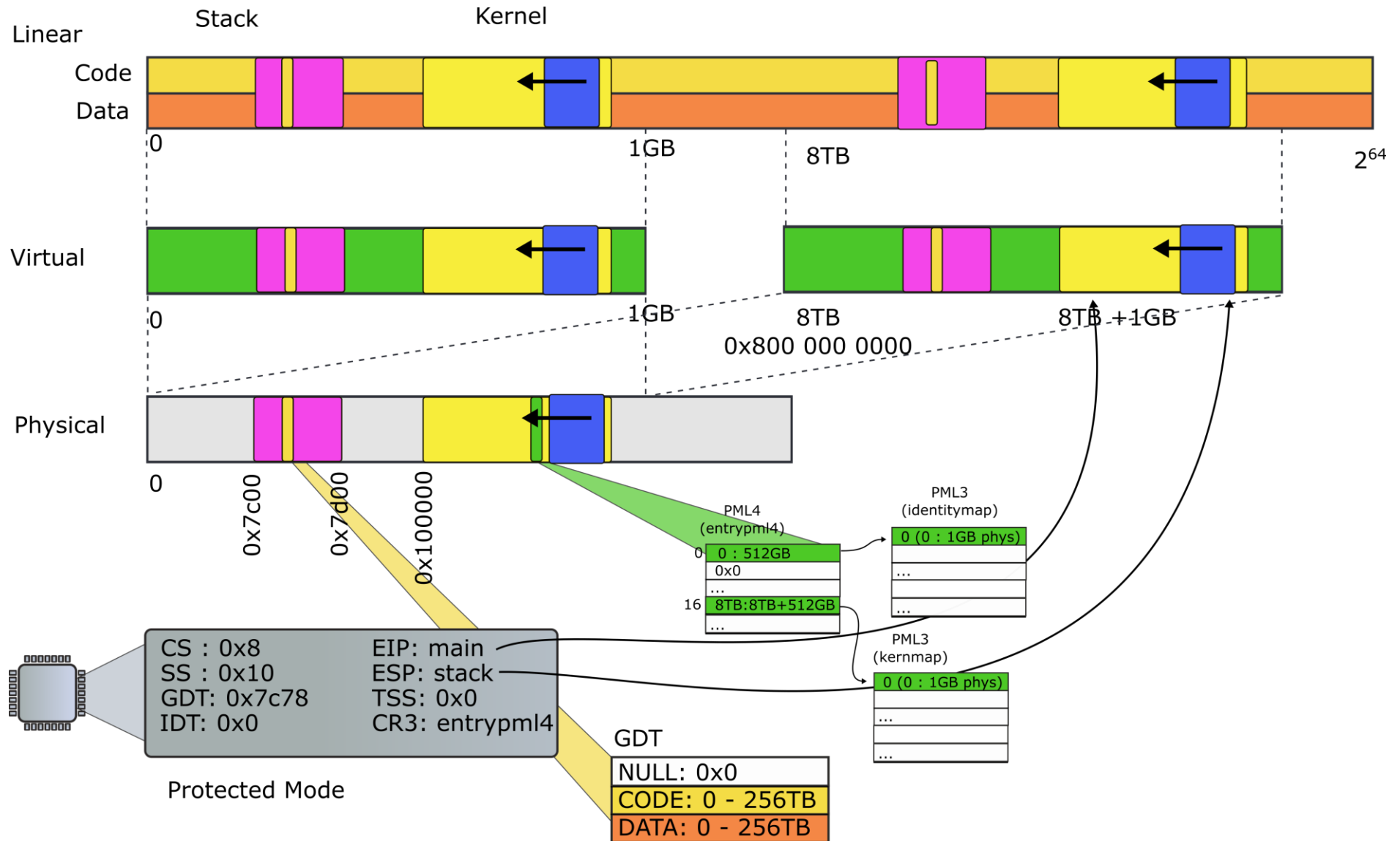# cs5460/6460: Operating Systems

# Lecture 08: System Init
# (Kernel Memory Allocator and Page Table)

Anton Burtsev

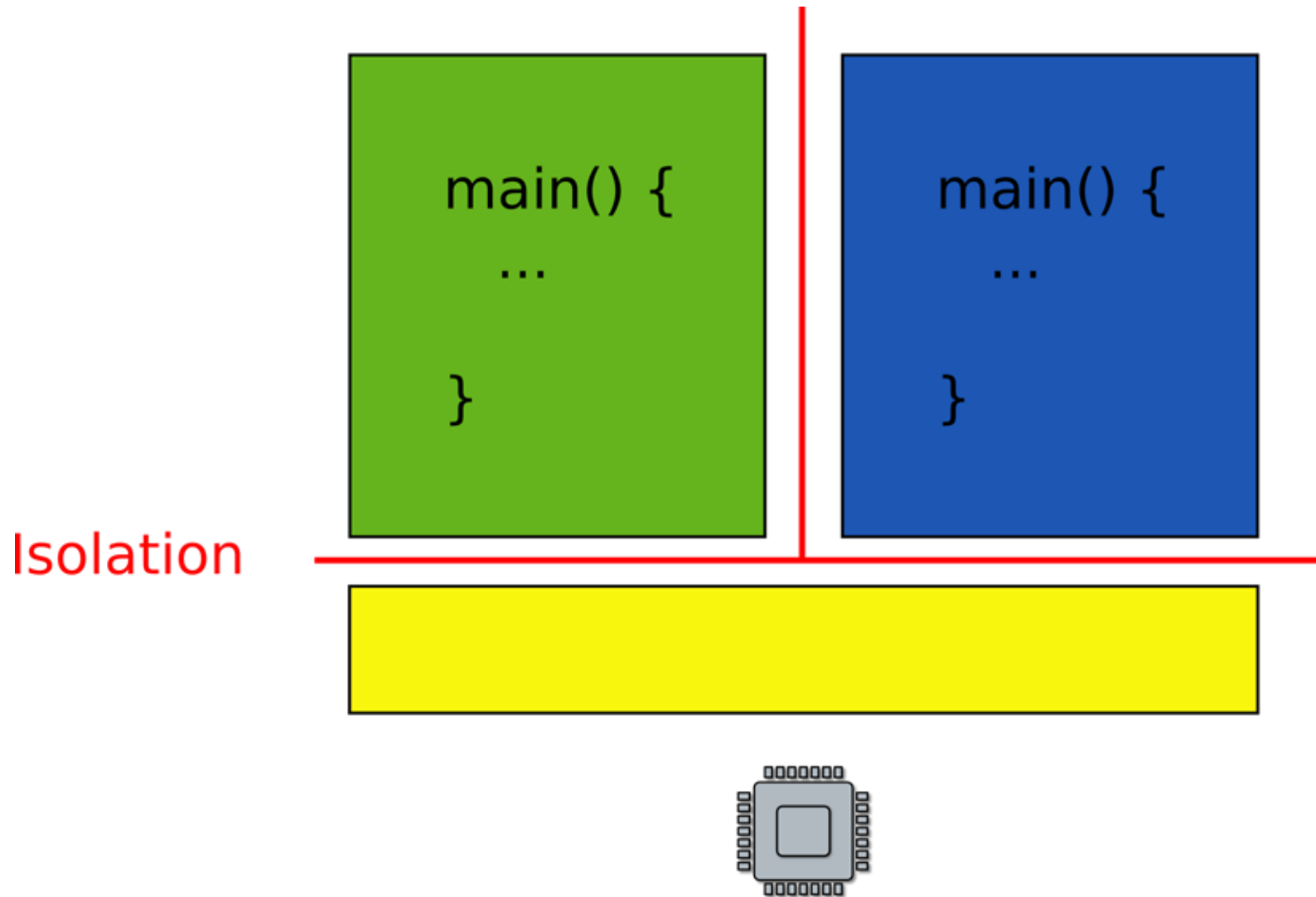February, 2026

# State of the system after boot

# Running in main()

```
1313 // Bootstrap processor starts running C code here.
1314 // Allocate a real stack and switch to it, first
1315 // doing some setup required for memory allocator to work.
1316 int
1317 main(void)
1318 {
1319   kinit1(end, P2V(4*1024*1024)); // phys page allocator
1320   kvmalloc(); // kernel page table
1321   mpinit(); // detect other processors
1322   lapicinit(); // interrupt controller
1323   seginit(); // segment descriptors
1324   cprintf("\ncpu%d: starting xv6\n\n", cpunum());
...
1340 }
```

# What's next?

# We want to run multiple processes

# But what is a process?

# A couple of requirements

- Each process is a collection of resources

- Memory

  - E.g., text, stack, heap

- In-kernel state

  - E.g., open file descriptors, network sockets (connections)
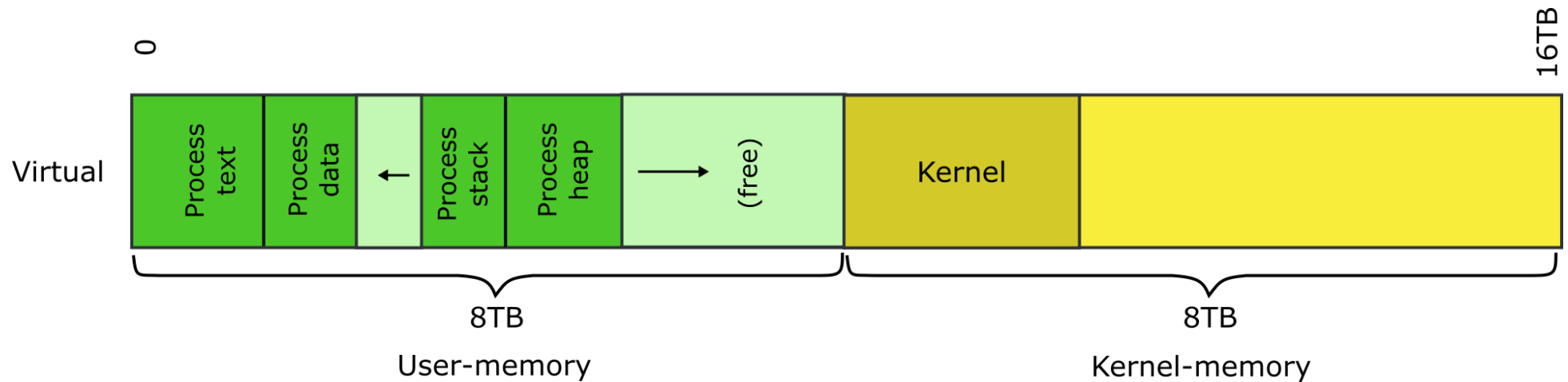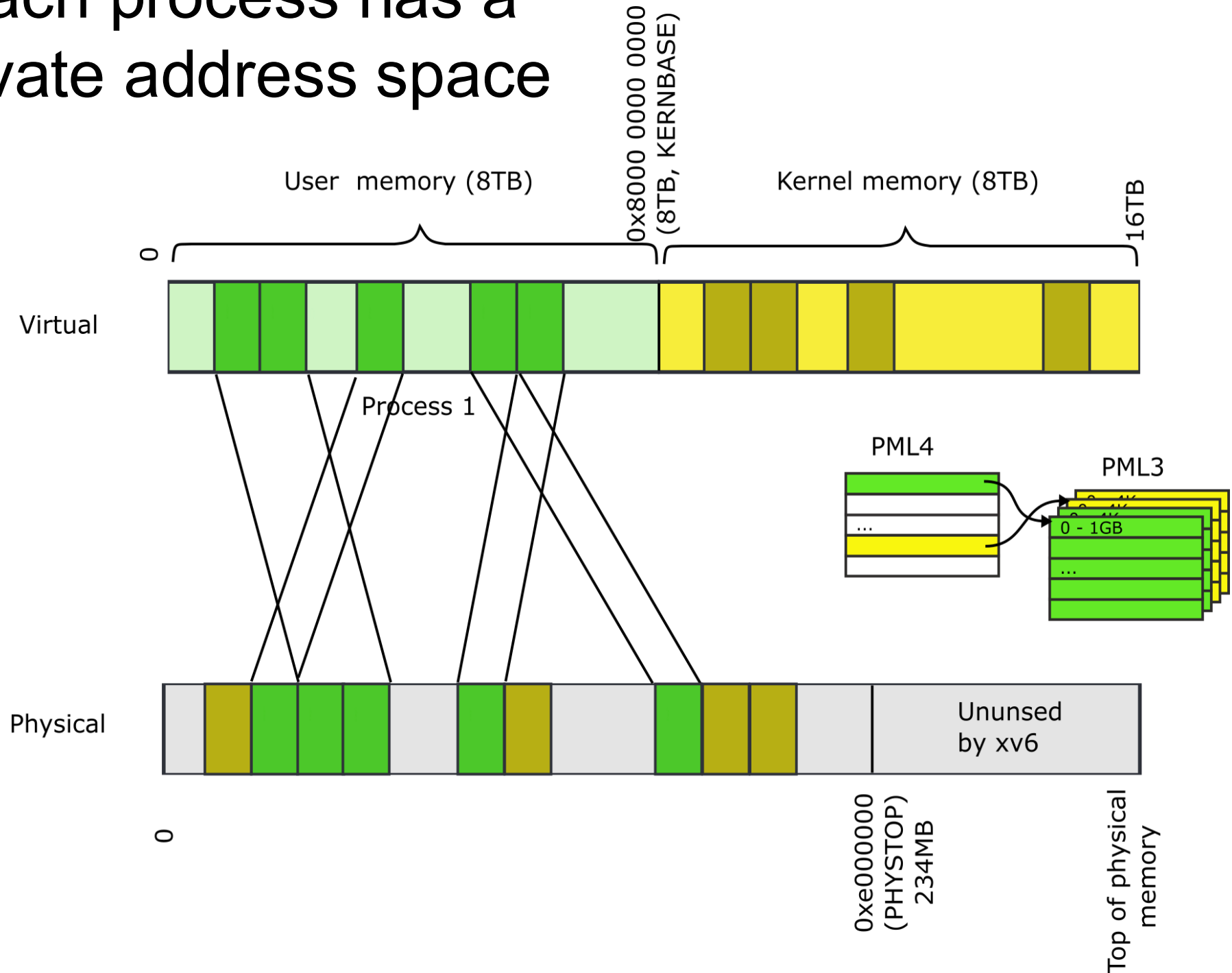
# A couple of requirements

