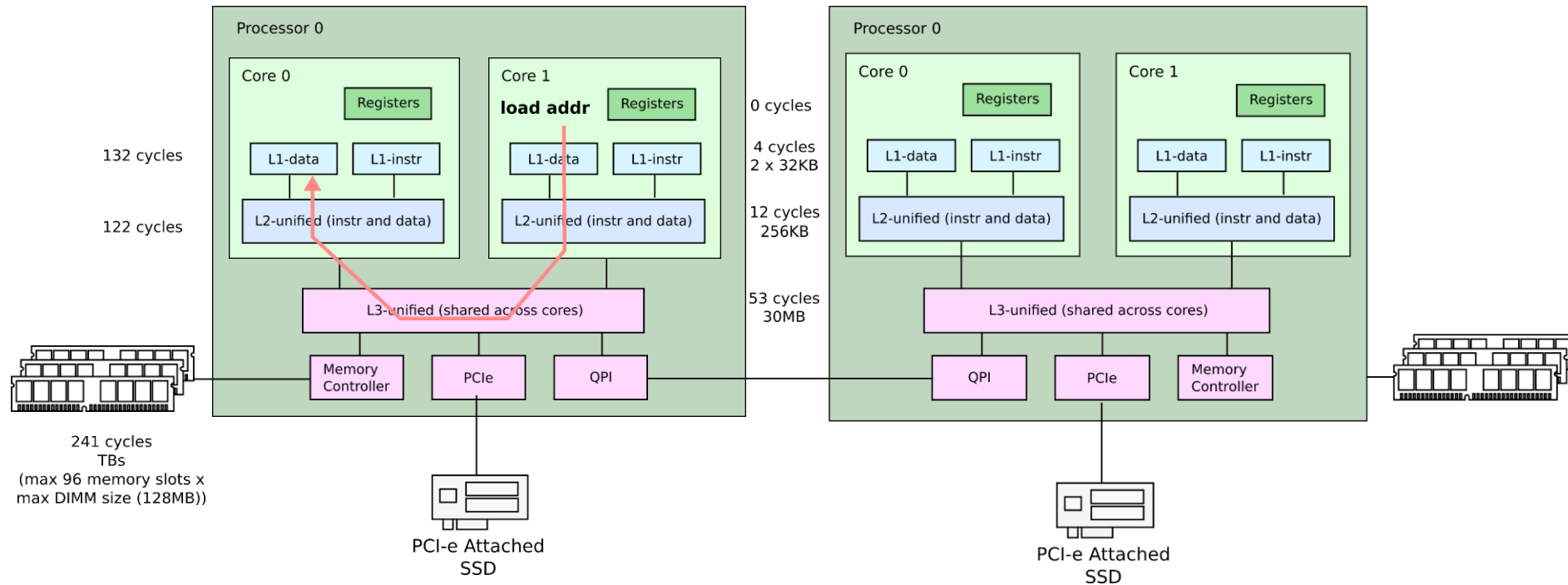
Latencies: load from local memory



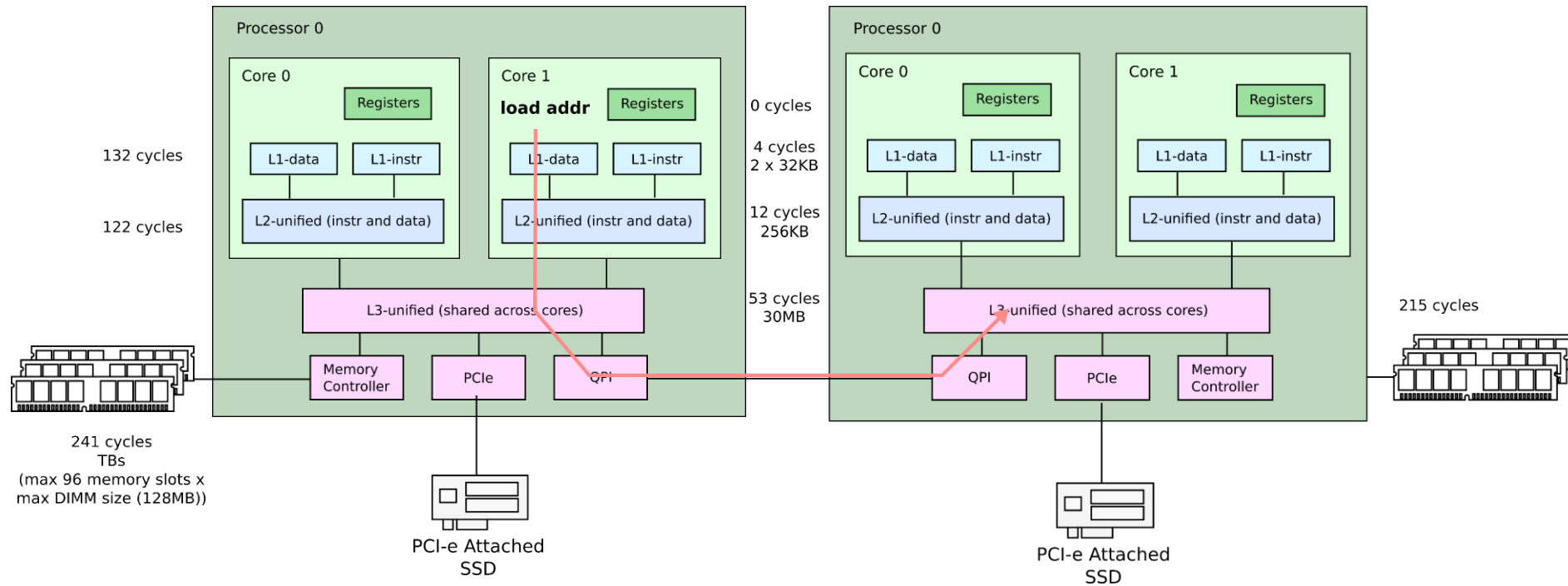
Latencies: load from same die core's L2



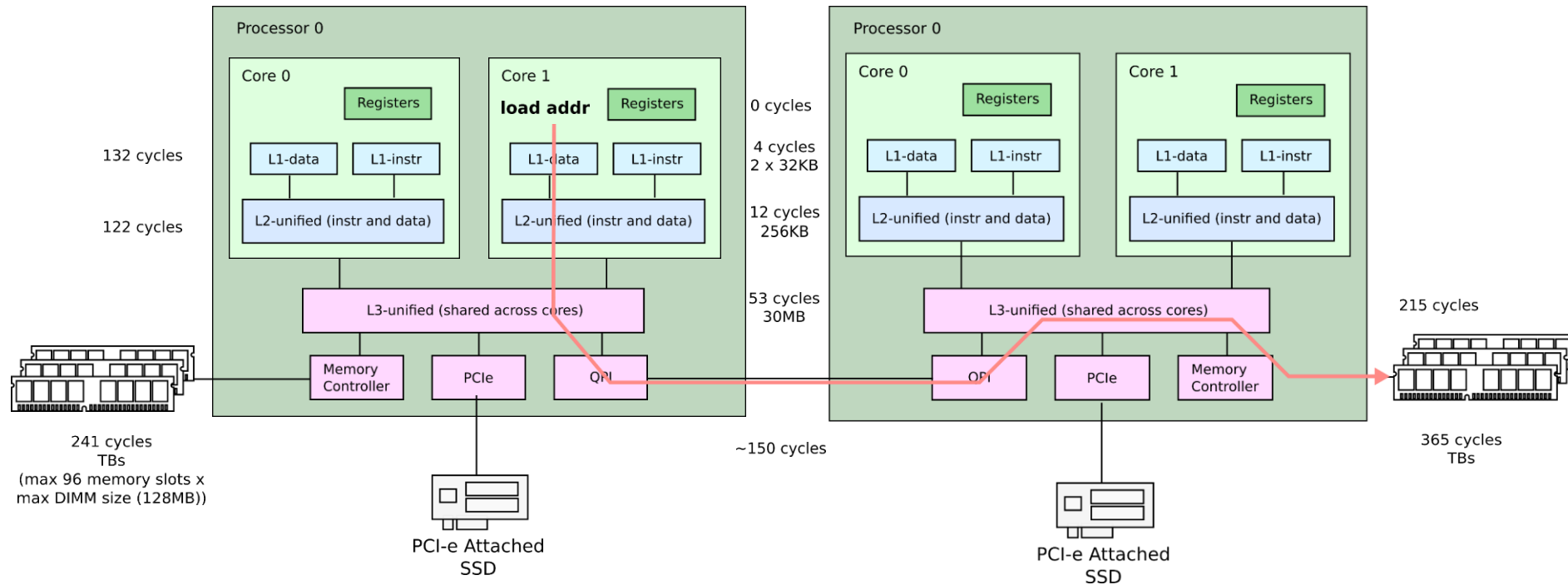
