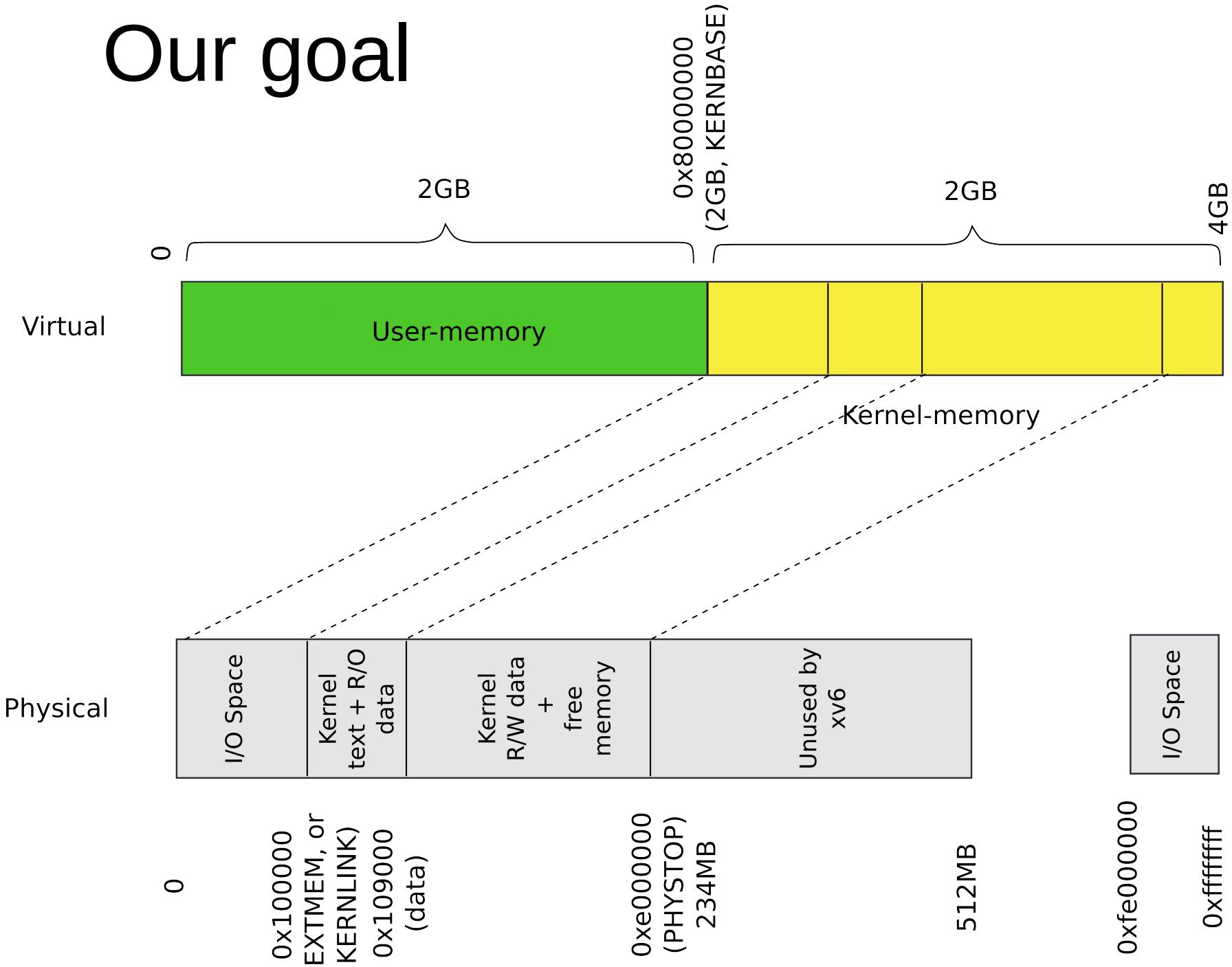


143A: Principles of Operating Systems

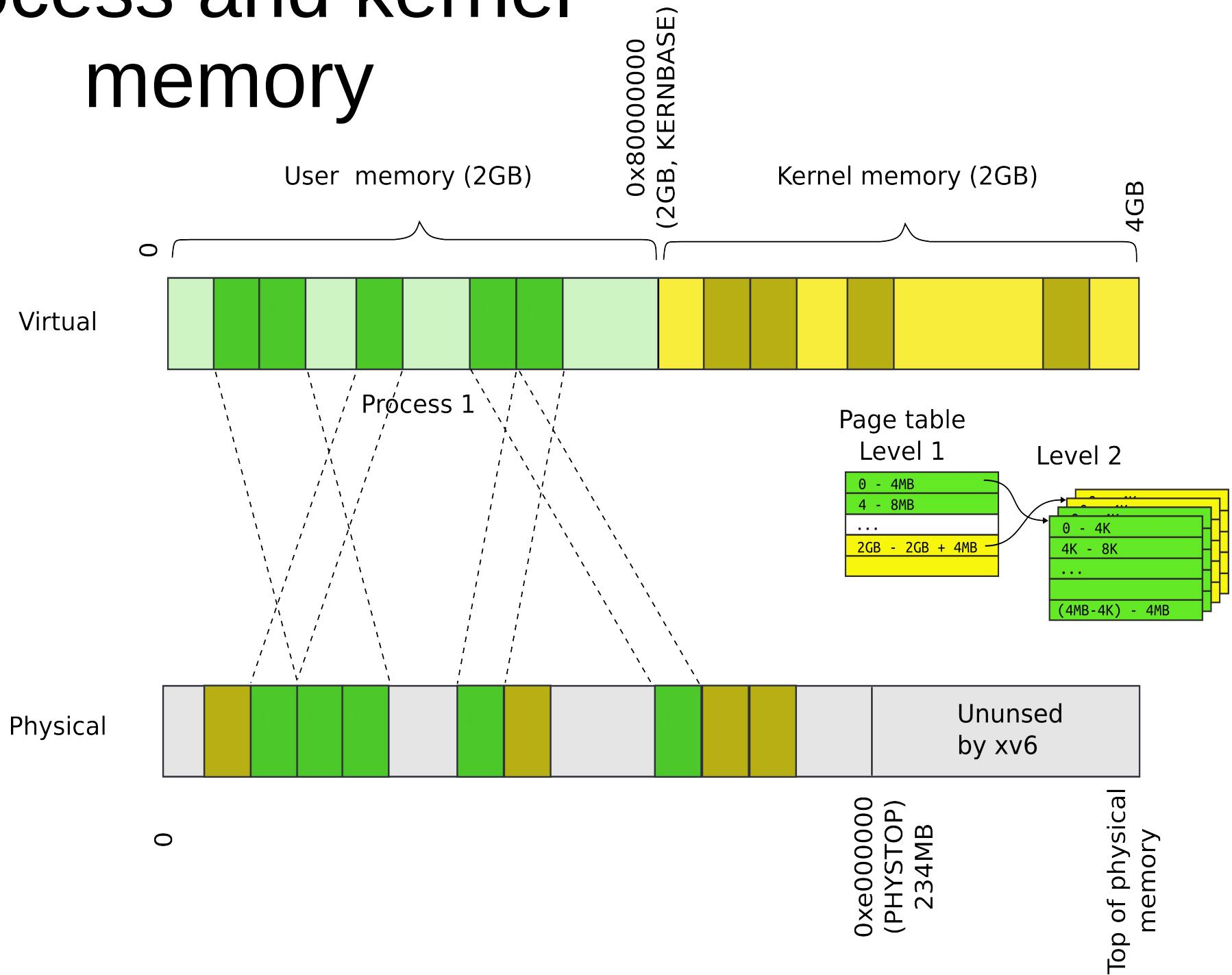
Lecture 09: Address spaces

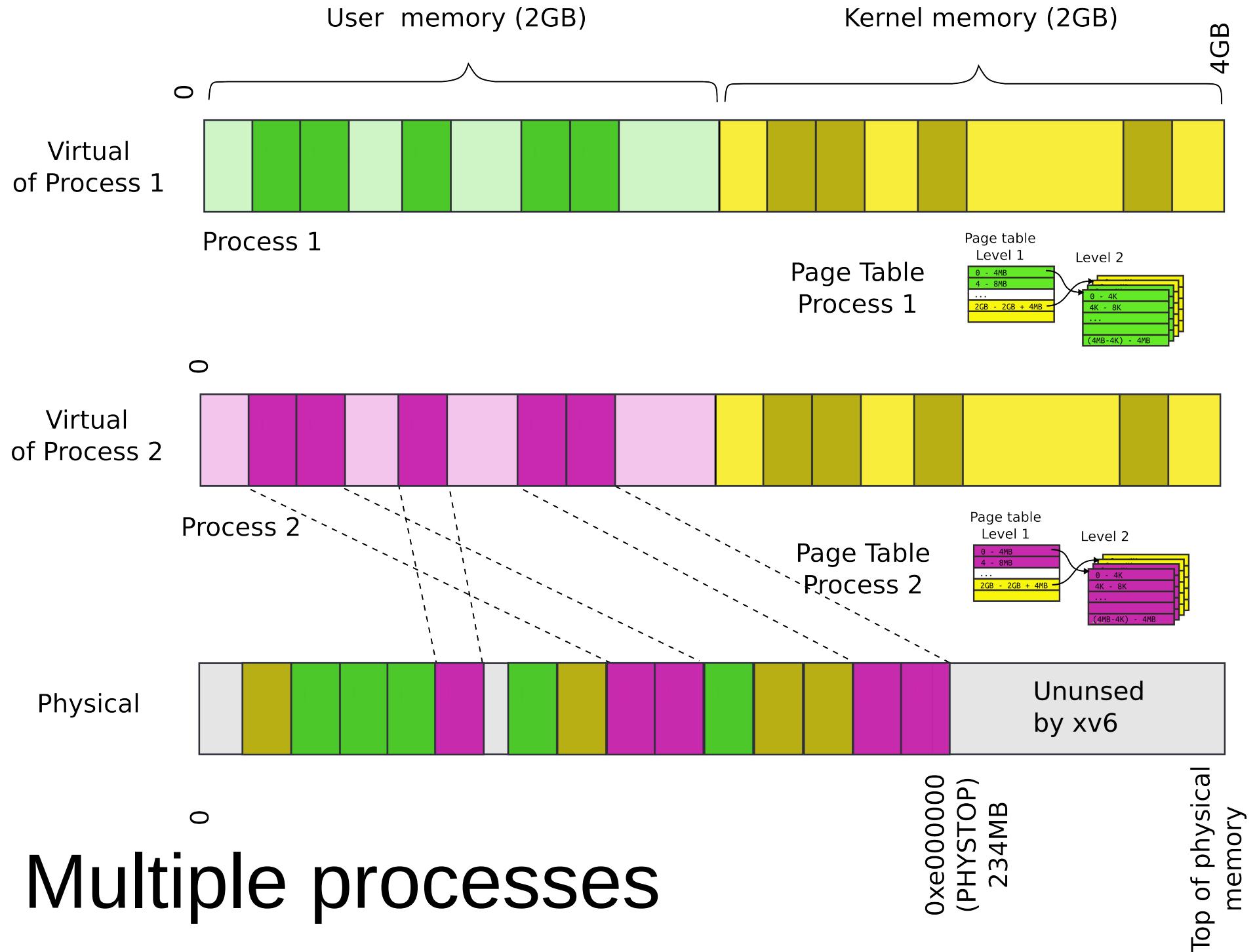
Anton Burtsev
January, 2017

Our goal

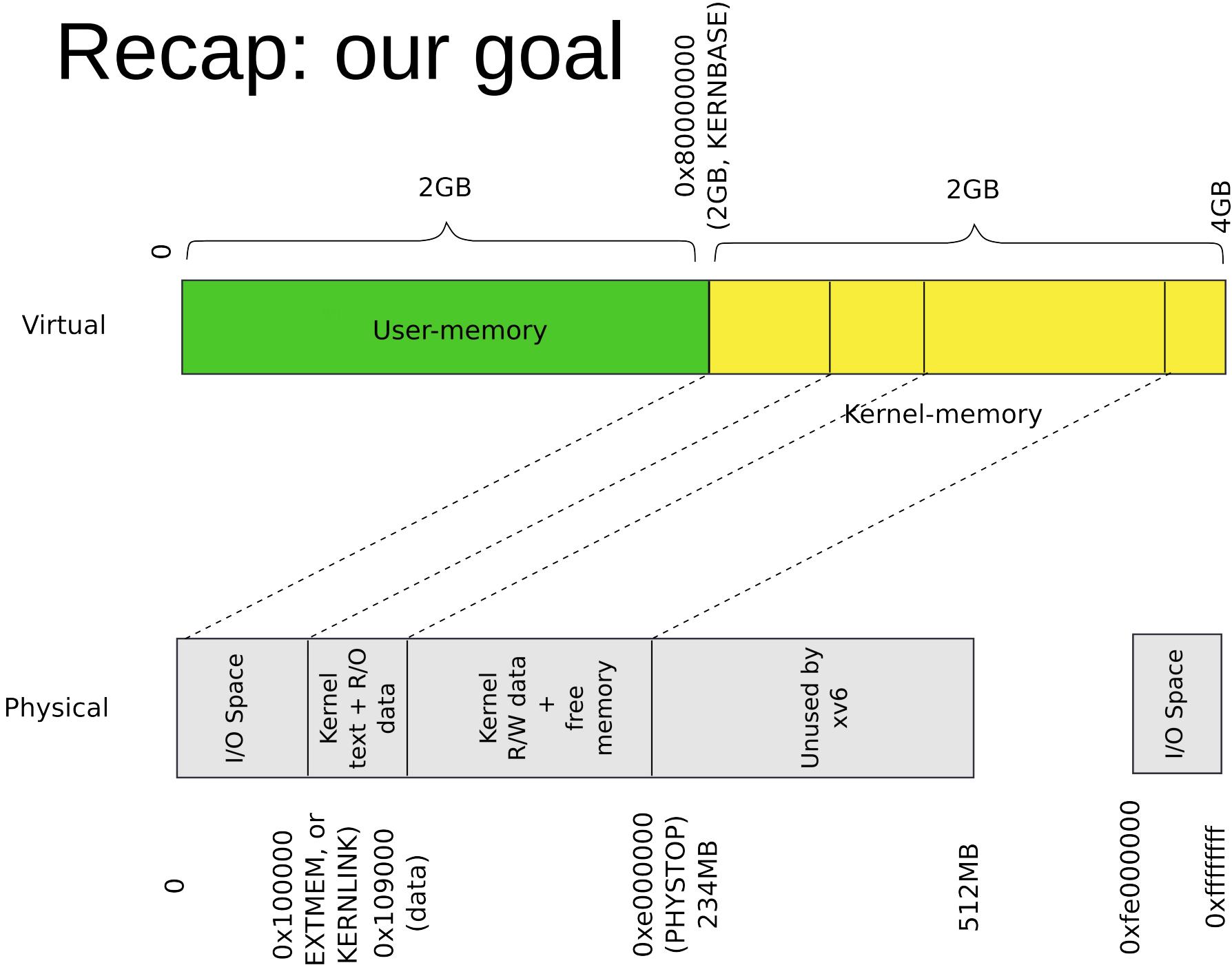


Process and kernel memory

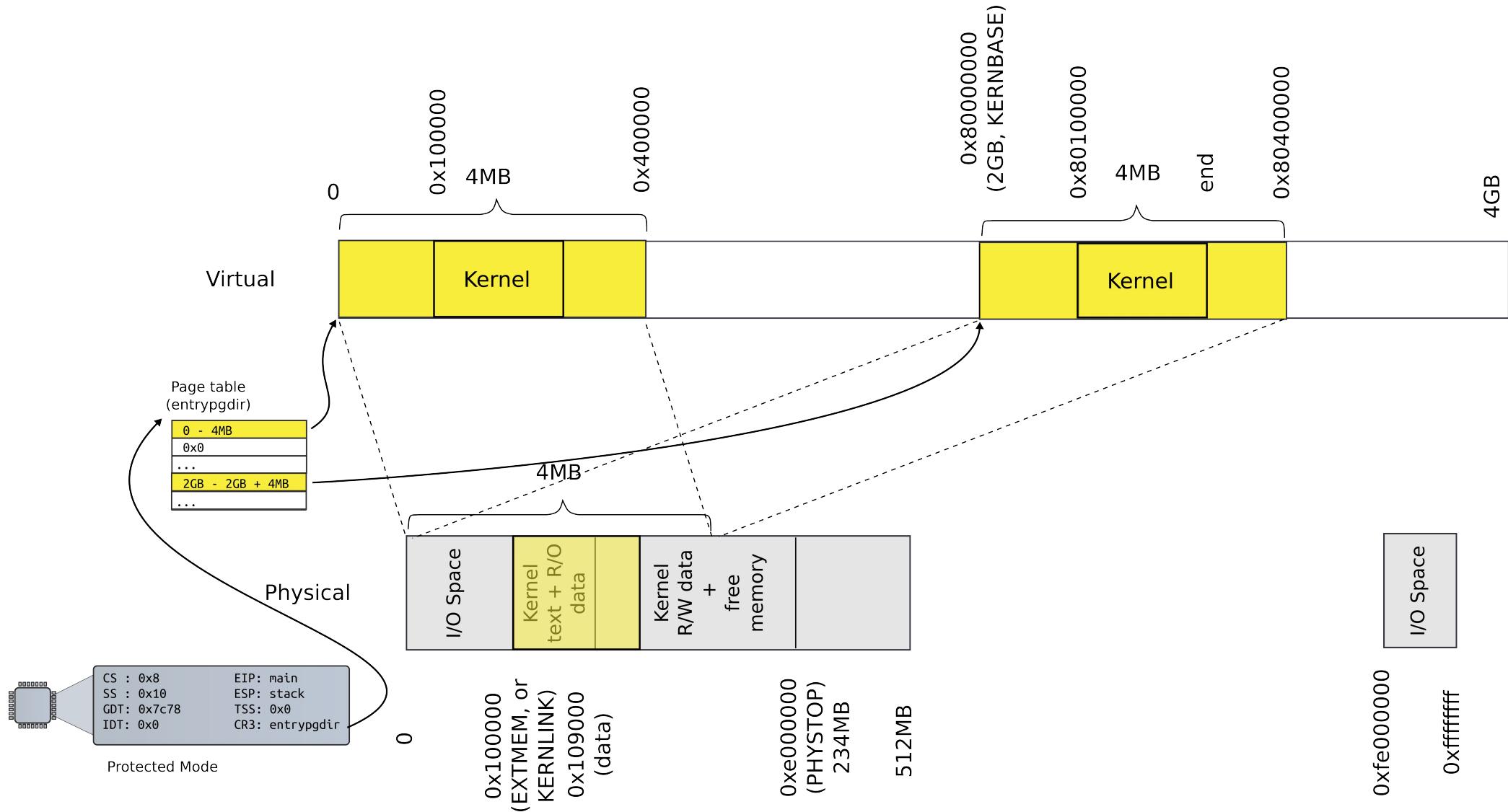