- Each process is a collection of resources
- Memory
  - E.g., text, stack, heap
- In-kernel state
  - E.g., open file descriptors, network sockets (connections)
- Processes are isolated from each other
- Processes don't trust each other
  - Individual users, some privileged
- Can't interfere with other processes
- Can't change kernel (to affect other processes)

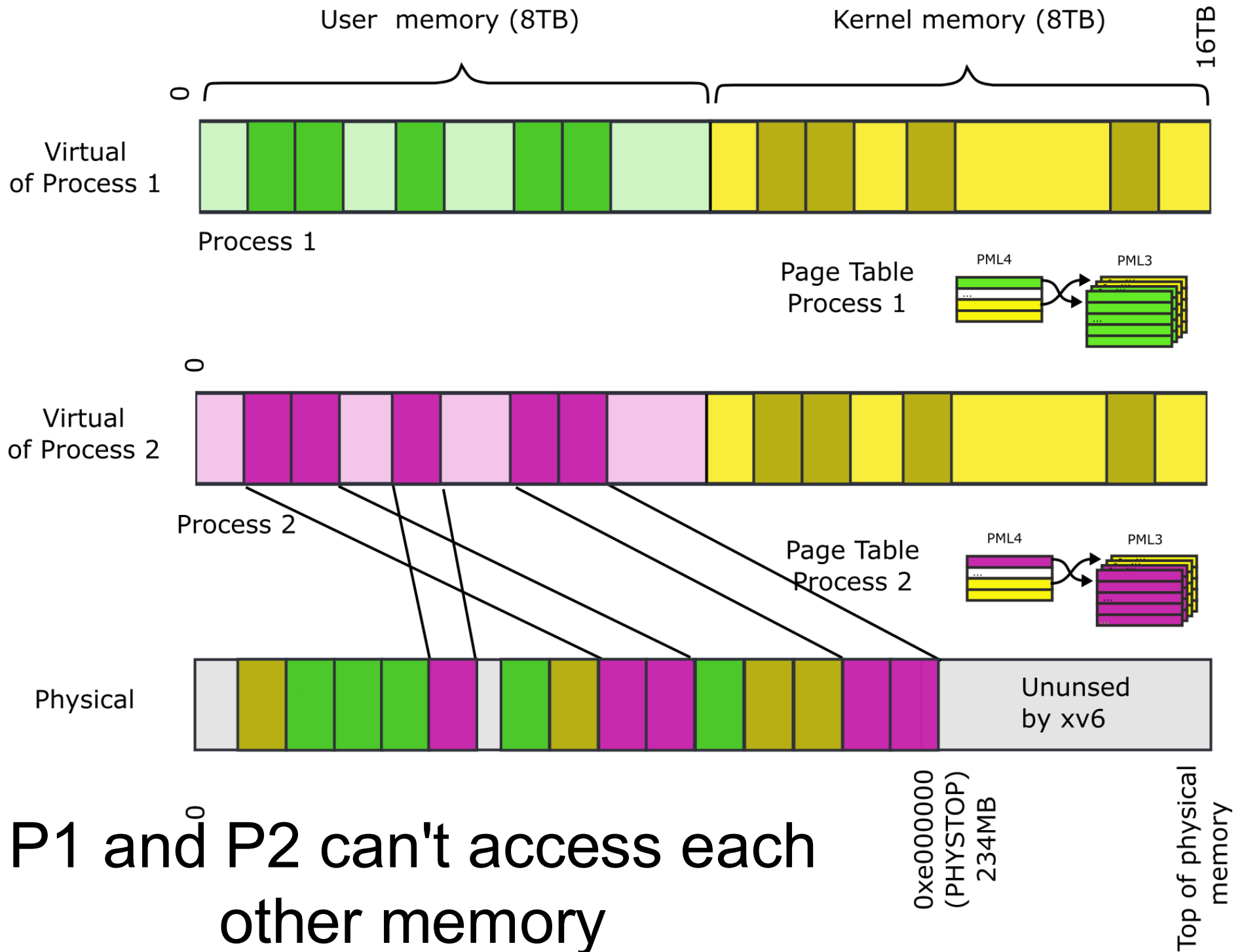# Each process will have 8TB private address space

# Each process has a private address space



User memory (8TB)

0x8000 0000 0000 (8TB, KERNBASE)

Kernel memory (8TB)

16TB

0

Virtual

Process 1

PML4

PML3

... 4K
0 - 1GB
...

Physical

Ununsed by xv6

0

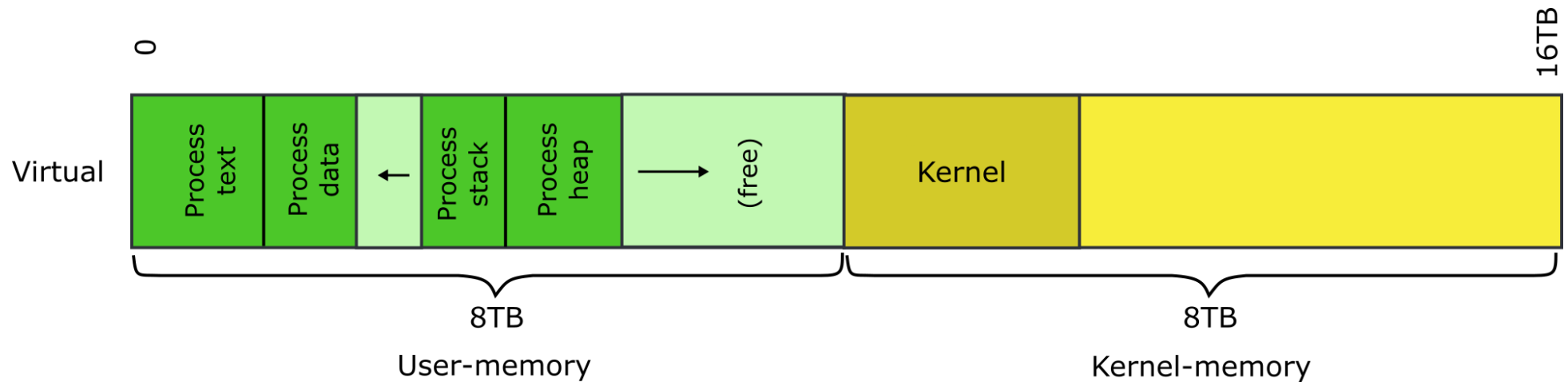0xe000000 (PHYSTOP) 234MB

Top of physical memory
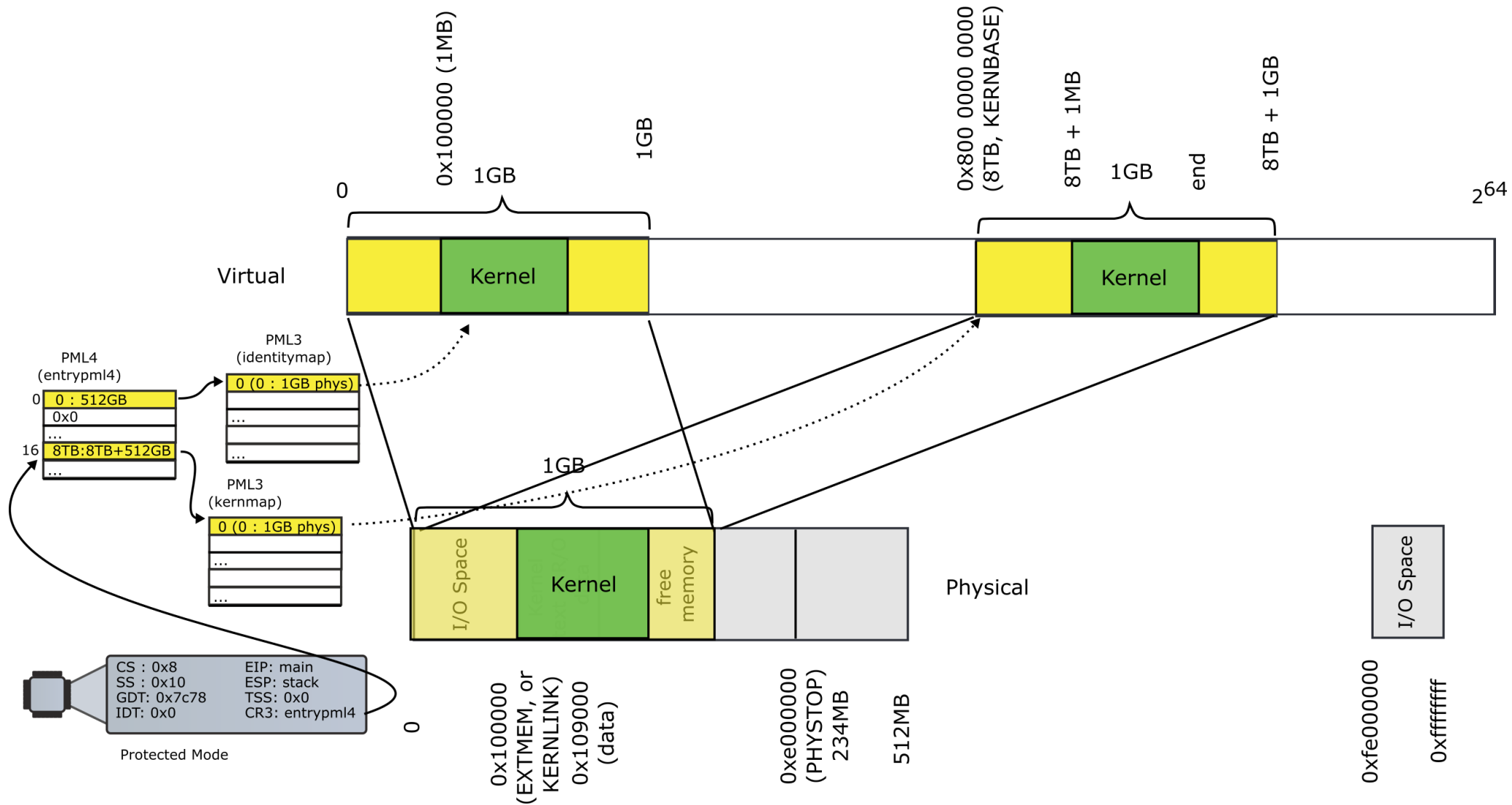
# Each process maps the kernel

- It's not strictly required

- But convenient for system calls

- No need to change the page table when process enters the kernel with a system call

- **Things are much faster!**

User memory (8TB)  Kernel memory (8TB)  16TB

0

Virtual of Process 1

Process 1

Page Table Process 1

PML4  PML3

0

Virtual of Process 2

Process 2

Page Table Process 2

PML4  PML3

Physical

Ununsed by xv6

0

0xe000000 (PHYSTOP) 234MB

Top of physical memory

P1 and P2 can't access each other memory

# Our goal: split address space

# Memory after boot

# Outline

- Create the kernel address space

- Create kernel memory allocator

- Allocate memory for page tables

  - Page table directory and page table

# Kernel memory allocator

- Kernel needs normal 4-level, 4KB page table
- Right now we have:
  - One (statically allocated) page table
  - That has only two 1GB entries
- 4KB page table is a better choice
  - Xv6 processes are small
  - Wasting 1GB or 2MB on a program that fits into 1KB is absurd

- But to create page tables we need memory
- Where can it come from?