Latencies: load from same die core's L1



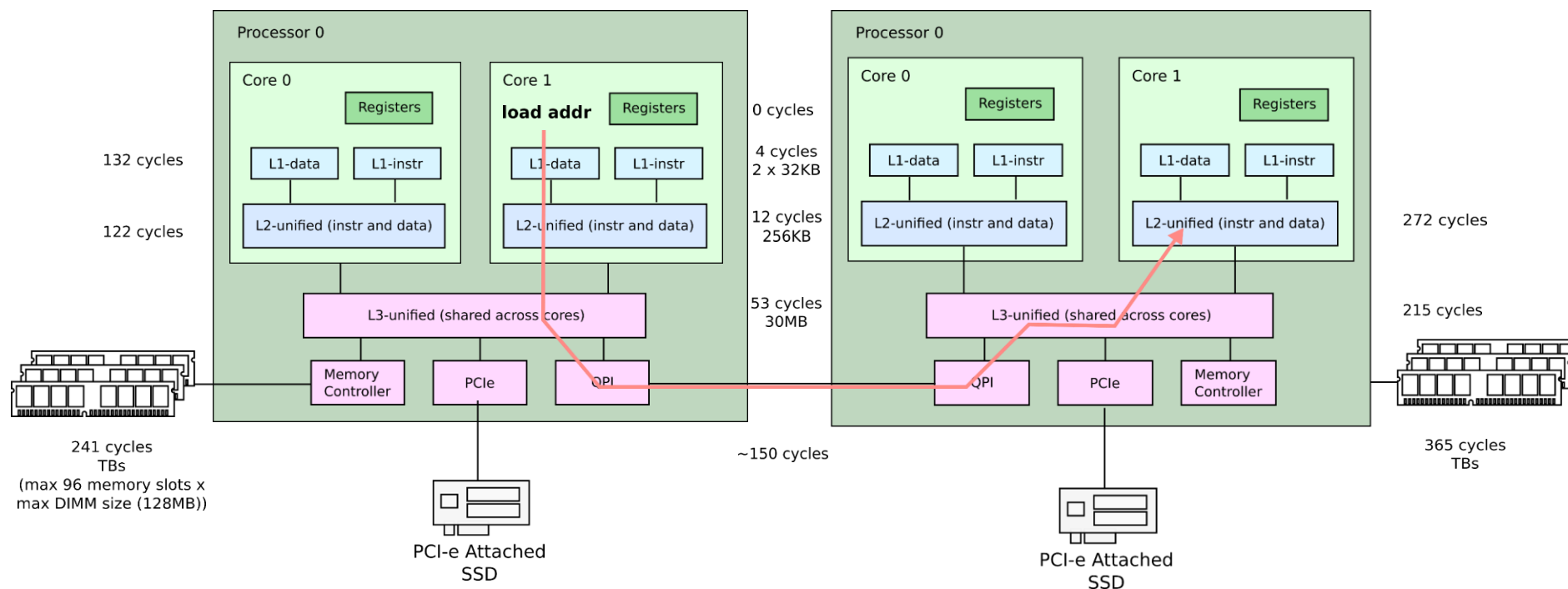
Latencies: load from remote L3



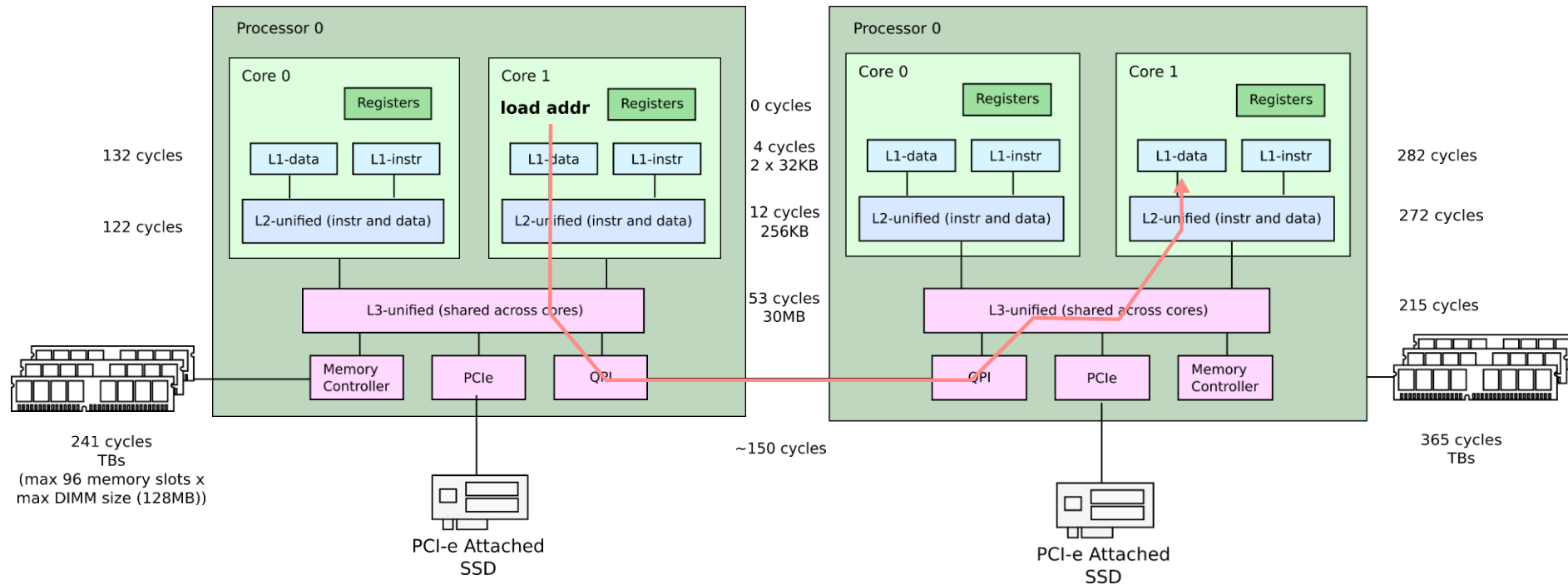
Latencies: load from remote memory



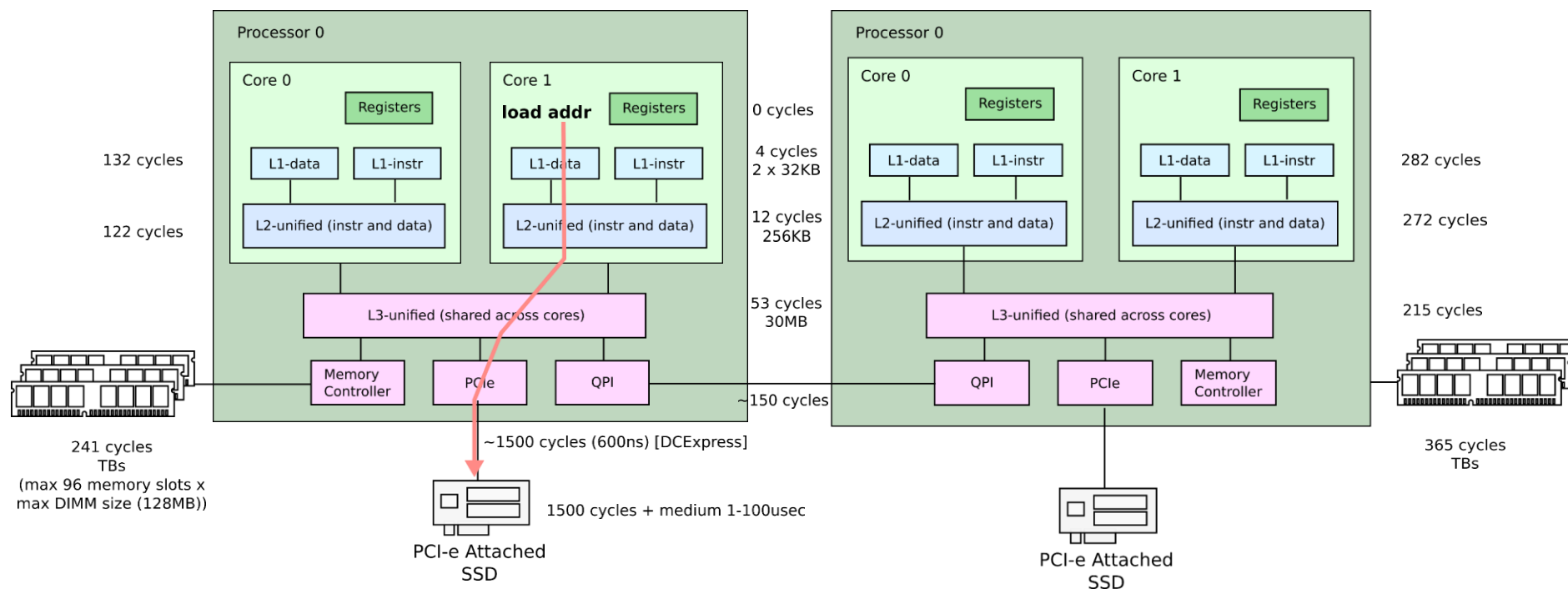