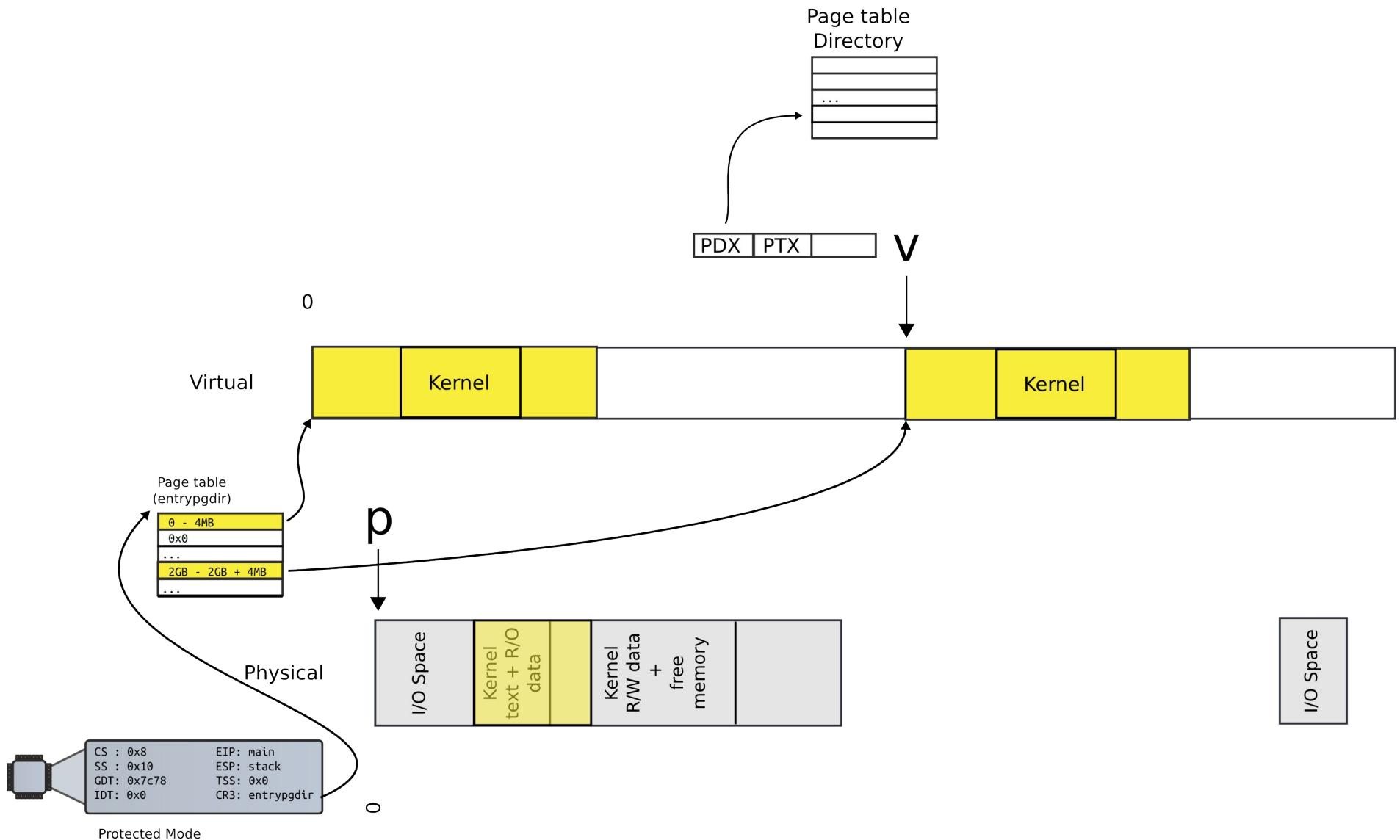
Recap: our goal



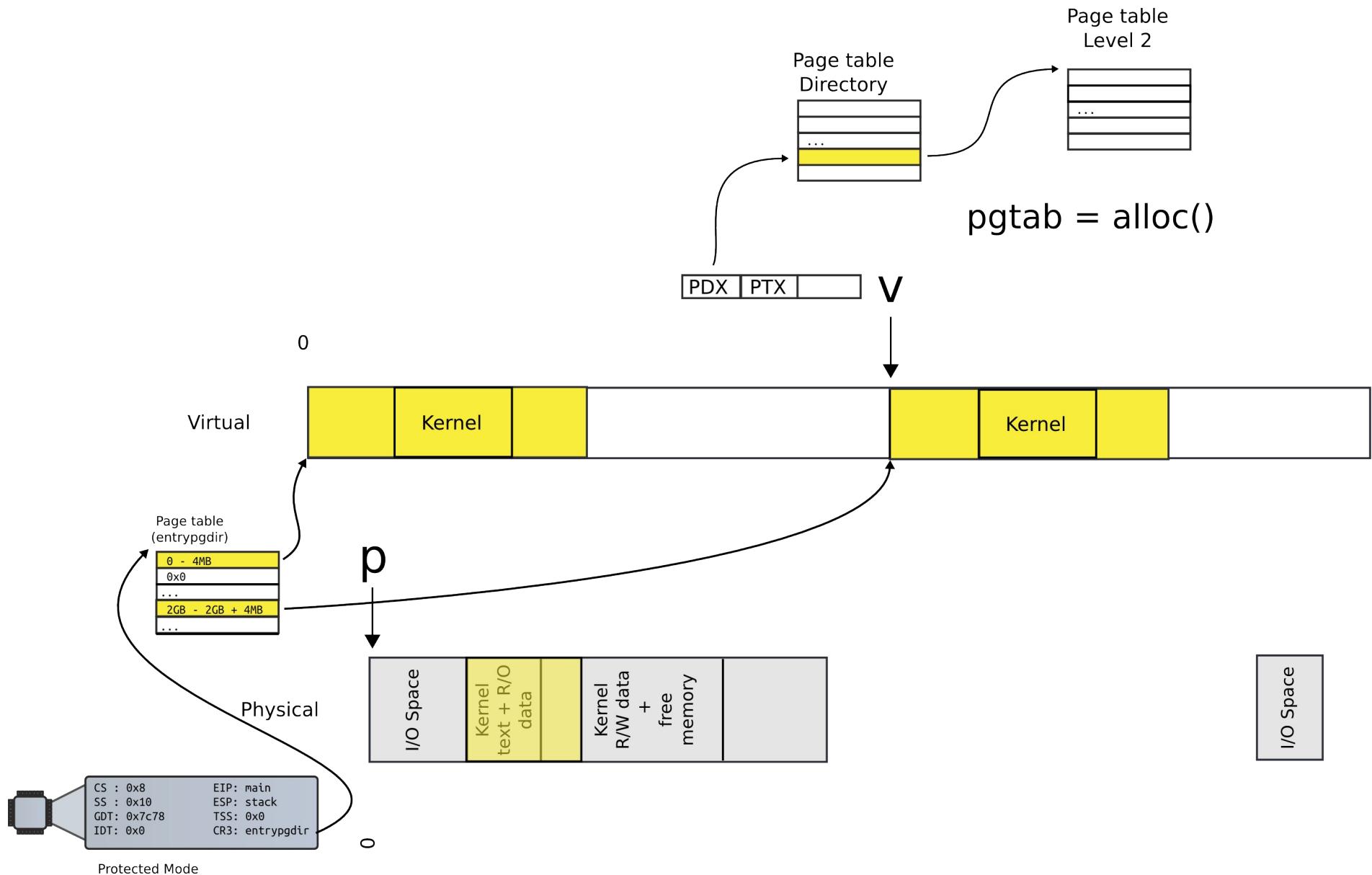
Memory after boot



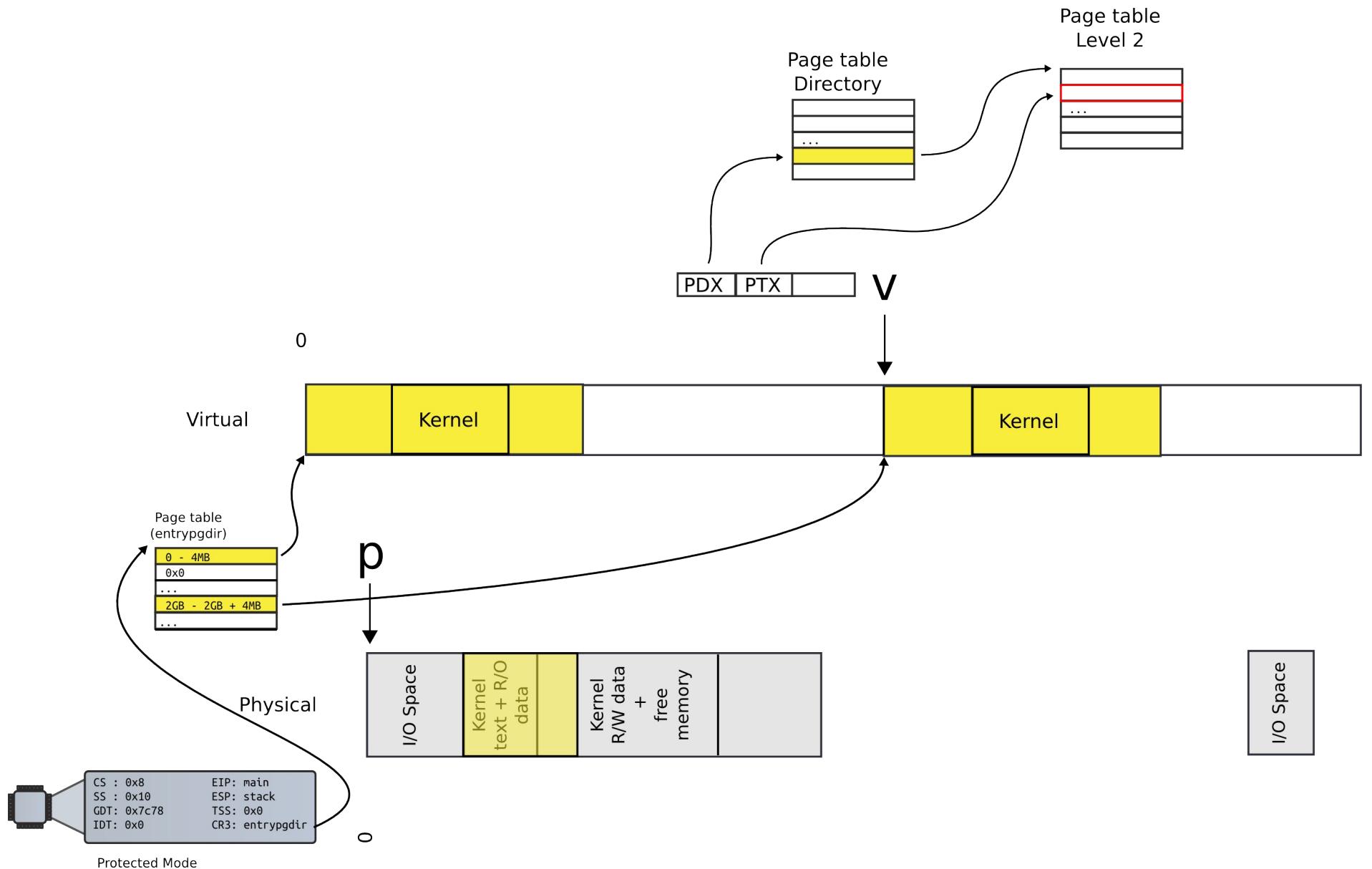
Locate page table directory entry



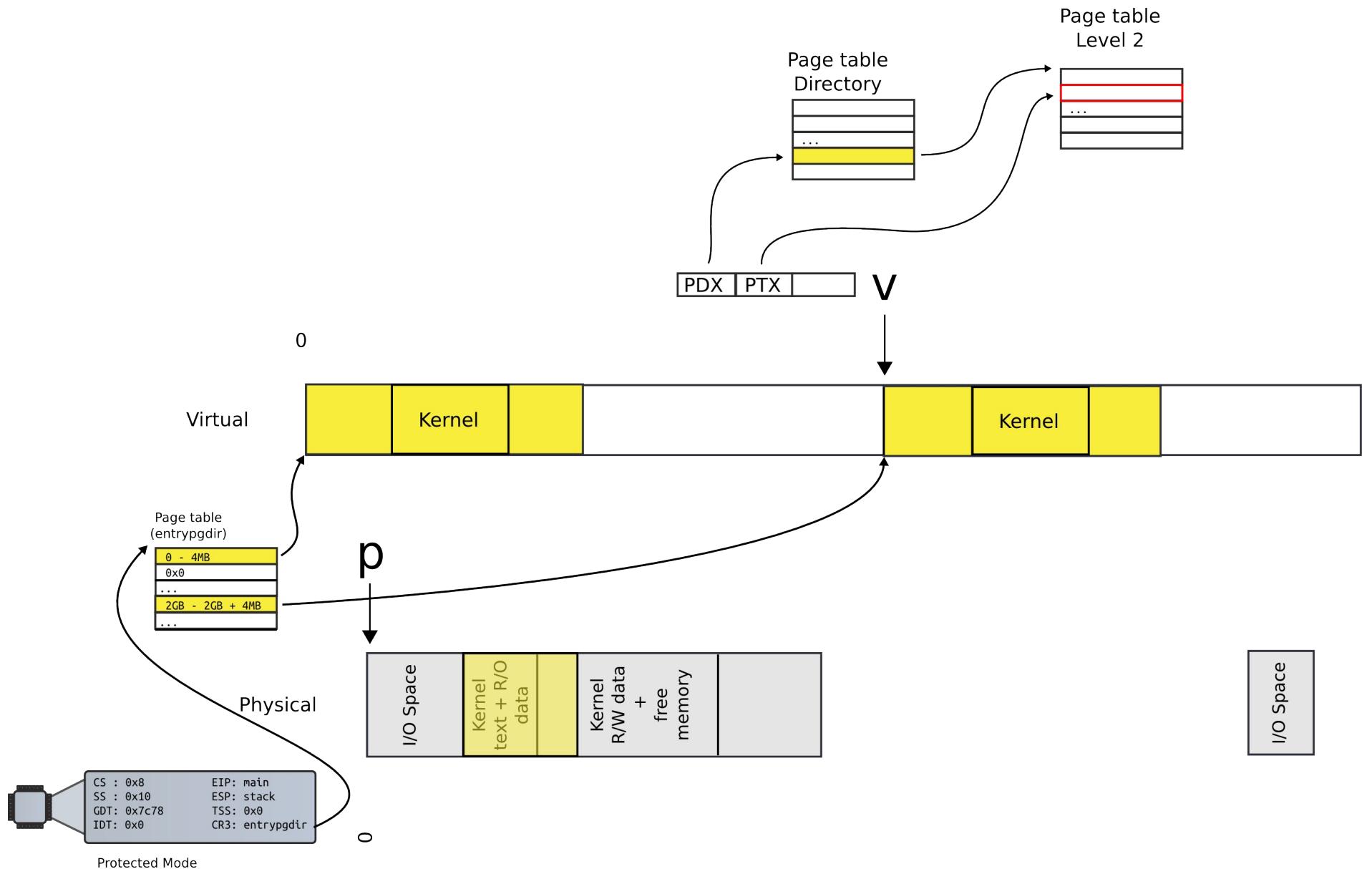
Allocate next level page table



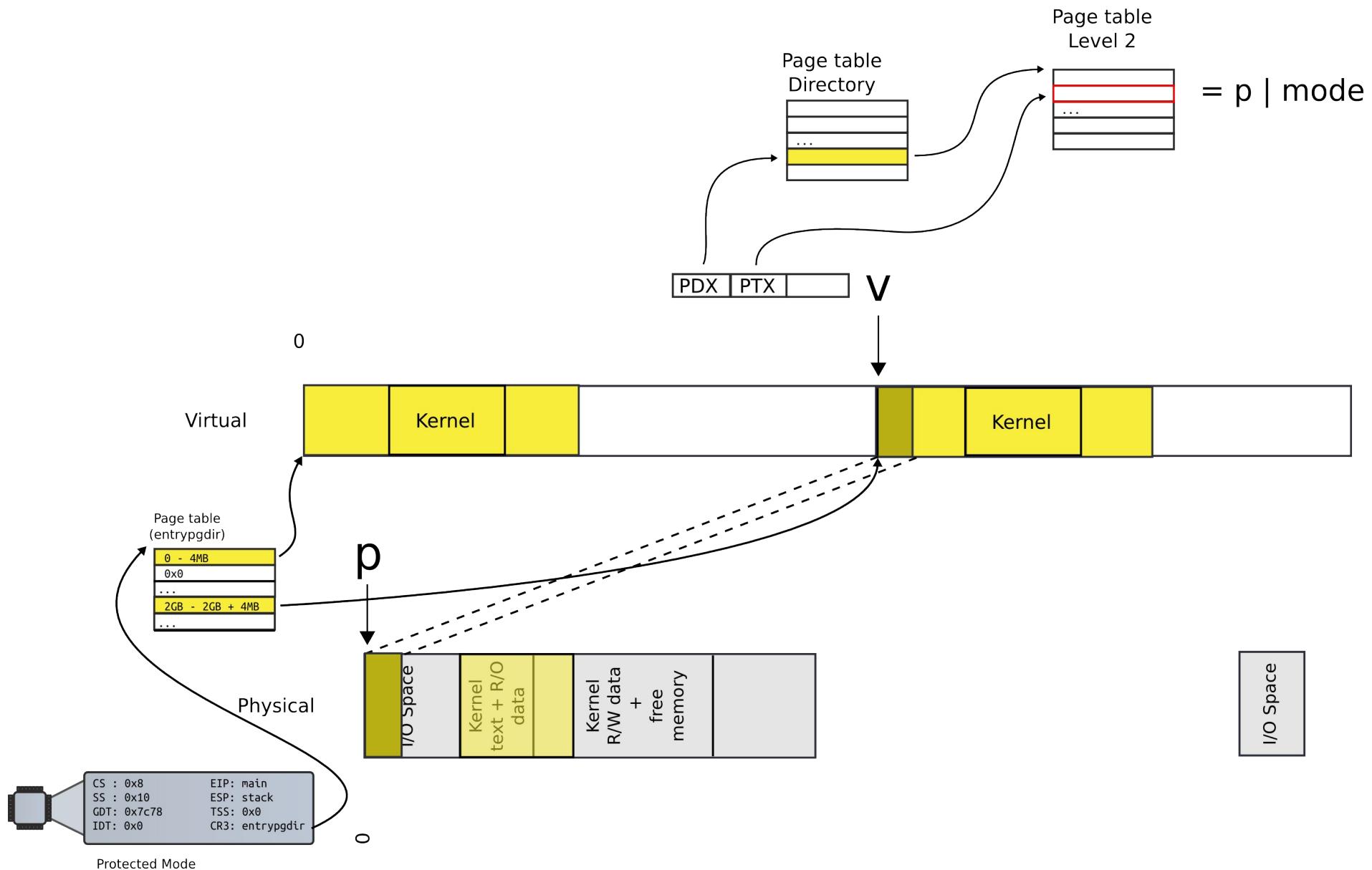