# Simple memory allocator

- Goal:
  - alloc() and free()
  - To allocate page tables, stacks, data structures, etc.

## Question 6

**Points**

10

✕ Delete Question

Title
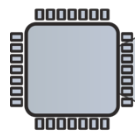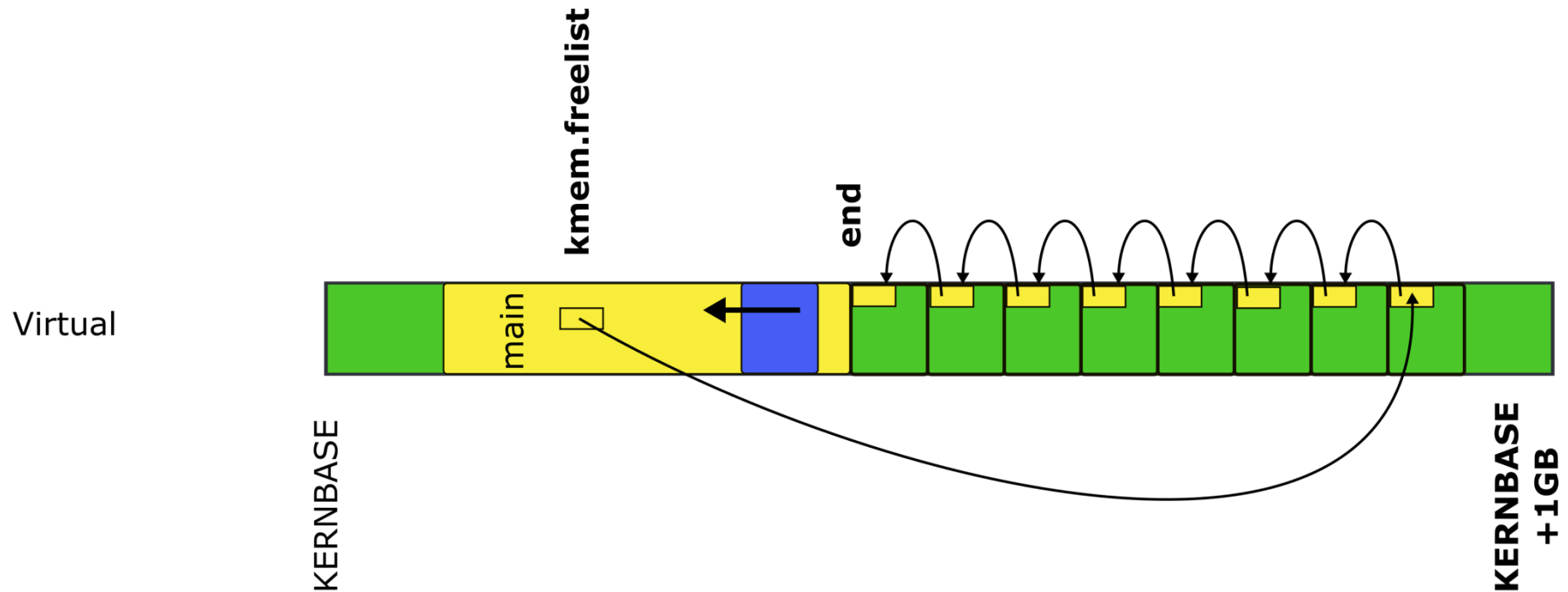
**Problem**          🖼 Insert Image          ≣ Insert Field

What is the virtual address of the very first page added to the kernel memory allocator, considering that the kernel memory allocator has just been initialized? Assume that the kernel is linked such that the "end" symbol is at virtual address 0x8000066beef.

( ) 0x800006steak
( ) 0x80000660000
( ) 0x8000066b000
( ) 0x8000066beef
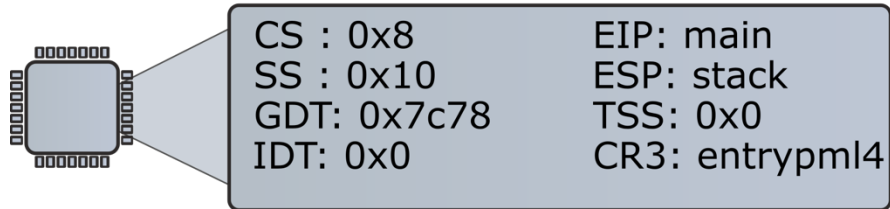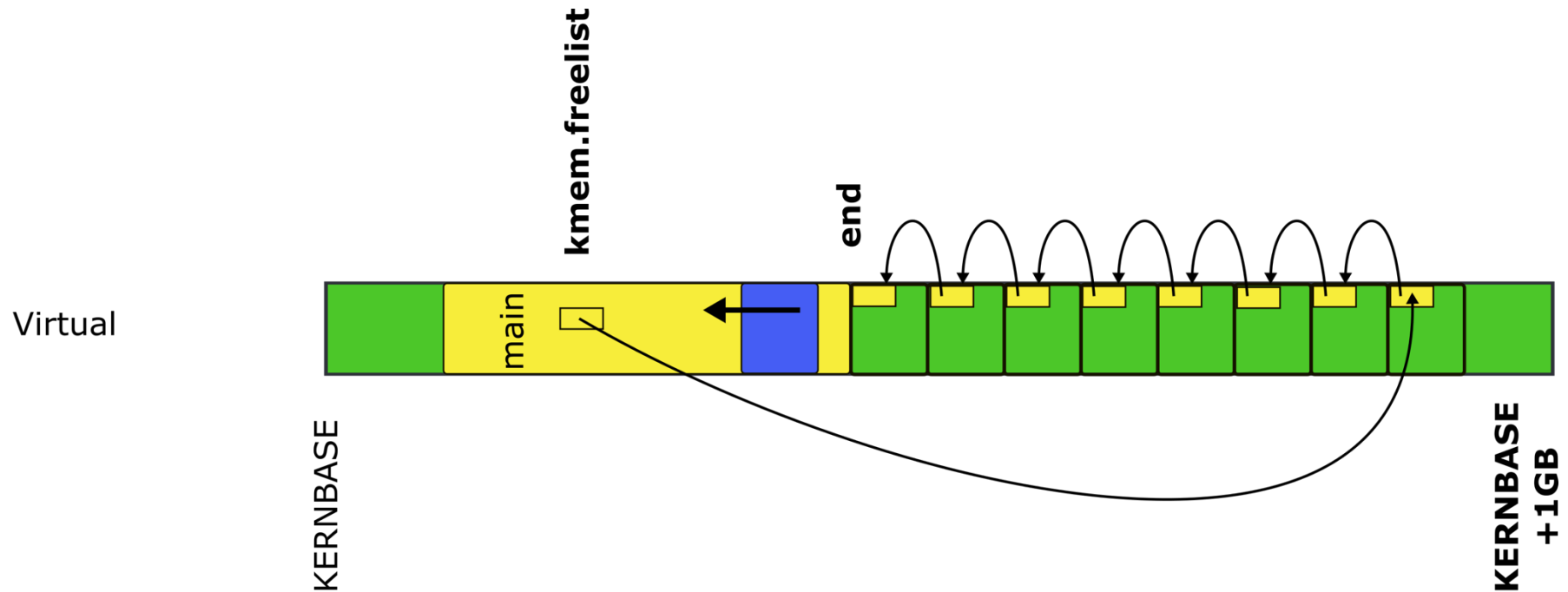( ) 0x8000066pork
( ) 0x8000066c000

# Page allocator



Virtual

kmem.freelist

end

main

KERNBASE

KERNBASE +1GB

CS : 0x8          EIP: main
SS : 0x10         ESP: stack
GDT: 0x7c78       TSS: 0x0
IDT: 0x0          CR3: entrypml4

Protected Mode

# Page allocator



Virtual

kmem.freelist

end

main

KERNBASE

KERNBASE +1GB

CS : 0x8    EIP: main
SS : 0x10   ESP: stack
GDT: 0x7c78 TSS: 0x0
IDT: 0x0    CR3: entrypml4
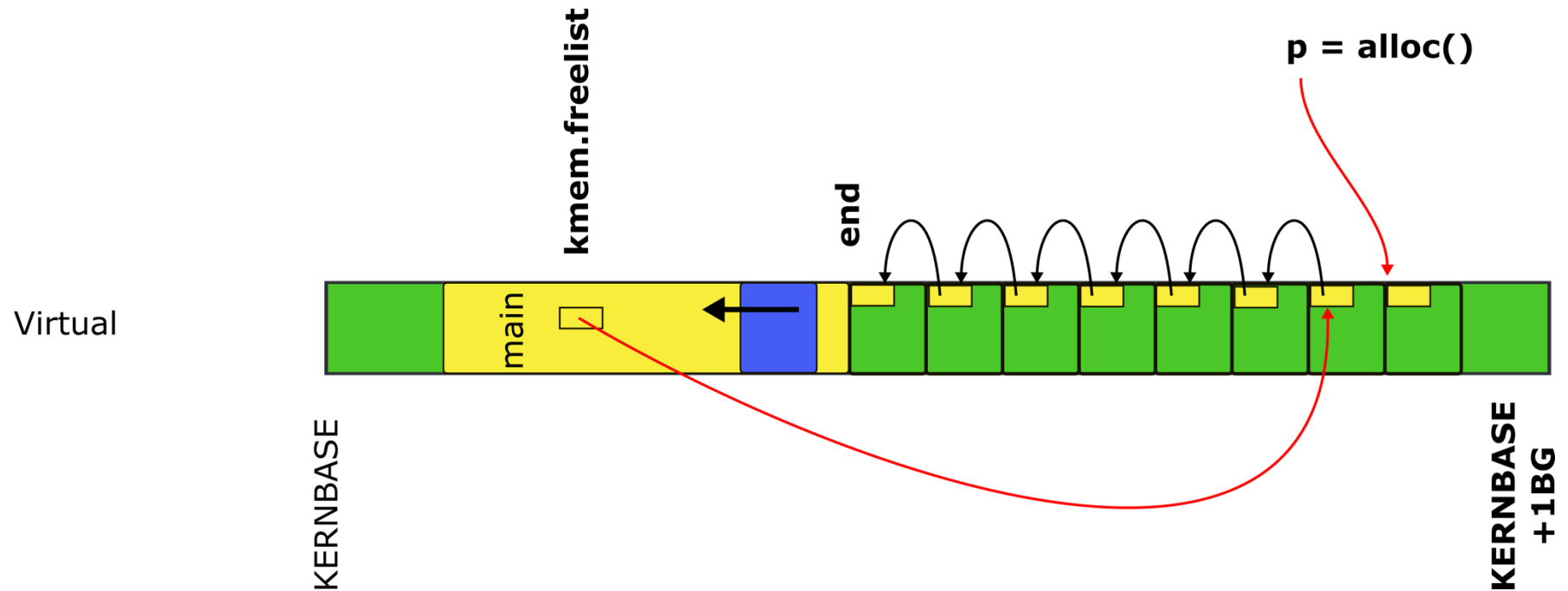
Protected Mode