Latencies: load from remote L2



Latencies: load from remote L2



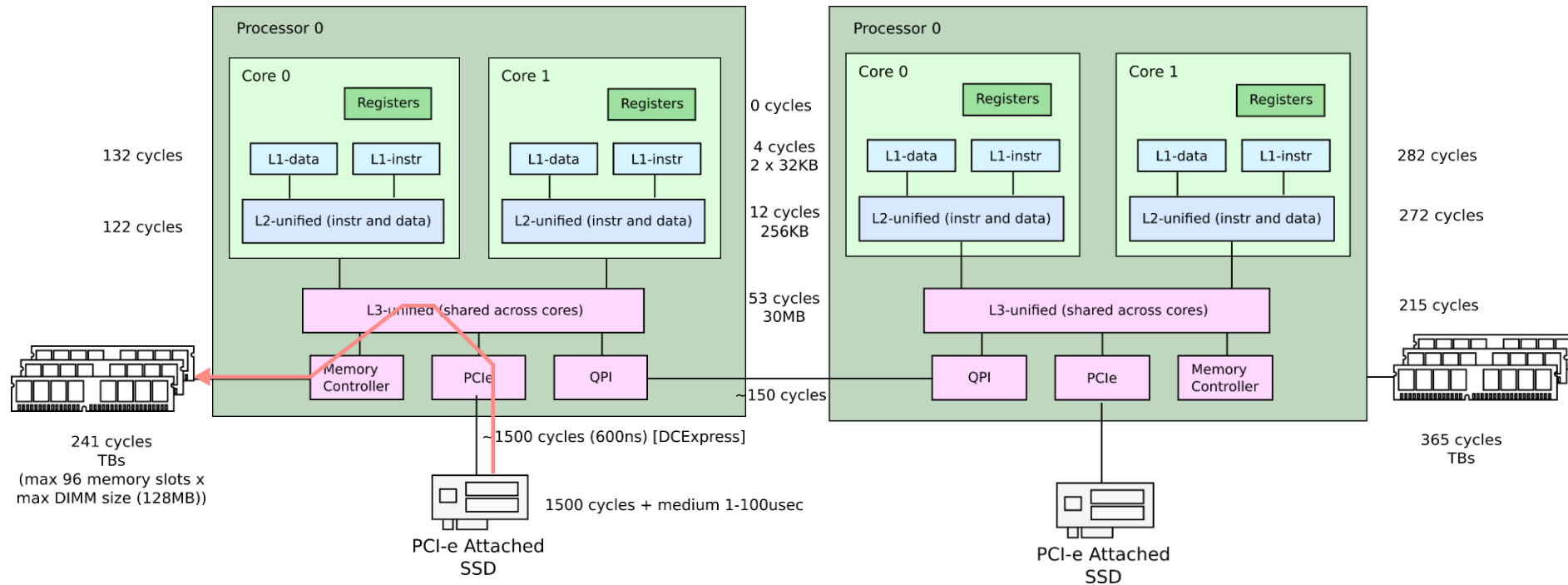
Latencies: PCIe round-trip



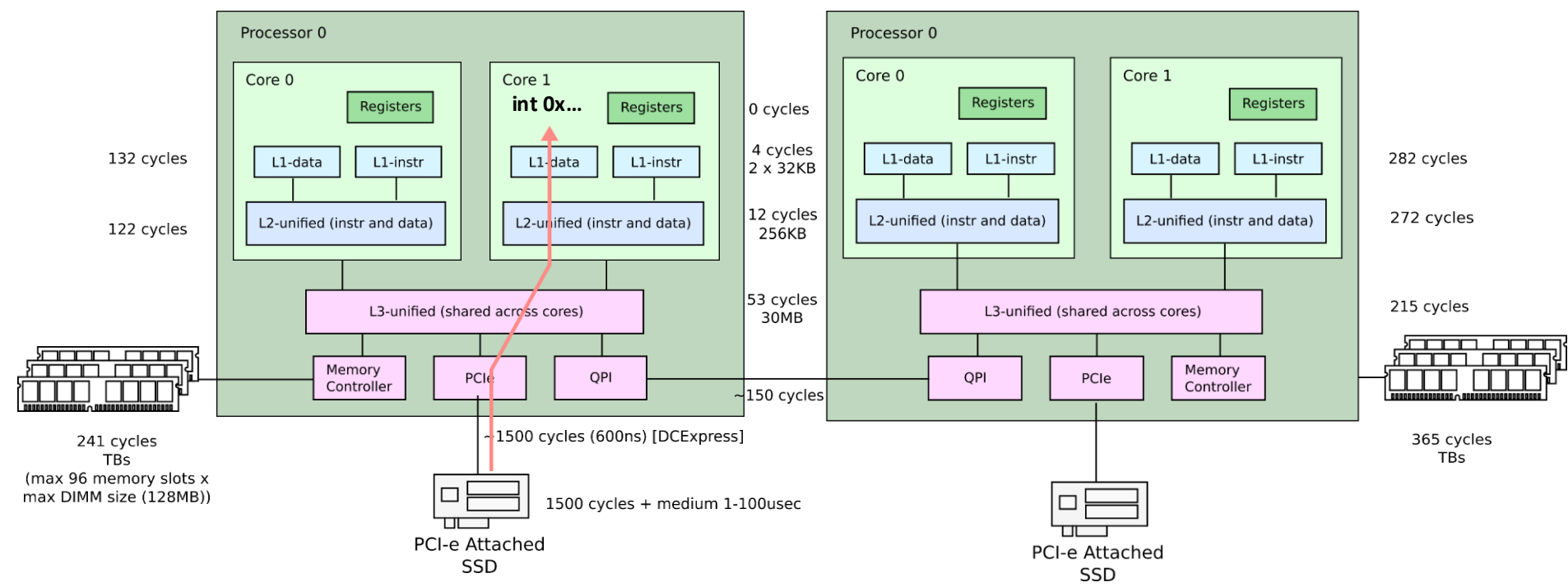
Device I/O

- Essentially just sending data to and from external devices
- Modern devices communicate over PCIe
 - Well there are other popular buses, e.g., USB, SATA (disks), etc.
 - Conceptually they are similar
- Devices can
 - Read memory
 - Send interrupts to the CPU