Locate PTE entry



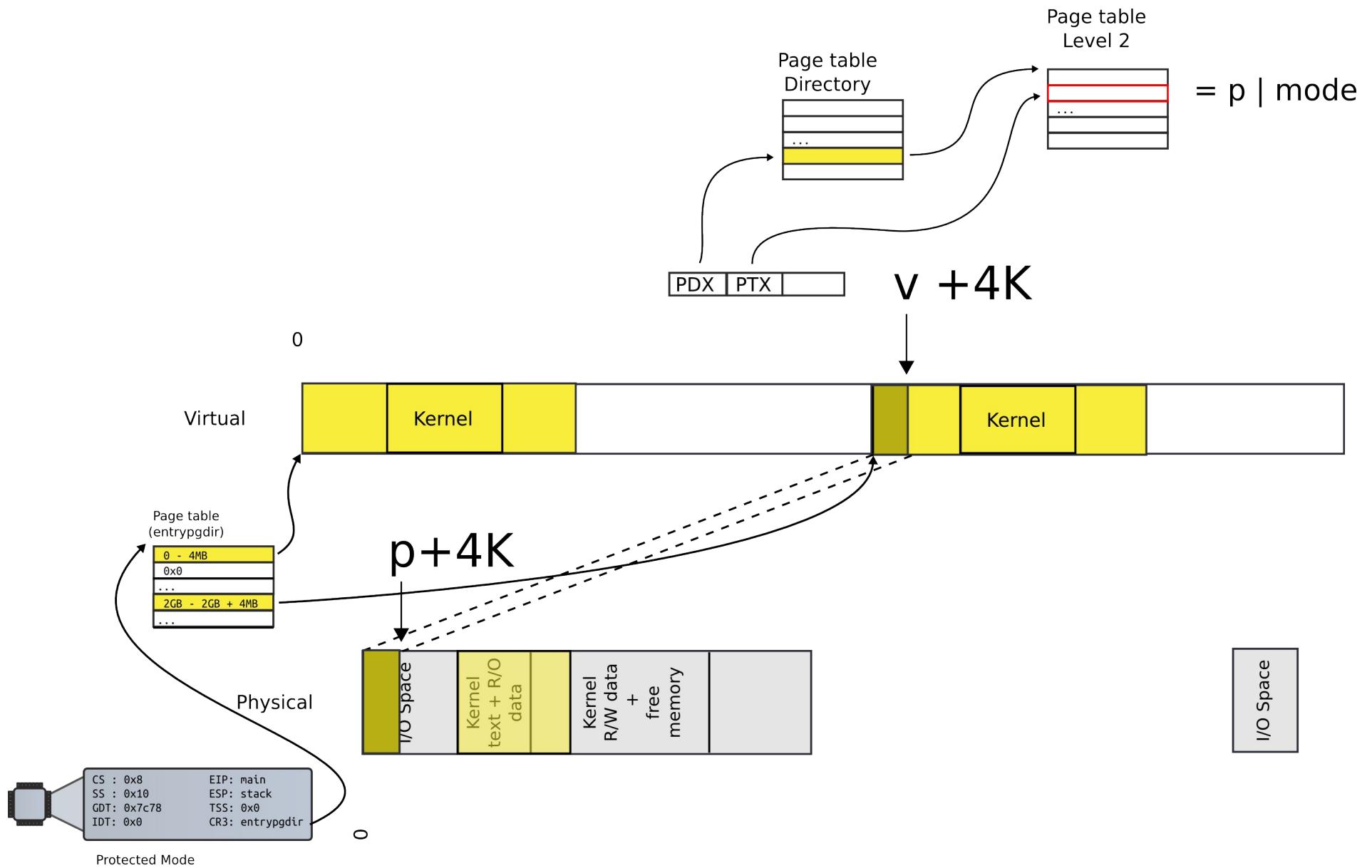
Locate PTE entry



Update mapping with physical addr



Move to next page



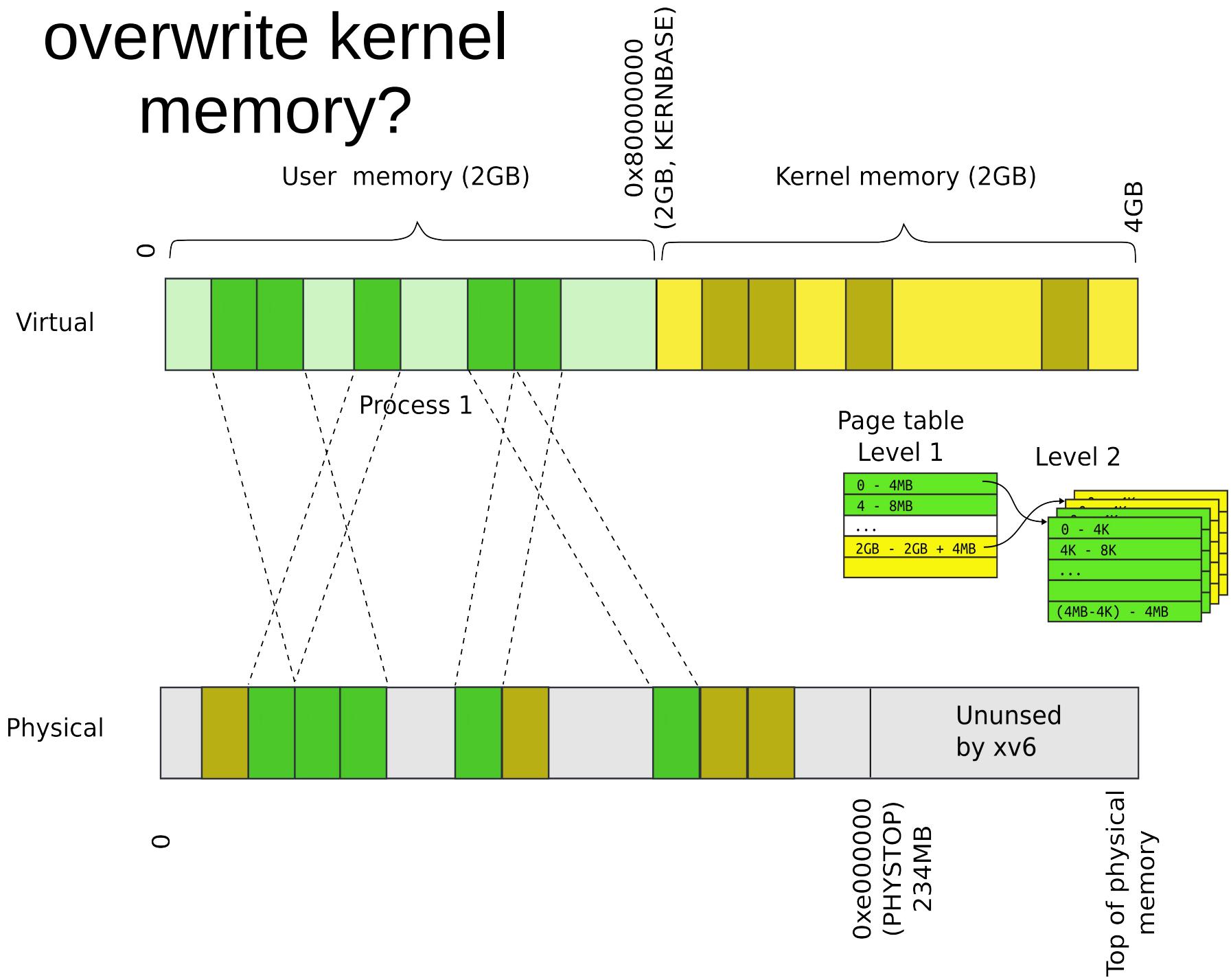
```
1754 walkpgdir(pde_t *pgdir, const void *va, int alloc)
1755 {
1756     pde_t *pde;
1757     pte_t *pgtab;
1758
1759     pde = &pgdir[PDX(va)];
1760     if(*pde & PTE_P){
1761         pgtab = (pte_t*)P2V(PTE_ADDR(*pde));
1762     } else {
1763         if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
1764             return 0;
1765         // Make sure all those PTE_P bits are zero.
1766         memset(pgtab, 0, PGSIZE);
...
1770         *pde = V2P(pgtab) | PTE_P | PTE_W | PTE_U;
1771     }
1772     return &pgtab[PTX(va)];
1773 }
```

Walk page table

```
1779 mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
1780 {
1781     char *a, *last;
1782     pte_t *pte;
1783
1784     a = (char*)PGROUNDDOWN((uint)va);
1785     last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
1786     for(;;){
1787         if((pte = walkpgdir(pgdir, a, 1)) == 0)
1788             return -1;
1789         if(*pte & PTE_P)
1790             panic("remap");
1791         *pte = pa | perm | PTE_P;
1792         if(a == last)
1793             break;
1794         a += PGSIZE;
1795         pa += PGSIZE;
1796     }
1797     return 0;
1798 }
```