3014 struct run {
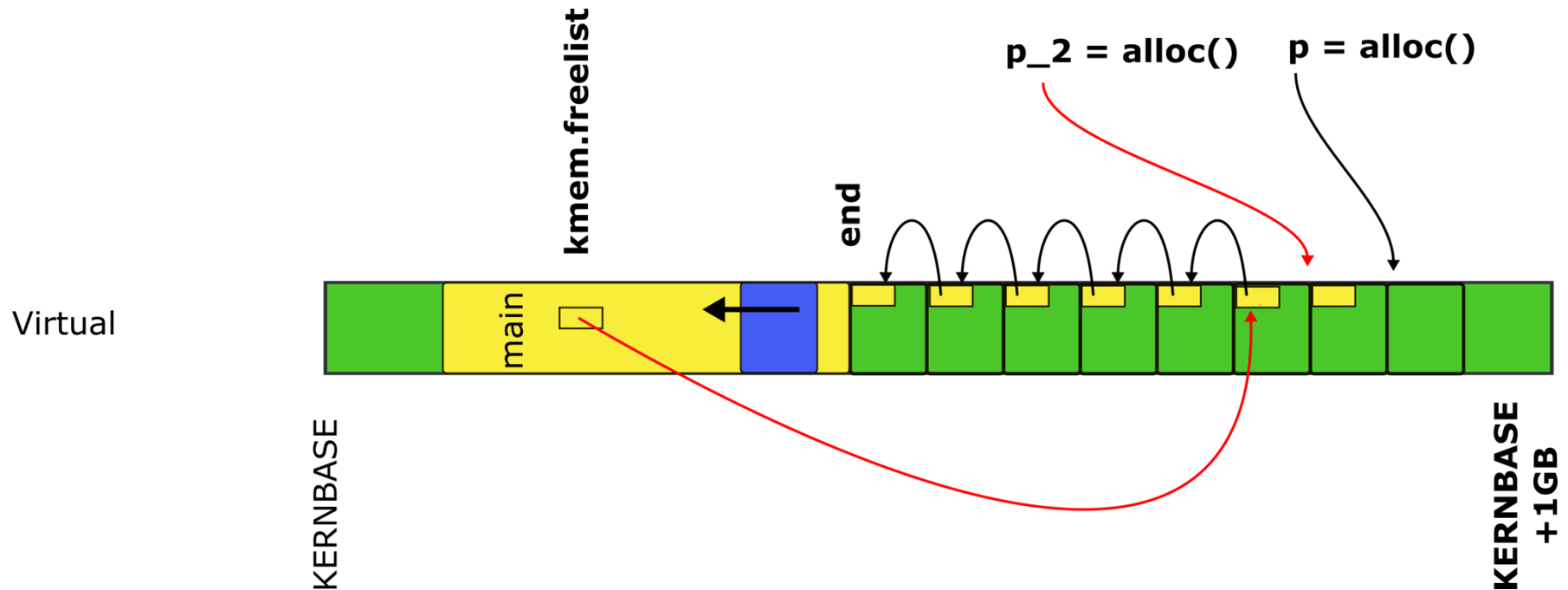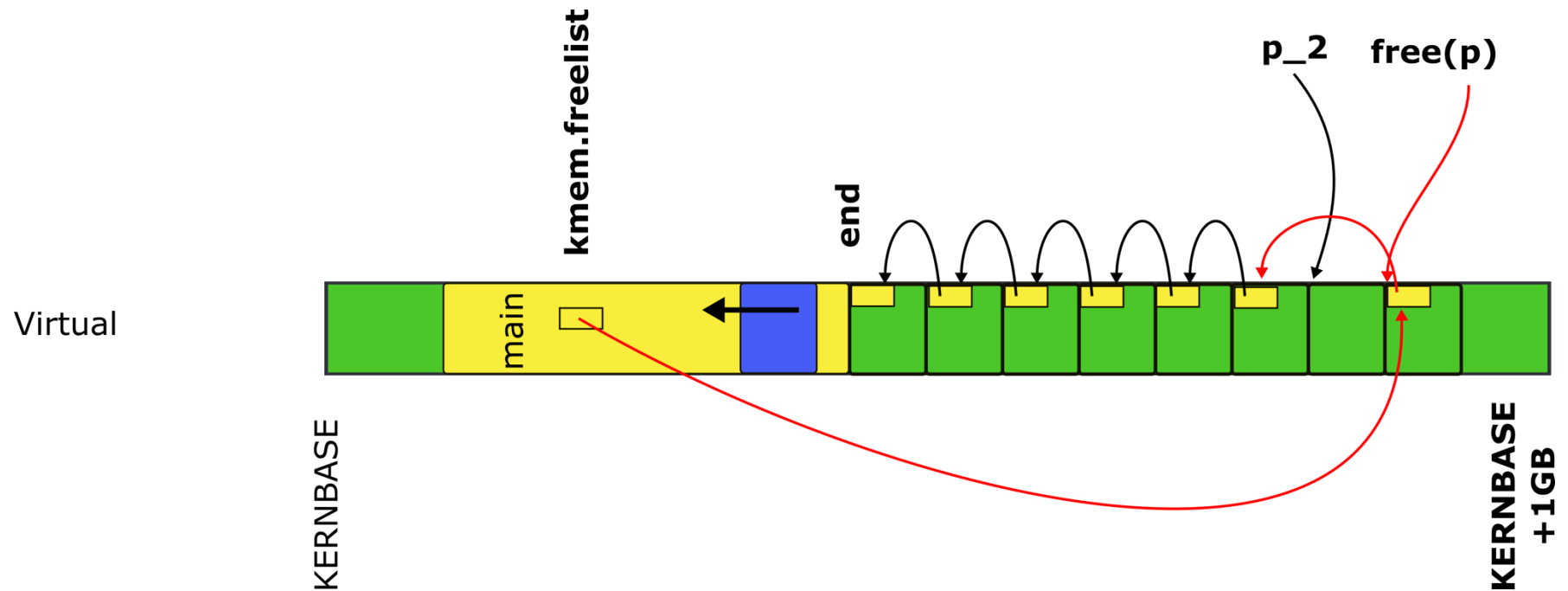
3015   struct run *next;

3016 };

# Page allocator



3014 struct run {

3015    struct run *next;

3016 };

# Page allocator



3014 struct run {

3015   struct run *next;

3016 };

# Page allocator



3014 struct run {

3015   struct run *next;

3016 };

# kalloc() - kernel allocator

3087 char*

3088 kalloc(void)

3089 {

3080   struct run *r;

...

3094   r = kmem.freelist;

3095   if(r)

3096     kmem.freelist = r−>next;

...

3099   return (char*)r;

3099 }



kmem.freelist

p = alloc()

main

3065 kfree(char *v)

3066 {

3067   struct run *r;

...

3077   r = (struct run*)v;

3078   r−>next = kmem.freelist;

3079   kmem.freelist = r;

...

2832 }

kmem.freelist

- Where can we get memory to keep the list itself?

- Note, the list is maintained within each page

- It has to write each page though to update the "next" pointer

# There is free memory in the 1GB page we've mapped

# Donate this free memory to the allocator



Virtual

kmem.freelist

end

main

KERNBASE

KERNBASE +1GB

CS : 0x8        EIP: main
SS : 0x10      ESP: stack
GDT: 0x7c78   TSS: 0x0
IDT: 0x0        CR3: entrypml4

Protected Mode

- Take memory from the end of the kernel binary
- To the end of the 4MB page

# kinit1(): initialize the allocator with free memory

```
1366 int
1367 main(void)
1368 {
1369   void *kinit1_end = (1024*1024*1024 > PHYSTOP) ? P2V(PHYSTOP)
1370                           : P2V(1024*1024*1024);
1371   kinit1(end, kinit1_end); // phys page allocator
1372.  kvmalloc(); // kernel page table
1373   mpinit(); // detect other processors
1374   lapicinit(); // interrupt controller
1375   seginit(); // segment descriptors
1376   picinit(); // disable pic
1377   ioapicinit(); // another interrupt controller

   ...
```

# Freerange()

3030 kinit1(void *vstart, void *vend)

3031 {

...

3034  **freerange(vstart, vend);**

3035 }

- Free range of memory from vstart to vend giving it to the allocator
- i.e., adding pages to the list

# freerange()

3051 freerange(void *vstart, void *vend)

3052 {

3053   char *p;

3054   p = (char*)PGROUNDUP((uint)vstart);

3055   for(; p + PGSIZE <= (char*)vend; p += PGSIZE)

3056     **kfree(p);**

3057 }

- freerange() internally simply frees the pages from vstart to vend
- kfree() adds them to the allocator list

# Where do we start?

1366 int

1367 main(void)

1368 {

1369   void *kinit1_end = (1024*1024*1024 > PHYSTOP) ? P2V(PHYSTOP)

1370                              : P2V(1024*1024*1024);

1371   kinit1(end, kinit1_end); // phys page allocator

- ## What is this **end**?

1311 extern char end[];

# Where do we start?

1366 int

1367 main(void)

1368 {

1369   void *kinit1_end = (1024*1024*1024 > PHYSTOP) ? P2V(PHYSTOP)

1370                           : P2V(1024*1024*1024);

1371   kinit1(end, kinit1_end); // phys page allocator

- ## What is this **end**?