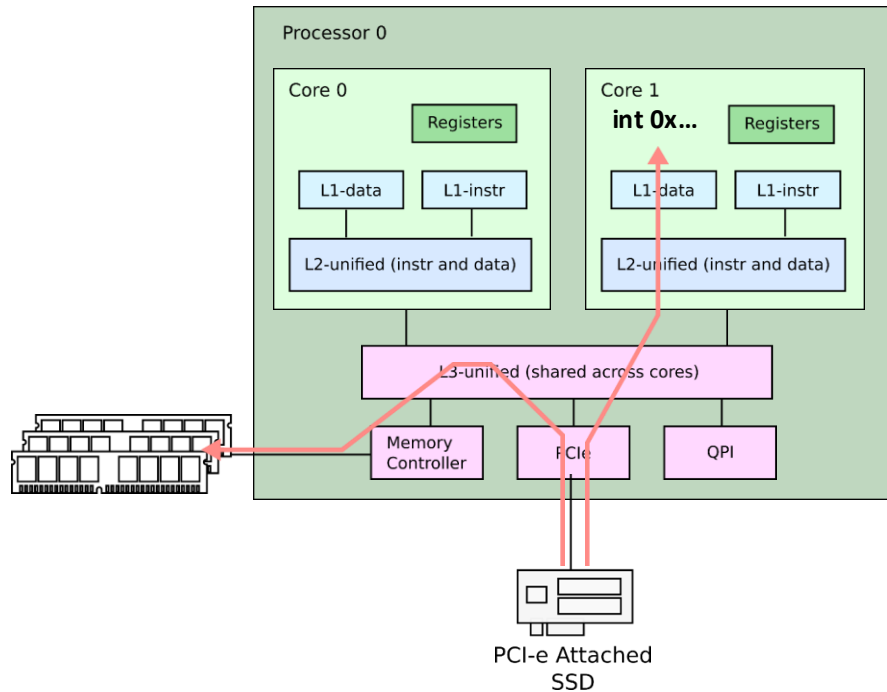
Direct memory access



Interrupts

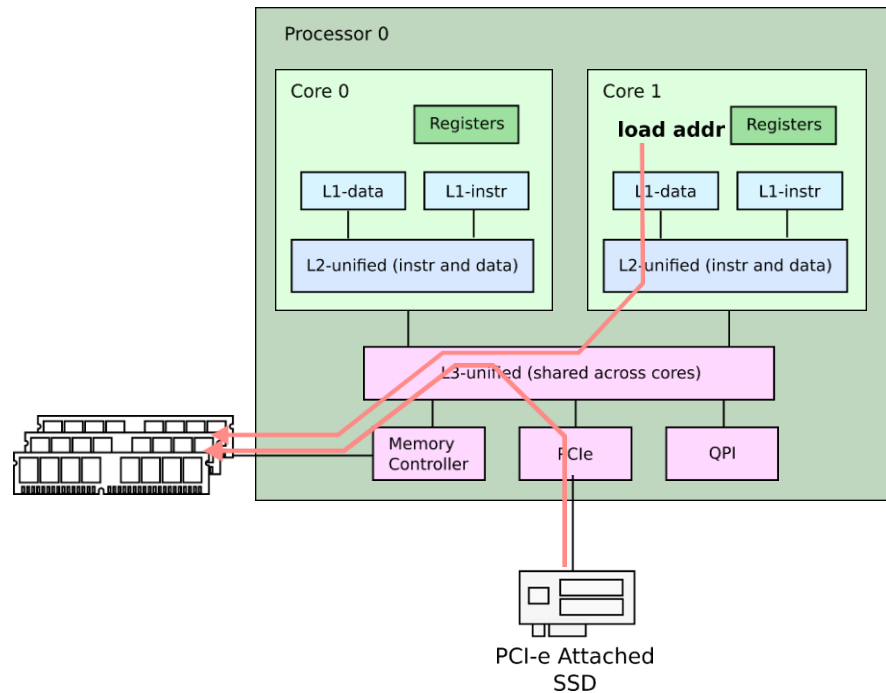


Device I/O



- Write incoming data in memory, e.g.,
 - Network packets
 - Disk requests, etc.
- Then raise an interrupt to notify the CPU
 - CPU starts executing interrupt handler
 - Then reads incoming packets from memory

Device I/O (polling mode)



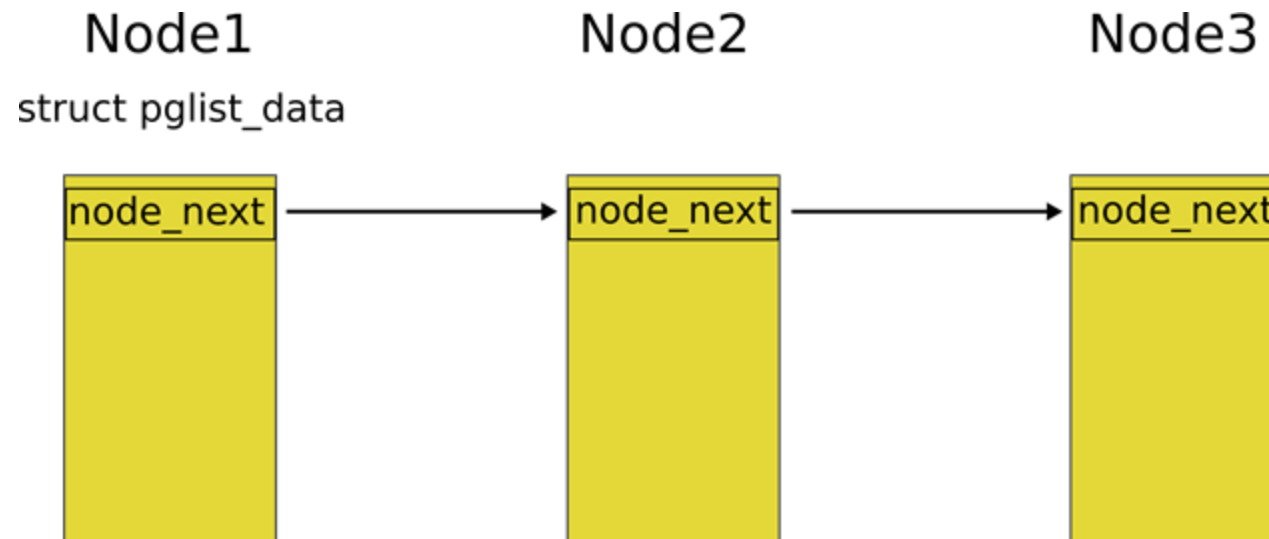
- Alternatively the CPU has to check for incoming data in memory periodically
 - Or poll
- Rationale
 - Interrupts are expensive

Uniform and non-uniform memory access

- Parts of memory can be faster than others

Nodes

- Attempt to allocate memory from the current node
- Fall back to the next node in list
 - If ran out of local memory