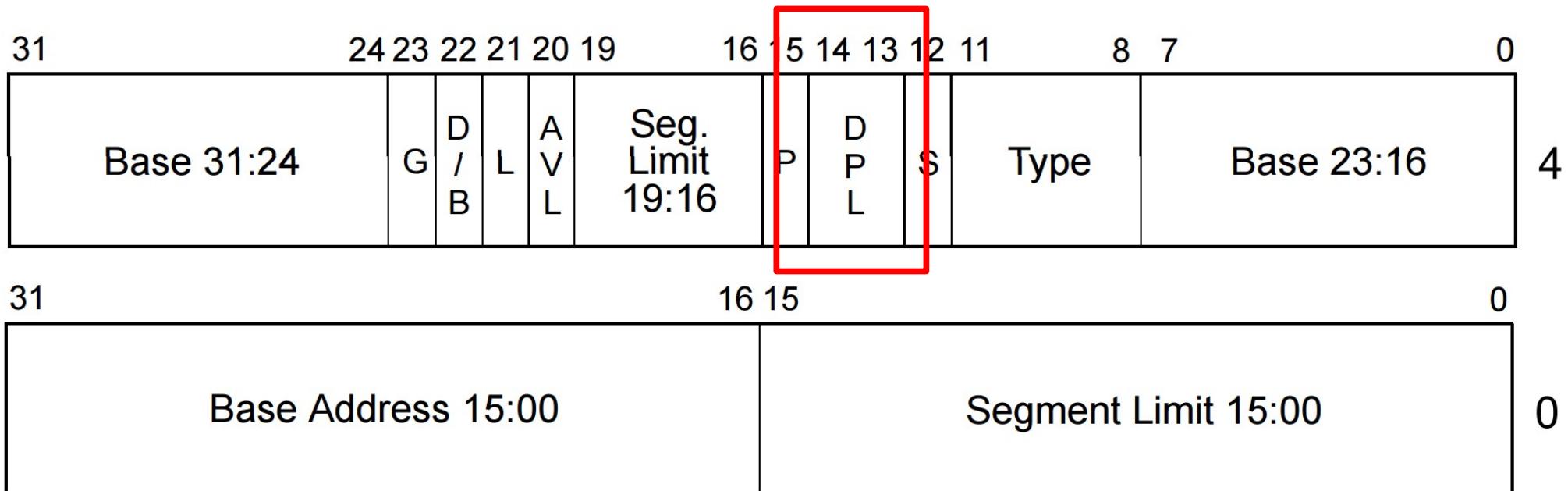
Create page table entries

Can a process overwrite kernel memory?



Privilege levels

- Each segment has a privilege level
 - DPL (descriptor privilege level)
 - 4 privilege levels ranging 0-3

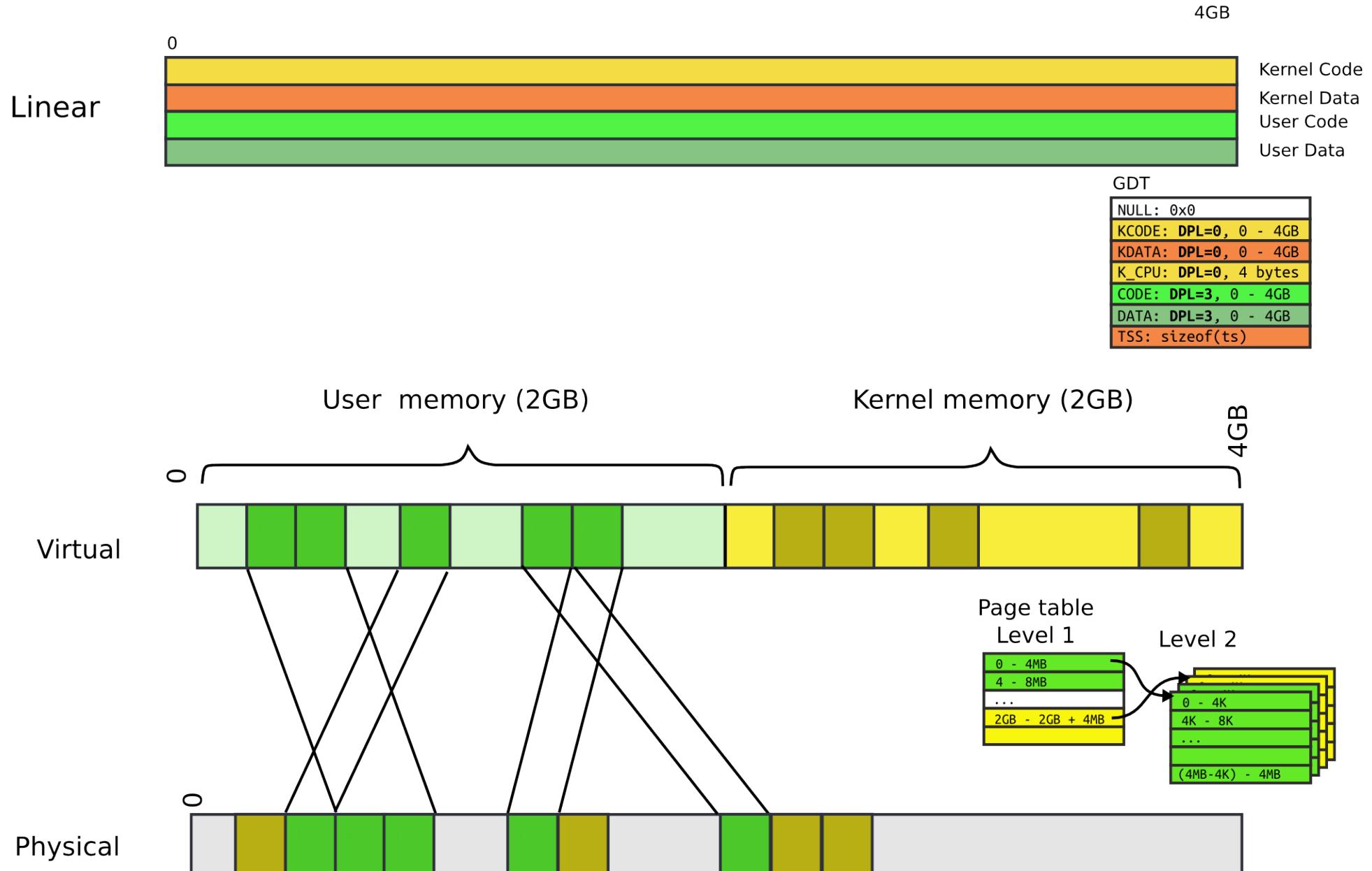


Privilege levels

- Currently running code also has privilege level
 - “Current privilege level” (CPL): 0-3
 - Can access only less privileged segments
 - E.g., 0 can access 1, 2, 3
- Some instructions are “privileged”
 - Can only be invoked at CPL = 0
 - Examples:
 - Load GDT
 - MOV <control register>
 - E.g. reload a page table by changing CR3