1311 extern char end[]; // first address after

          kernel loaded from ELF file

# Donate this free memory to the allocator

# Recap

- Kernel has a memory allocator

- It allocates memory in chunks of 4KB

- Good enough to maintain kernel data structures

# Kernel page table
# (for 4KB page tables)

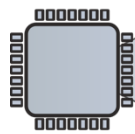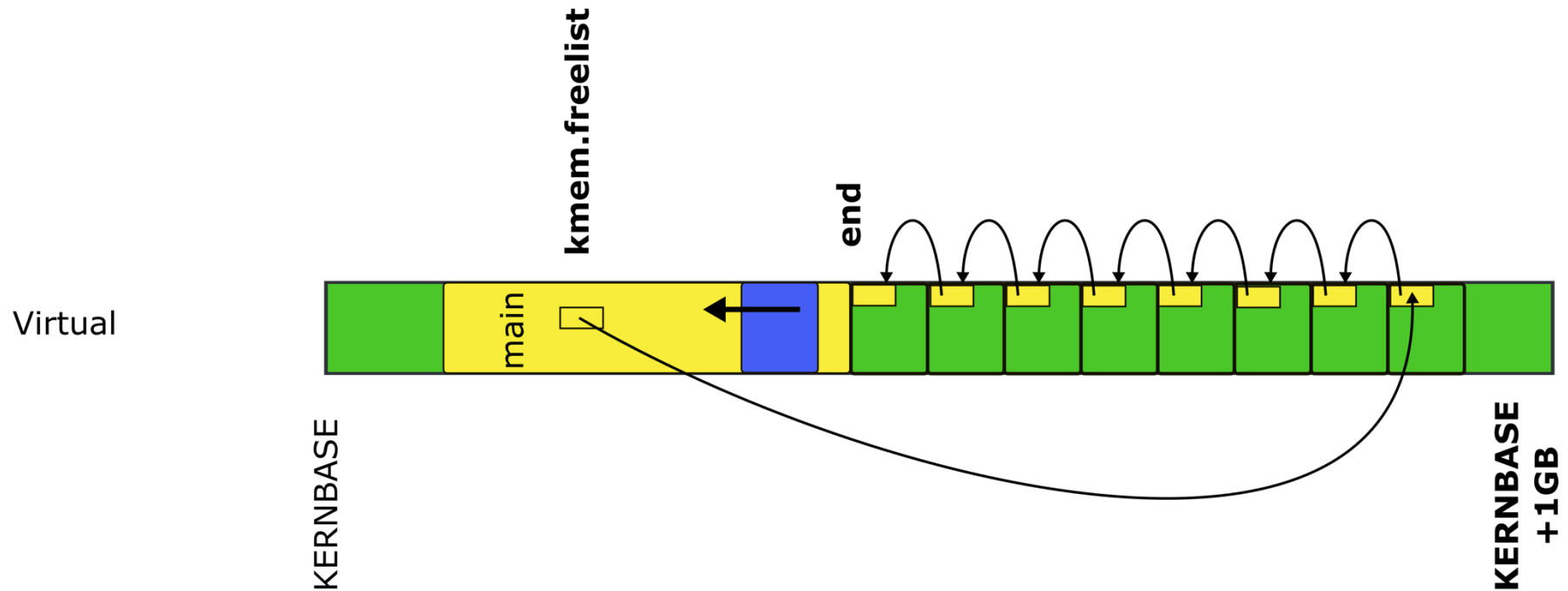# Back to main(): Kernel address space

```
1366 int
1367 main(void)
1368 {
1369   void *kinit1_end = (1024*1024*1024 > PHYSTOP) ? P2V(PHYSTOP)
1370                                     : P2V(1024*1024*1024);
1371   kinit1(end, kinit1_end); // phys page allocator
1372   kvmalloc(); // kernel page table
1373   mpinit(); // detect other processors
1374   lapicinit(); // interrupt controller
1375   seginit(); // segment descriptors
1376   picinit(); // disable pic
1377   ioapicinit(); // another interrupt controller
...
```

- What do you think has to happen?
  - i.e., how do we construct a kernel page table?

# Recap: our goal



**Virtual**

| | | | | | |
|---|---|---|---|---|---|
| User-memory (8TB) | | | | | |

0 — 8TB — 0x800 0000 0000 (8TB, KERNBASE) — Kernel-memory — $2^{64}$

**Physical**

| I/O Space | Kernel text + R/O data | Kernel R/W data + free memory | Unused by xv6 | I/O Space |
|---|---|---|---|---|

0
0x100000 (EXTMEM, or KERNLINK)
0x109000 (data)
0xe000000 (PHYSTOP) 234MB
512MB
0xfe000000
0xffffffff
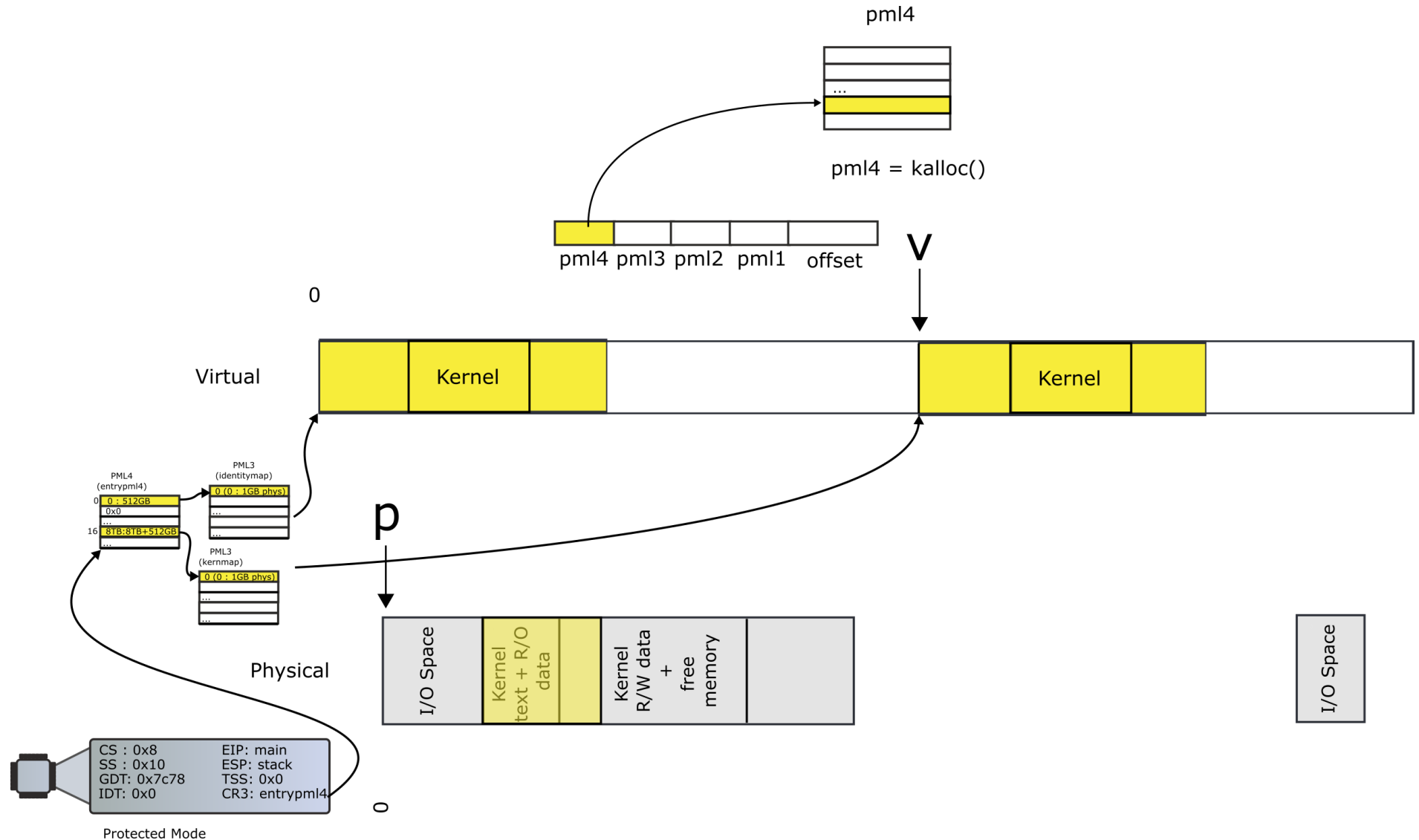
# Outline

- Map a region of virtual memory into page tables

  - Start from 2GBs

  - Iterate memory page by page

  - Allocate page table directory and page tables as we go

  - Fill in page table entries with proper physical addresses
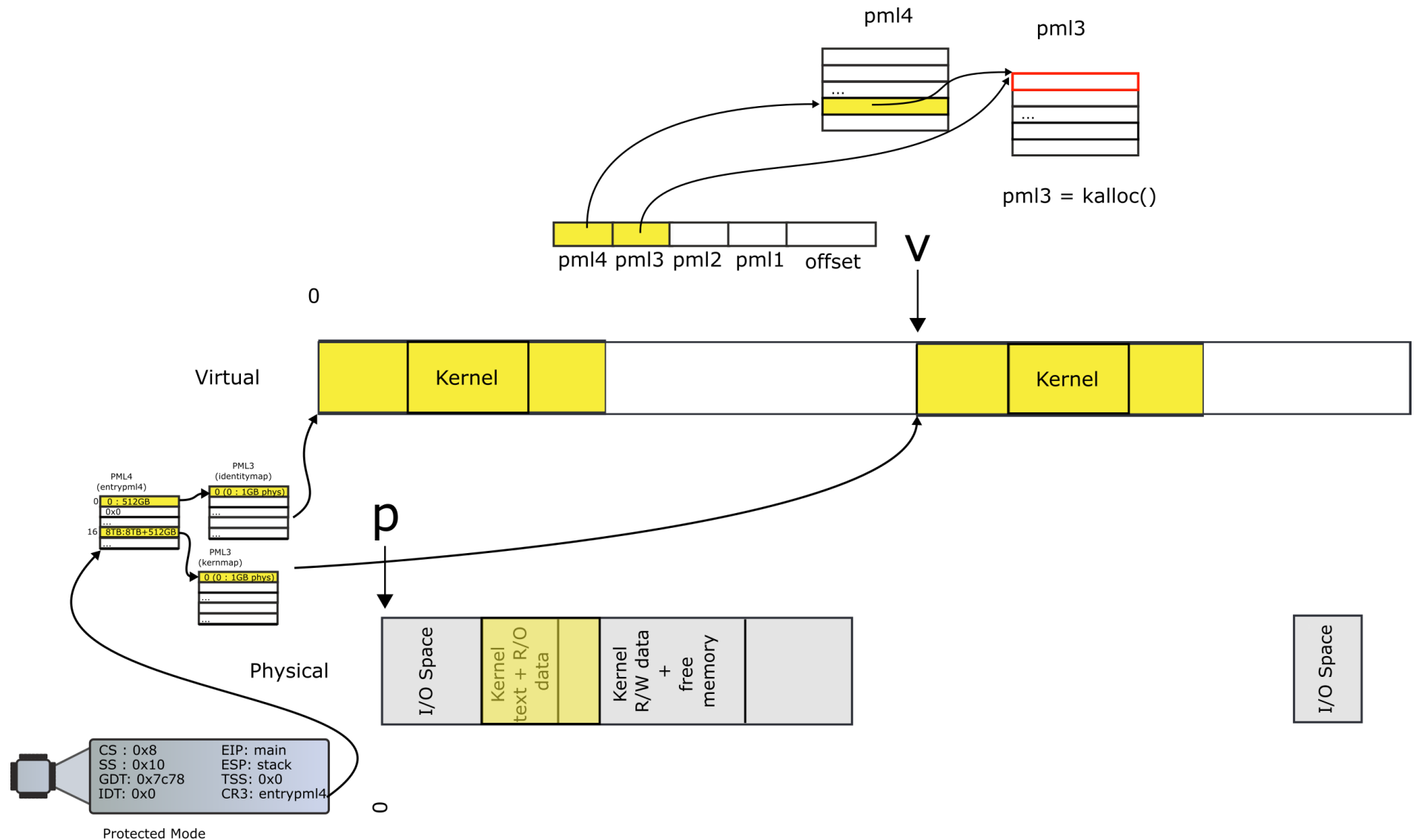
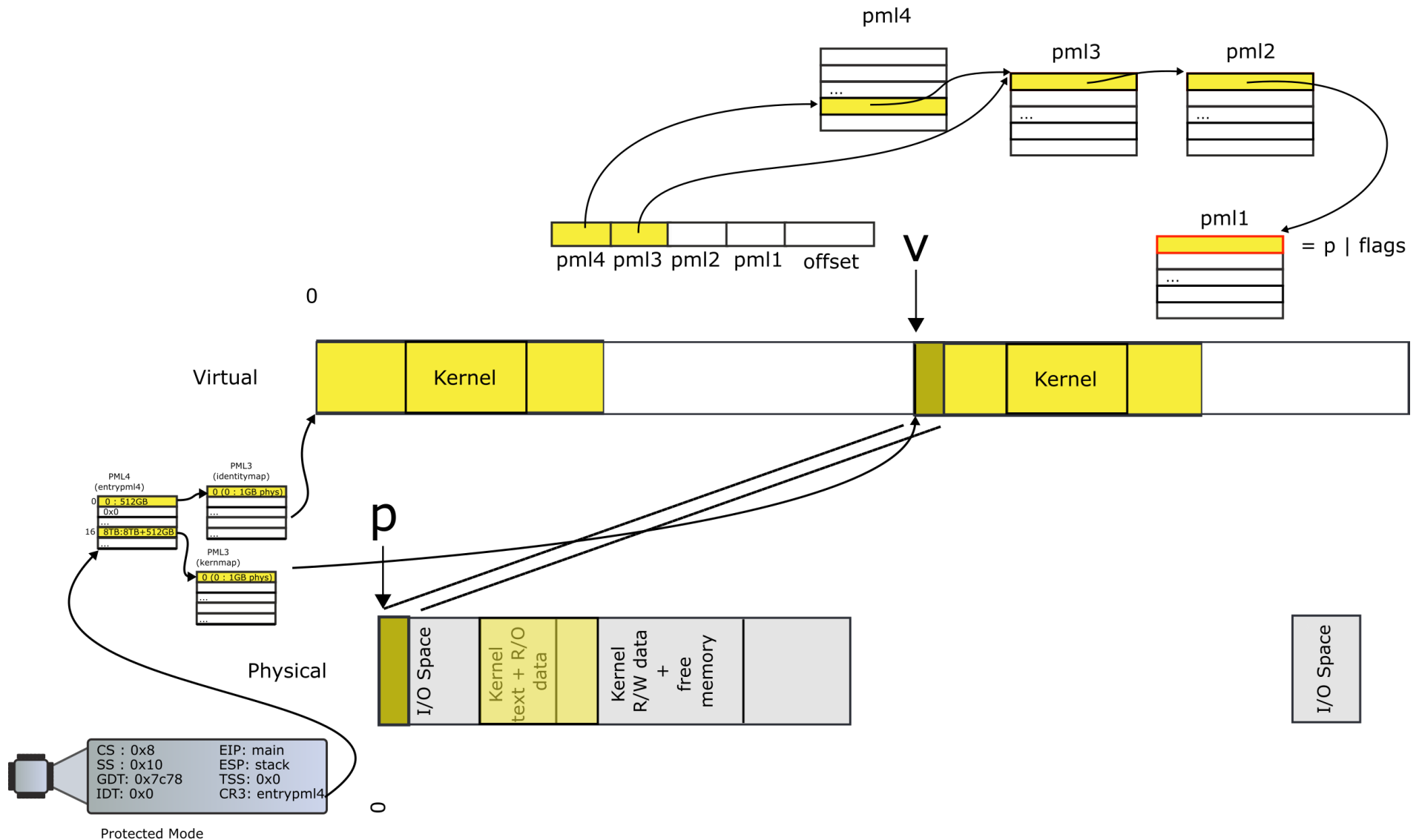# Allocate page table directory entry