pglist_data represents a node

<mmzone.h>

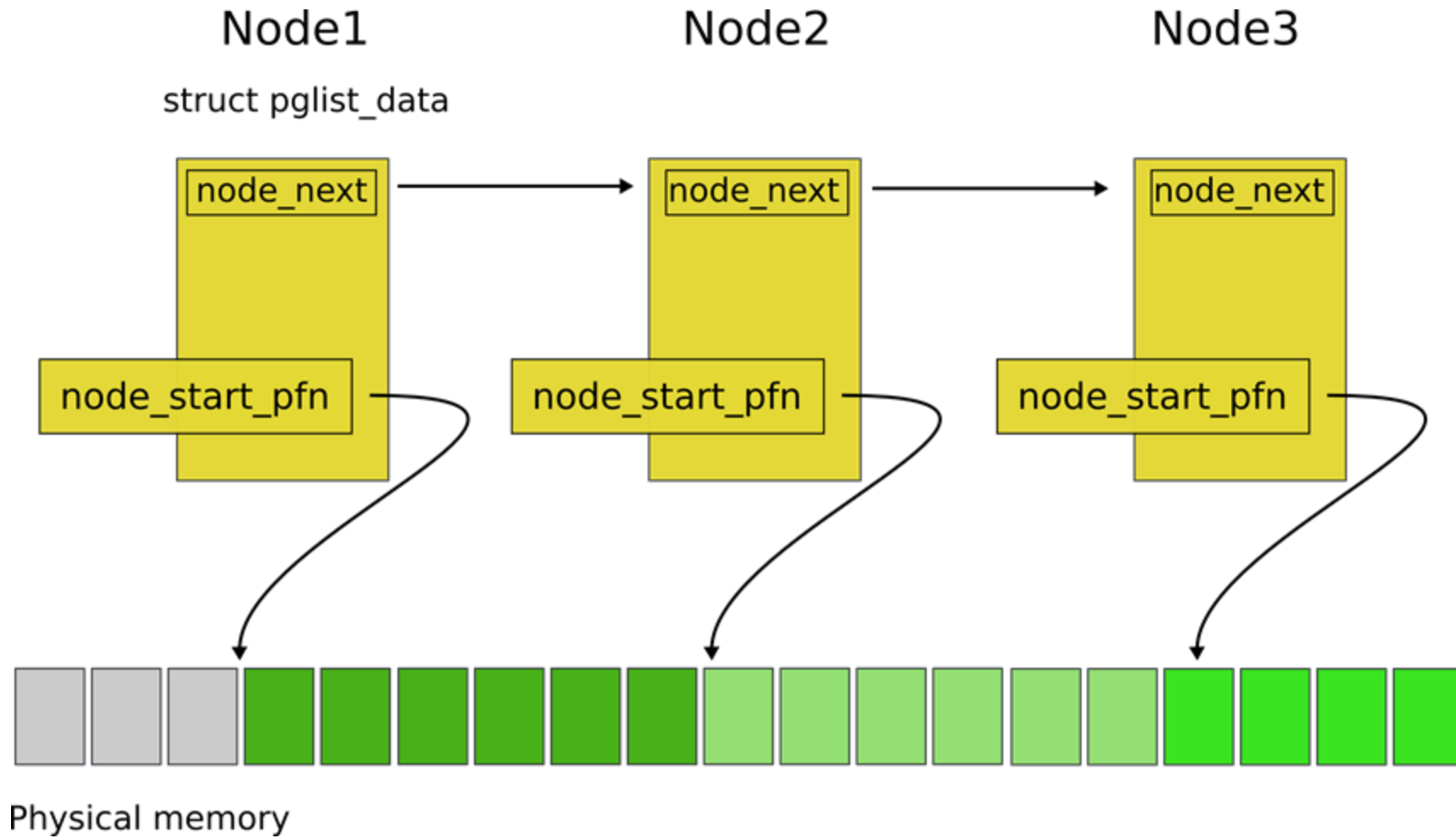
```
typedef struct pglist_data {
    struct zone node_zones[MAX_NR_ZONES];
    struct zonelist node_zonelists[MAX_ZONELISTS];
    int nr_zones;
    struct page *node_mem_map;
    struct bootmem_data *bdata;

    unsigned long node_start_pfn;
    unsigned long node_present_pages; /* total number of physical pages */