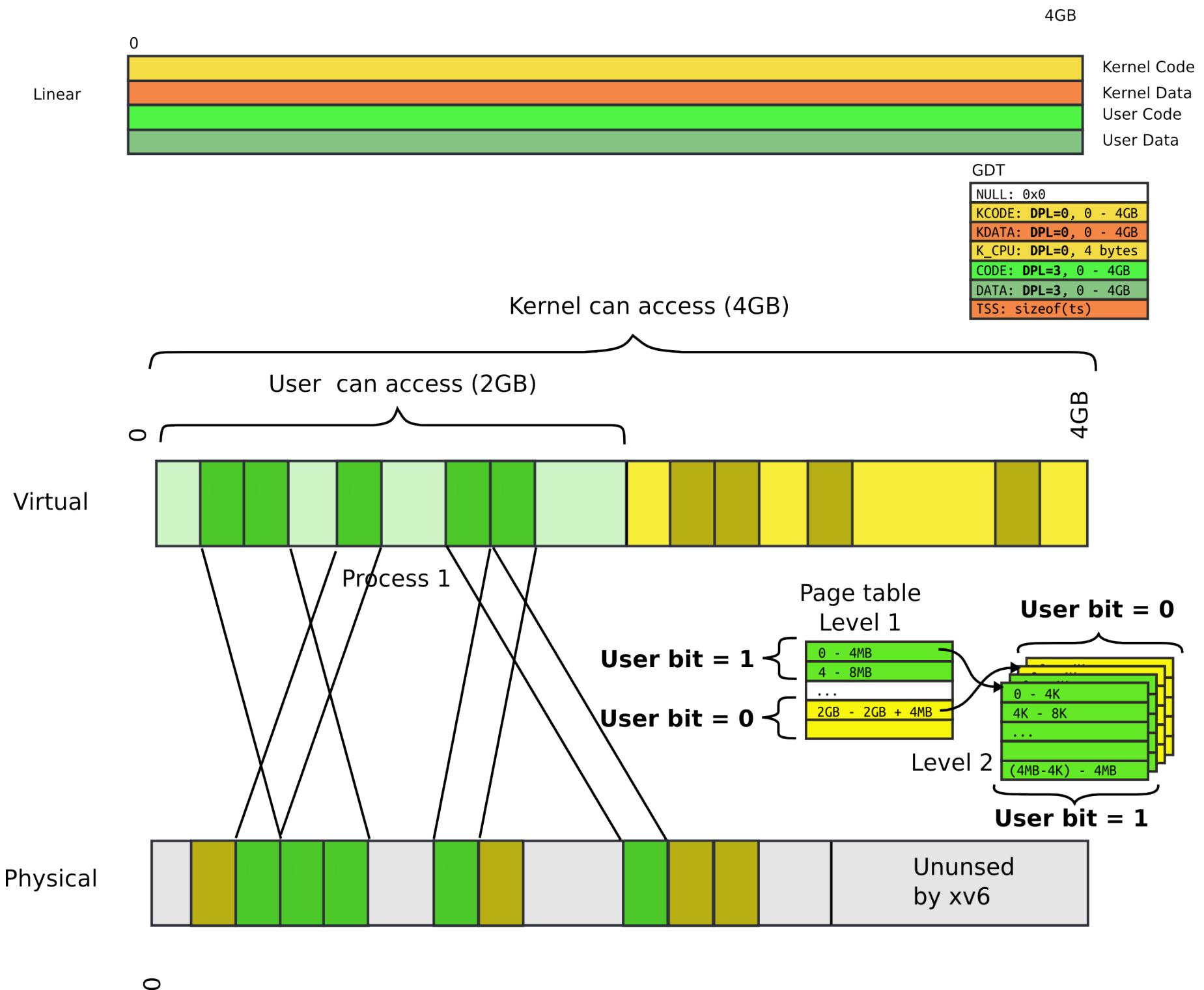
Real world

- Only two privilege levels are used in modern OSes:
 - OS kernel runs at 0
 - User code runs at 3
- This is called “flat” segment model
 - Segments for both 0 and 3 cover entire address space
- But then... how the kernel is protected?
 - Page tables



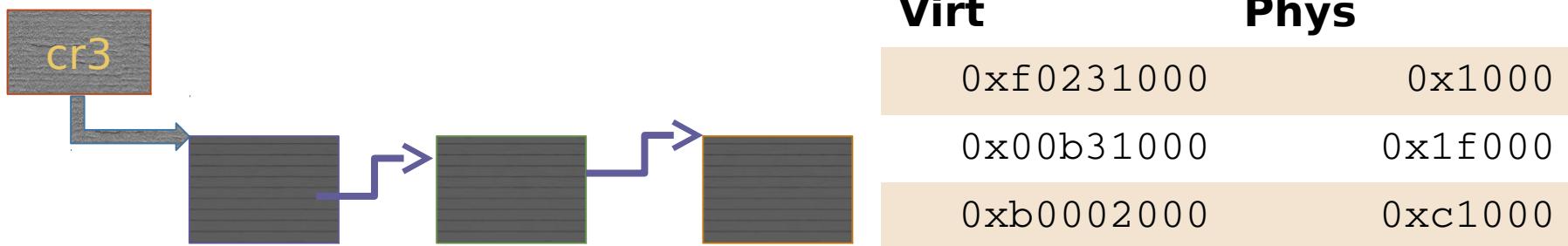
Page table: user bit

- Each entry (both Level 1 and Level 2) has a bit
 - If set, code at privilege level 3 can access
 - If not, only levels 0-2 can access
- Note, only 2 levels, not 4 like with segments
- All kernel code is mapped with the user bit clear
 - This protects user-level code from accessing the kernel



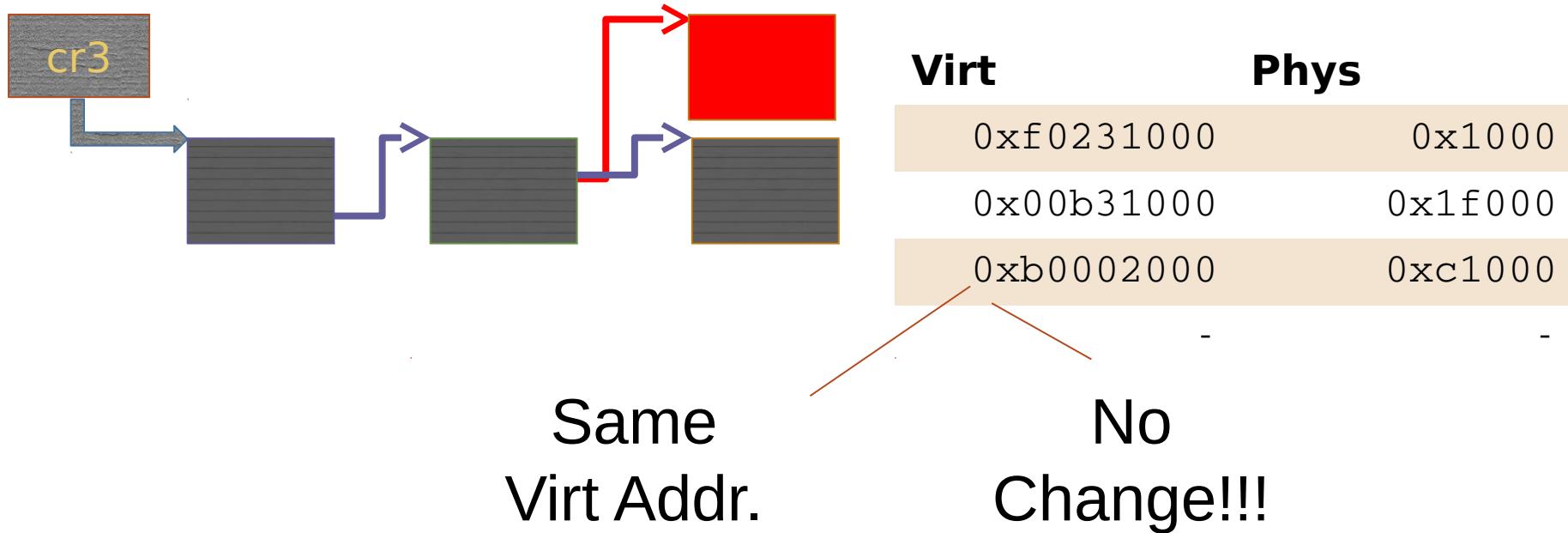
TLB

- CPU caches results of page table walks
 - In translation lookaside buffer (TLB)
- Walking page table is slow
 - Each memory access is 200-300 cycles on modern hardware
 - L3 cache access is 70 cycles



TLB

- TLB is a cache (in CPU)
 - It is not coherent with memory
 - If page table entry is changes, TLB remains the same and is out of sync



Invalidating TLB

- After every page table update, OS needs to manually invalidate cached values
- Modern CPUs have “tagged TLBs”,
 - Each TLB entry has a “tag” – identifier of a process
 - No need to flush TLBs on context switch
- On Intel this mechanism is called
 - Process-Context Identifiers (PCIDs)

Questions?