# Allocate next level page table

# Locate PTE entry

# Update mapping with physical addr

# Move to next page

This is exactly what kernel is doing

(let's read the source code)

```
1363 // Bootstrap processor starts running C code here.
1364 // Allocate a real stack and switch to it, first
1365 // doing some setup required for memory allocator to work.
1366 int
1367 main(void)
1368 {
1369   void *kinit1_end = (1024*1024*1024 > PHYSTOP) ? P2V(PHYSTOP)
1370                                   : P2V(1024*1024*1024);
1371   kinit1(end, kinit1_end); // phys page allocator
1372   kvmalloc(); // kernel page table
1373   mpinit(); // detect other processors
1374   lapicinit(); // interrupt controller
1375   seginit(); // segment descriptors
1376   picinit(); // disable pic
1377   ioapicinit(); // another interrupt controller
1378   consoleinit(); // console hardware
...
}
```

Allocate page tables

# kvmalloc()

1906 // Allocate one page table for the machine for the kernel address

1907 // space for scheduler processes.

1908 void

1909 kvmalloc(void)

1910 {

1911   kpml4 = setupkvm();

1912   switchkvm();

1913 }

main()

kvmalloc()

# Allocate page table directory

```
1887 setupkvm(void)
1888 {
1889   pml4e_t *pml4;
1890   struct kmap *k;
1891
1892   if((pml4 = (pml4e_t*)kalloc()) == 0)
1893     return 0;
1894   memset(pml4, 0, PGSIZE);
1895   if (P2V(PHYSTOP) > DEV_P2V(DEVSPACE))
1896     panic("PHYSTOP too high");
1897   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1898     if(mappages(pml4, k->virt, k->phys_end - k->phys_start,
1899         (uint64)k->phys_start, k->perm) < 0) {
1900       freevm(pml4);
1901       return 0;
1902     }
1903   return pml4;
1904 }
```

main()
kvmalloc()
setupkvm()

```
1887 setupkvm(void)

1888 {

1889   pml4e_t *pml4;

1890   struct kmap *k;

1891

1892   if((pml4 = (pml4e_t*)kalloc()) == 0)

1893     return 0;

1894   memset(pml4, 0, PGSIZE);

1895   if (P2V(PHYSTOP) > DEV_P2V(DEVSPACE))

1896     panic("PHYSTOP too high");

1897   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)

1898     if(mappages(pml4, k->virt, k->phys_end - k->phys_start,

1899          (uint64)k->phys_start, k->perm) < 0) {

1900       freevm(pml4);

1901       return 0;

1902     }

1903   return pml4;

1904 }
```

What is the address of this page?

e table directory

main()

kvmalloc()

setupkvm()

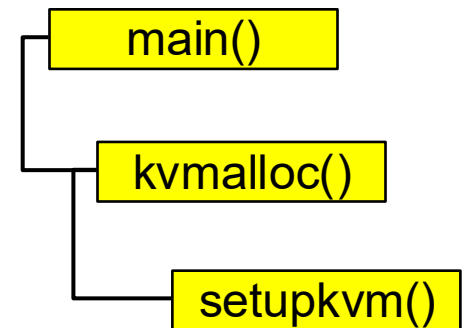https://pollev.com/cs5460

```
1887 setupkvm(void)

1888 {

1889   pml4e_t *pml4;

1890   struct kmap *k;

1891

1892   if((pml4 = (pml4e_t*)kalloc()) == 0)

1893     return 0;

1894   memset(pml4, 0, PGSIZE);

1895   if (P2V(PHYSTO

1896     panic("PHYSTO

1897   for(k = kmap; k

1898     if(mappages(p

1899         (uint64)

1900       freevm(pml4)

1901       return 0;

1902     }

1903   return pml4;

1904 }
```

What is the address of this page?

e table directory

main()

kvmalloc()

```
1887 setupkvm(void)

1888 {

1889   pml4e_t *pml4;

1890   struct kmap *k;

1891

1892   if((pml4 = (pml4e_t*)kalloc()) == 0)

1893     return 0;

1894   memset(pml4, 0, PGSIZE);

1895   if (P2V(PHYSTOP) > DEV_P2V(DEVSPACE))

1896     panic("PHYSTOP too high");

1897   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)

1898     if(mappages(pml4, k->virt, k->phys_end - k->phys_start,

1899          (uint64)k->phys_start, k->perm) < 0) {

1900       freevm(pml4);

1901       return 0;

1902     }

1903   return pml4;

1904 }
```

What is the address of this page?

table directory

main()

kvmalloc()

setupkvm()

https://pollev.com/cs5460

# Allocate page table directory

# Iterate in a loop: map physical pages
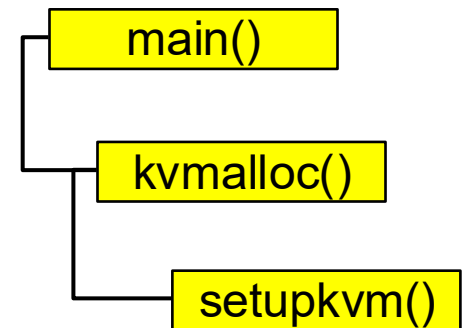
```
1887 setupkvm(void)
1888 {
1889   pml4e_t *pml4;
1890   struct kmap *k;
1891
1892   if((pml4 = (pml4e_t*)kalloc()) == 0)
1893     return 0;
1894   memset(pml4, 0, PGSIZE);
1895   if (P2V(PHYSTOP) > DEV_P2V(DEVSPACE))
1896     panic("PHYSTOP too high");
1897   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1898     if(mappages(pml4, k->virt, k->phys_end - k->phys_start,
1899            (uint64)k->phys_start, k->perm) < 0) {
1900       freevm(pml4);
1901       return 0;
1902     }
1903   return pml4;
1904 }
```

main()

kvmalloc()

setupkvm()

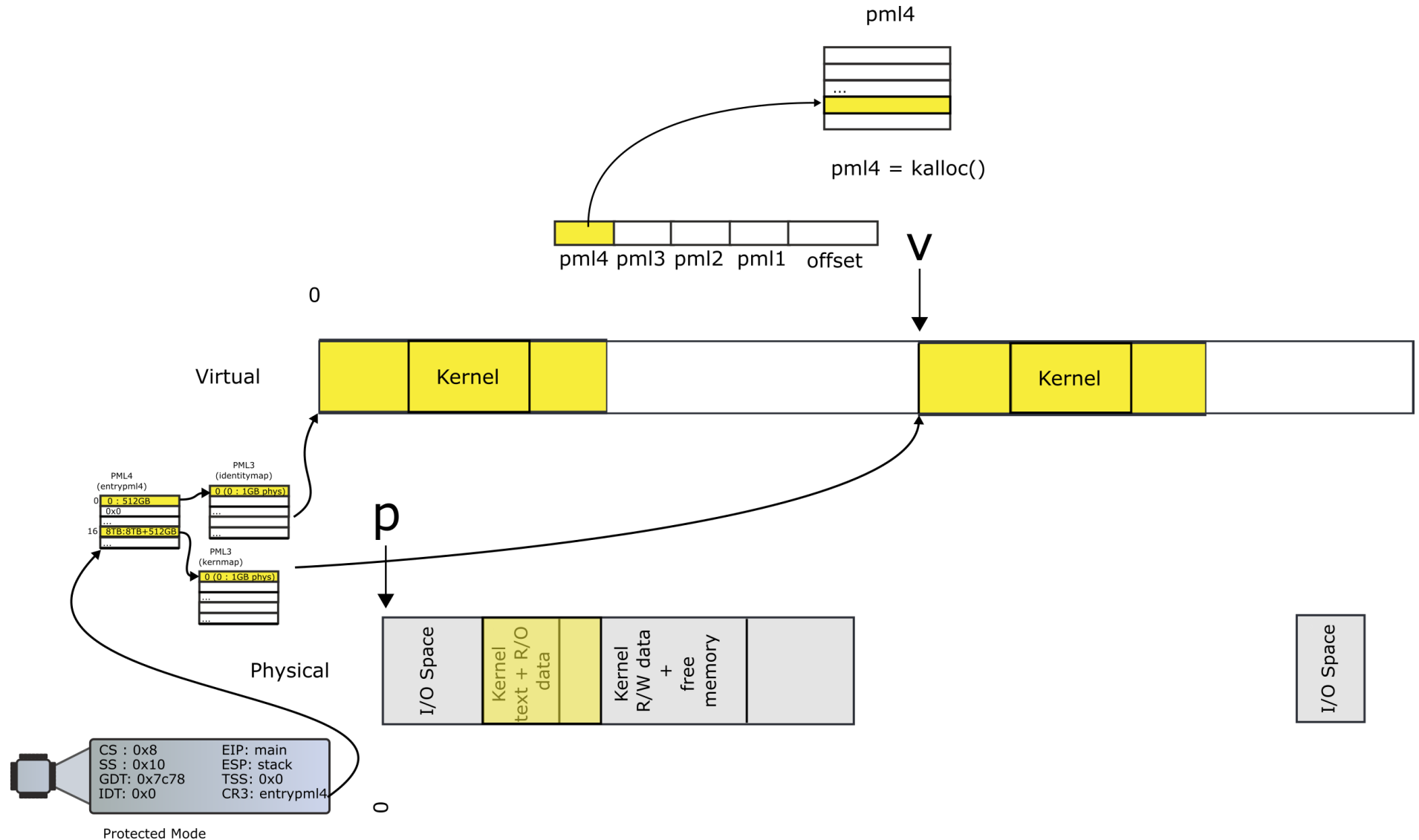# Iterate in a loop: map physical pages

```
1887 setupkvm(void)
1888 {
1889   pml4e_t *pml4;
1890   struct kmap *k;
1891
1892   if((pml4 = (pml4e_t*)kalloc()) == 0)
1893     return 0;
1894   memset(pml4, 0, PGSIZE);
1895   if (P2V(PHYSTOP) > DEV_P2V(DEVSPACE))
1896     panic("PHYSTOP too high");
1897   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1898     if(mappages(pml4, k->virt, k->phys_end - k->phys_start,
1899         (uint64)k->phys_start, k->perm) < 0) {
1900       freevm(pml4);
1901       return 0;
1902     }
1903   return pml4;
1904 }
```
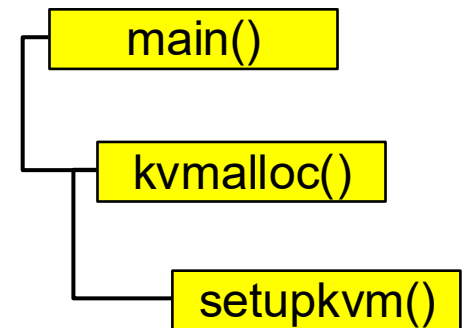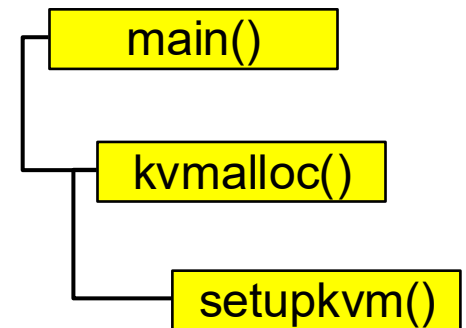
main()
kvmalloc()
setupkvm()