    unsigned long node_spanned_pages; /* total size of physical page
                                       range, including holes */

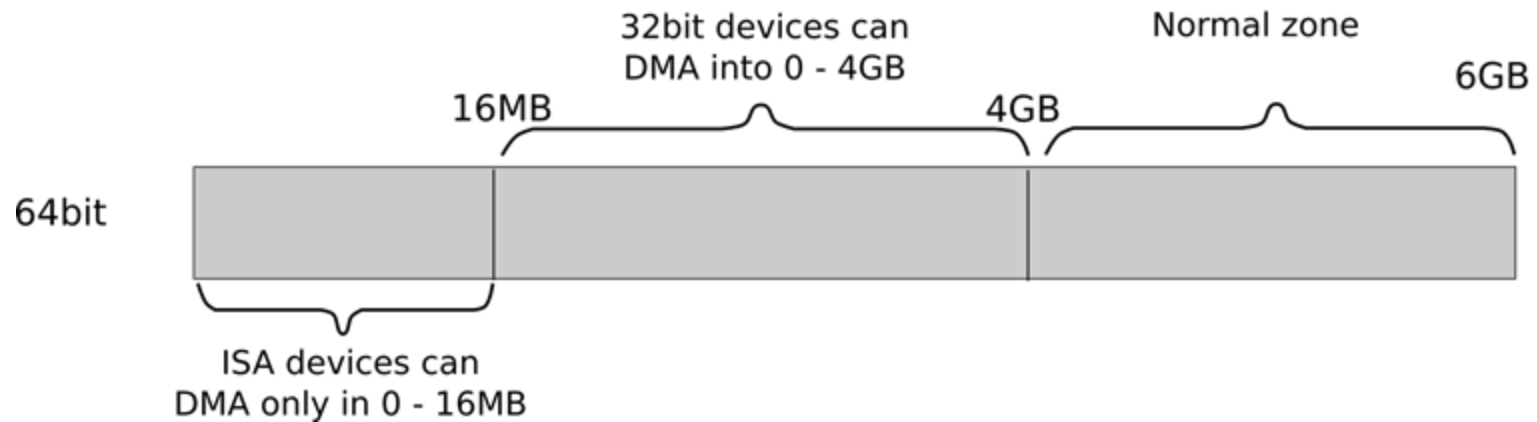
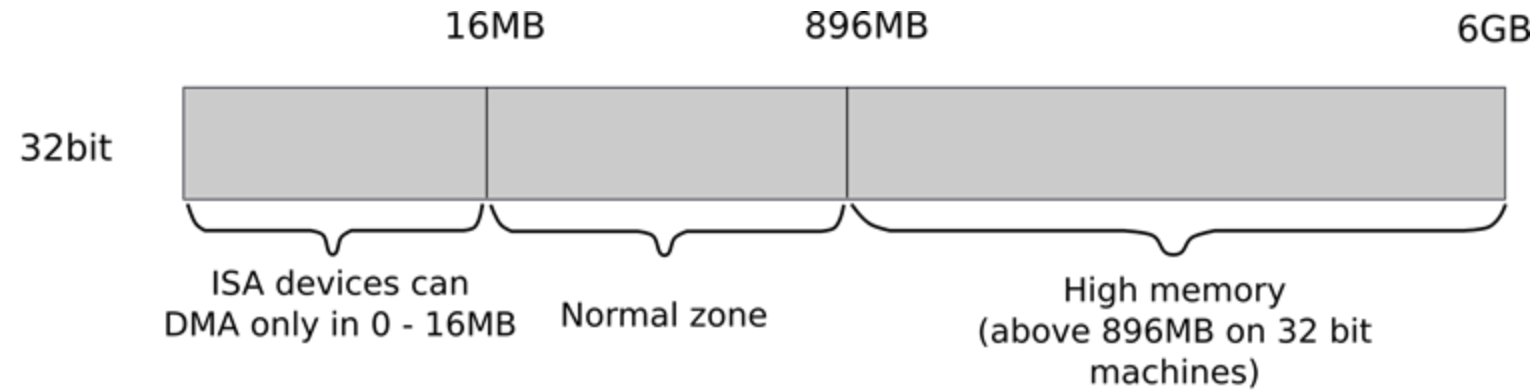
    int node_id;
    struct pglist_data *pgdat_next;
    wait_queue_head_t kswapd_wait;
    struct task_struct *kswapd;
    int kswapd_max_order;
} pg_data_t;
```