# Kernel map



Virtual

8TB

0

0x800 0000 0000
(8TB, KERNBASE)

$2^{64}$

User-memory

Kernel-memory

Physical

I/O Space

Kernel text + R/O data

Kernel R/W data + free memory

Unused by xv6

I/O Space

0

0x100000 (EXTMEM, or KERNLINK)

0x109000 (data)

0xe000000 (PHYSTOP) 234MB

512MB

0xfe000000

0xffffffff

# Kernel map



https://pollev.com/cs5460

# Kmap – kernel map

1823 static struct kmap {

1824    void *virt;

1825    uint phys_start;

1826    uint phys_end;

1827    int perm;

1828 } kmap[] = {

1829    { (void*)KERNBASE, 0, EXTMEM, PTE_W}, // I/O space

1830    { (void*)KERNLINK, V2P(KERNLINK), V2P(data), 0},//text+rodata

1831    { (void*)data, V2P(data), PHYSTOP, PTE_W}, // kern data+memory

1832    { (void*)DEVSPACE, DEVSPACE, 0, PTE_W}, // more devices
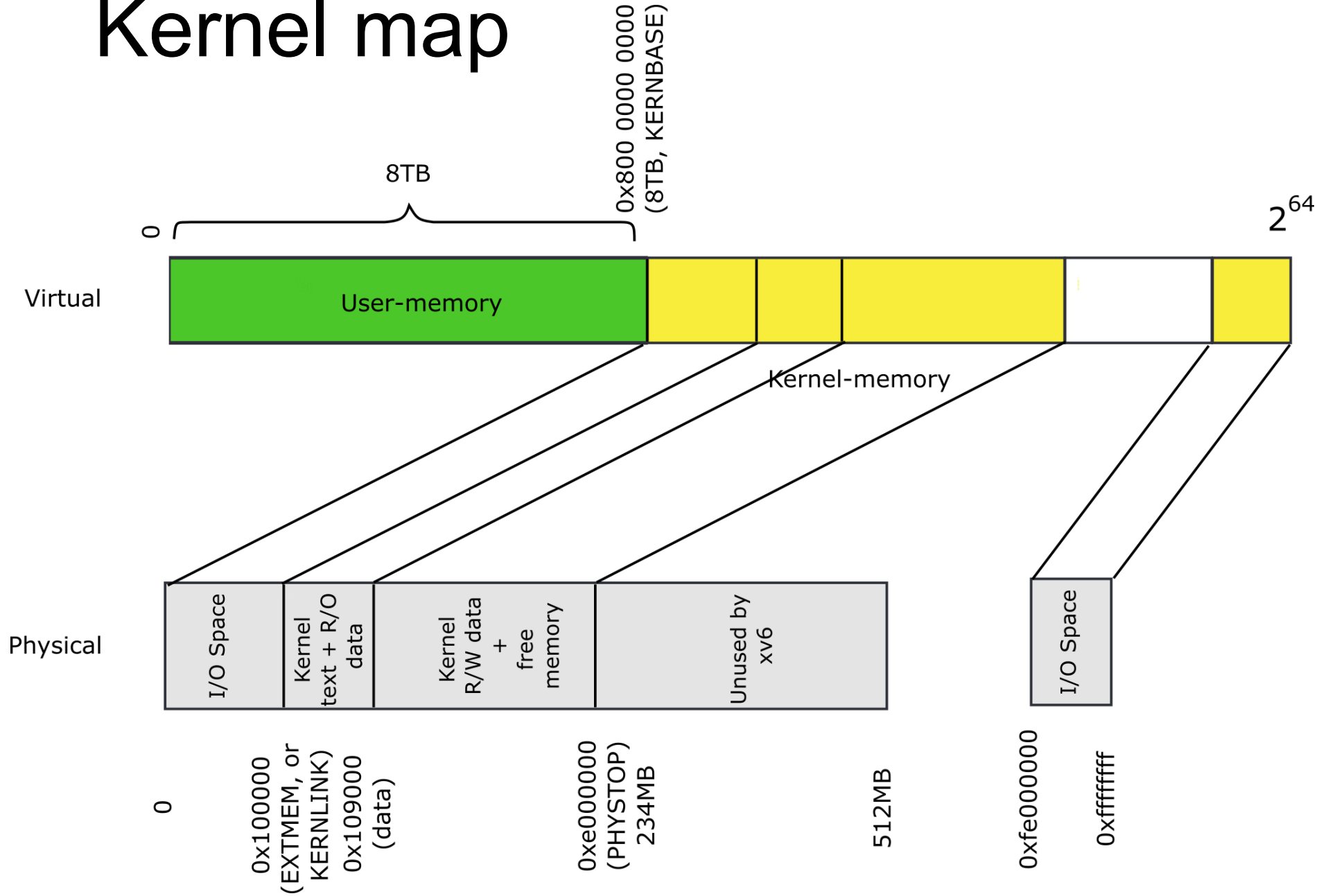
1833 };

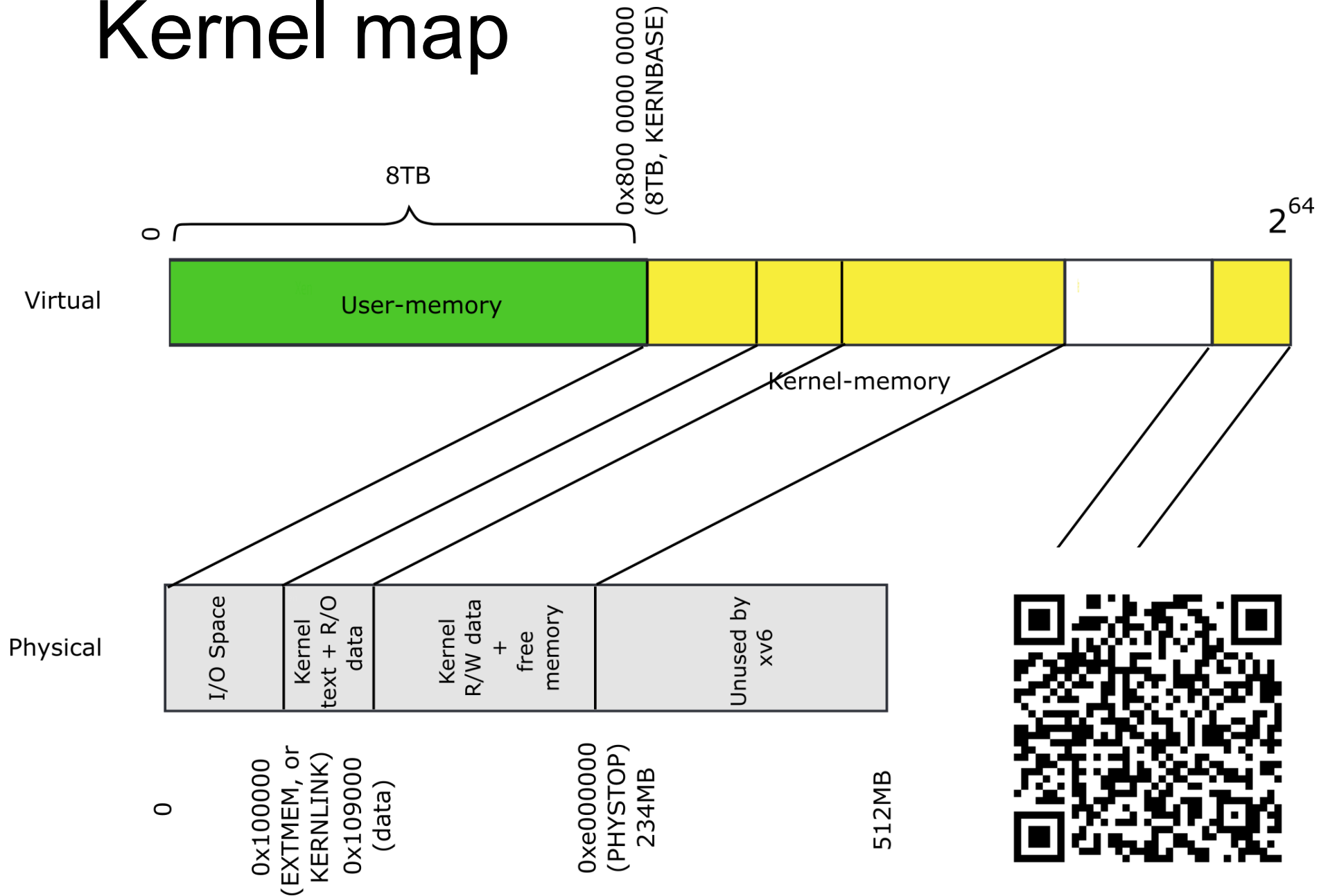# Iterate in a loop: map physical pages

```
1887 setupkvm(void)
1888 {
1889   pml4e_t *pml4;
1890   struct kmap *k;
1891
1892   if((pml4 = (pml4e_t*)kalloc()) == 0)
1893     return 0;
1894   memset(pml4, 0, PGSIZE);
1895   if (P2V(PHYSTOP) > DEV_P2V(DEVSPACE))
1896     panic("PHYSTOP too high");
1897   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1898     if(mappages(pml4, k->virt, k->phys_end - k->phys_start,
1899         (uint64)k->phys_start, k->perm) < 0) {
1900       freevm(pml4);
1901       return 0;
1902     }
1903   return pml4;
1904 }
```

main()

kvmalloc()

setupkvm()

# Inside mappages(

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```

main()
kvmalloc()
setupkvm()
mappages()

- Get the start (a) and end (last) pages fo the virtual address range we are mapping
- Then work in a loop mapping every page one by one

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```

main()

kvmalloc()

setupkvm()

mappages()

- First lookup the page table entry (pte) corresponding to the virtual address (a) we're mapping