<https://elixir.bootlin.com/linux/v6.10.6/source/include/linux/mmzone.h#L1277>

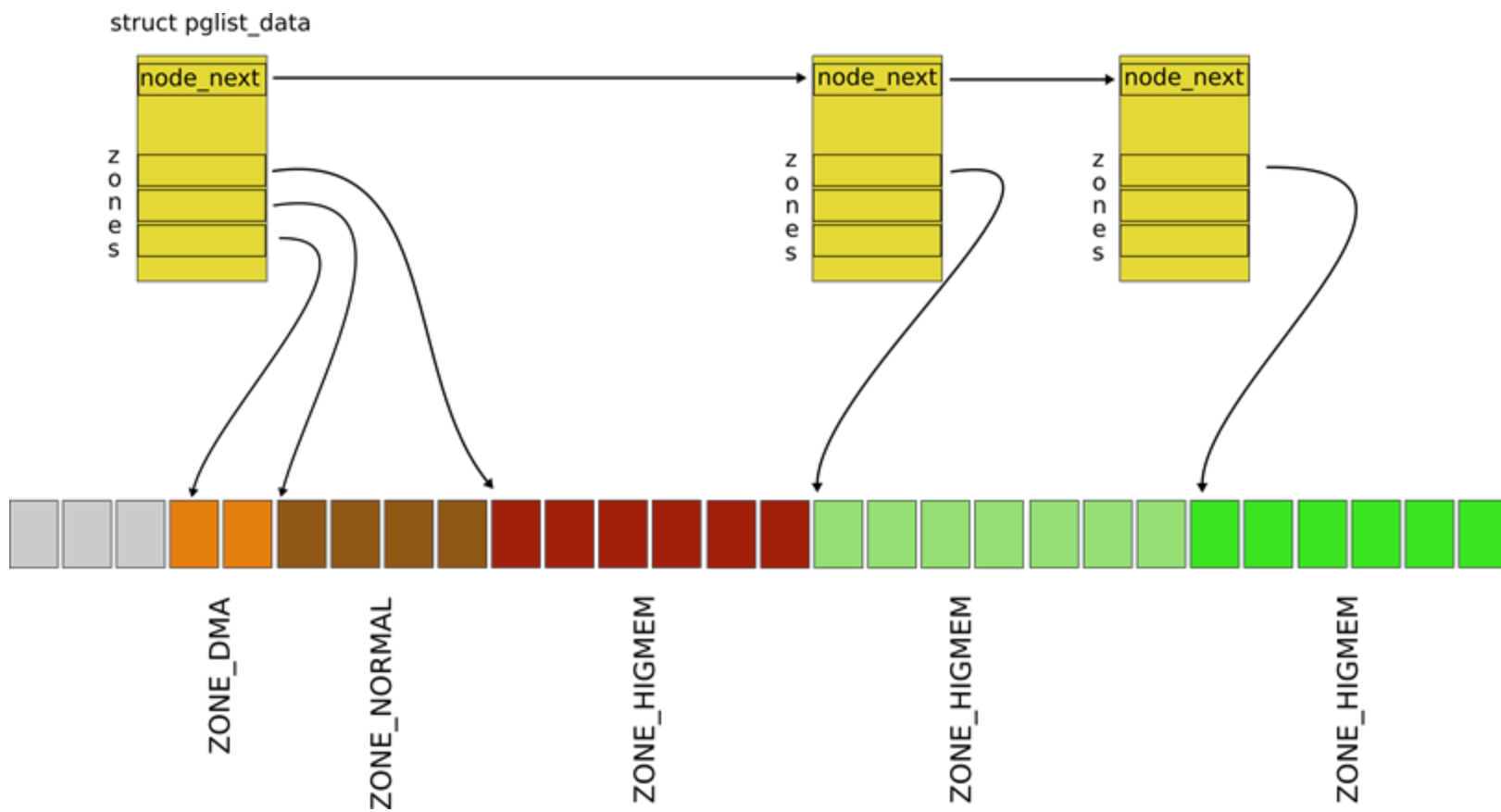
Nodes



Zones



Zones



Thank you!



Are you rewriting
our kernel in Rust?

Nah... let's just verify it