# Locate the page table entry

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```
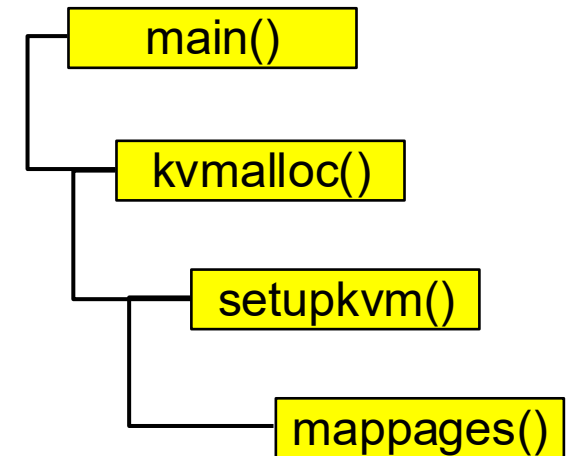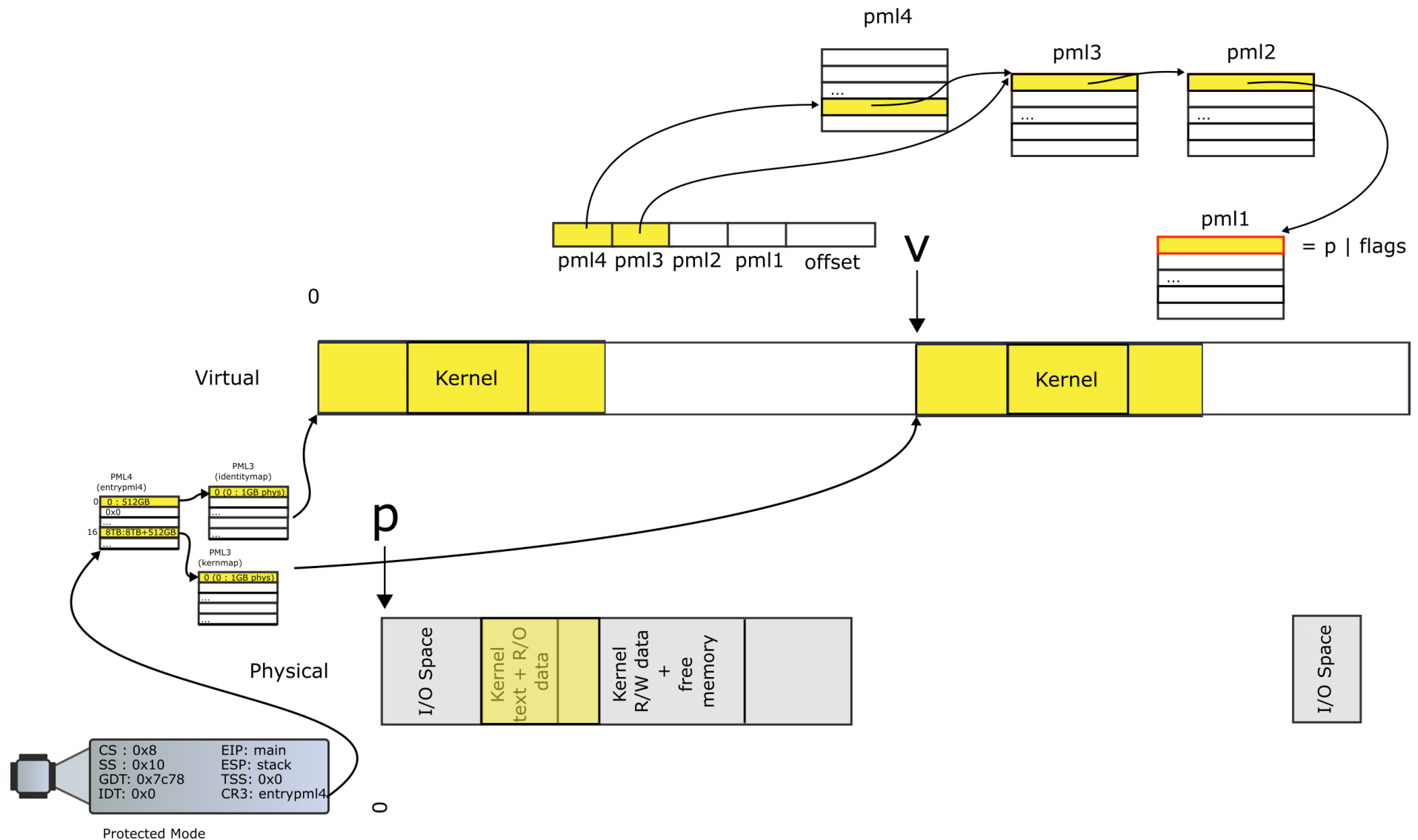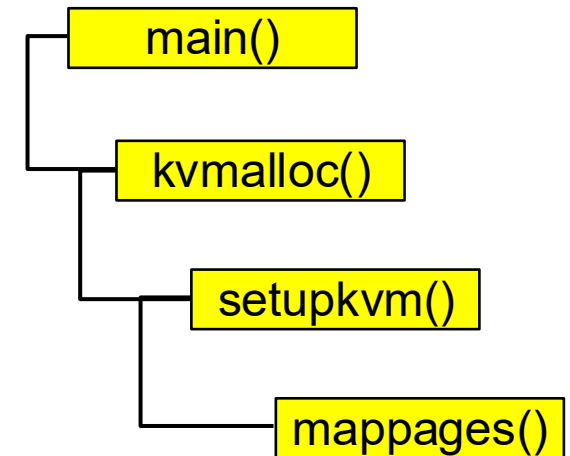


- Update the page directory entry (*pte) with the physical address (pa)

# Update mapping with physical addr



pml4

pml3

pml2

pml1

= p | flags

pml4  pml3  pml2  pml1    offset

v

0

Virtual

Kernel

Kernel

PML4
(entrypml4)

PML3
(identitymap)

0  0 : 512GB
   0x0
   ...
16 8TB:8TB+512GB

0 (0 : 1GB phys)
...

PML3
(kernmap)

0 (0 : 1GB phys)
...

p

Physical

I/O Space

Kernel
text + R/O
data

Kernel
R/W data
+
free
memory

I/O Space

0

CS : 0x8        EIP: main
SS : 0x10       ESP: stack
GDT: 0x7c78     TSS: 0x0
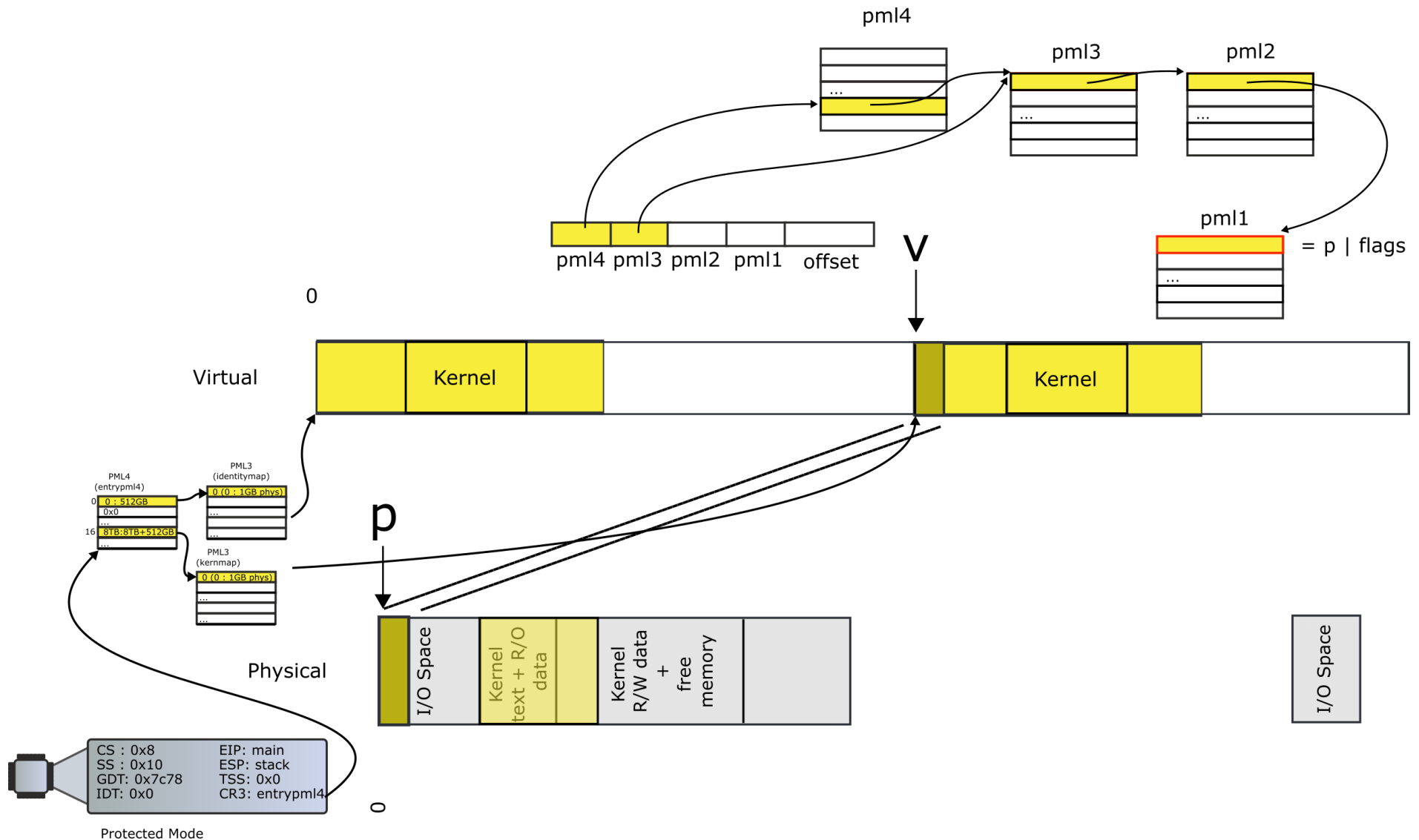IDT: 0x0        CR3: entrypml4

Protected Mode

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```
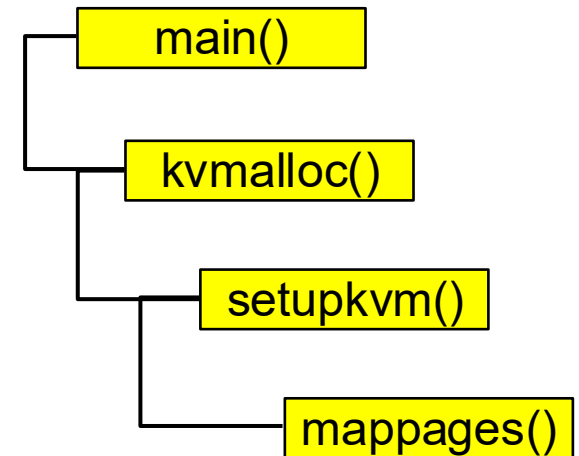


- But we need a function that locates the pte for us...

# What should it look like?

- A function takes a virtual address

- Returns a page table directory entry that maps it

Linear Address

| 47 | 39 38 | 30 29 | 21 20 | 12 11 | 0 |
|---|---|---|---|---|---|
| PML4 | Directory Ptr | Directory | Table | Offset | |

4-KByte Page

Physical Addr

PTE

Page Table

Page-Directory-Pointer Table

PDE with PS=0

Page-Directory

PDPTE

PML4E

CR3

```
1775 // Return the address of the PTE in page table pml4
1776 // that corresponds to virtual address va. If alloc!=0,
1777 // create any required page table pages.
1778 // On failure (return value equals 0), missing is set to the level
1779 // in which the table was not present.
1780 static pte_t *
1781 walkpml4(pml4e_t *pml4, const void *va, int alloc, int *missing)
1782 {
1783   return walkpglevel(pml4, va, alloc, 4, missing);
1784 }
```

walkpml4(): walk page table

- Locate the page directory entry  (*pde)

```
1734 static pte_t *
1735 walkpglevel(uint64 *pgdir, const void *va, int alloc, int level,
1736 int *missing)
1737 {
...
1741  switch (level){
1742  case 4:
1743    entry = &pgdir[PML4X(va)];
1744    break;
1745  case 3:
1746    entry = &pgdir[PDPTX(va)];
1747    break;
1748  case 2:
1749    entry = &pgdir[PDX(va)];
1750    break;
1751  case 1:
1752    return &pgdir[PTX(va)];
1753  default:
1754    panic("walkpglevel");
1755  }
1756
1757  if(*entry & PTE_P){
1758    next_pgdir = (uint64*)P2V(PTE_ADDR(*entry));
1759  } else {
...
1773 }
```
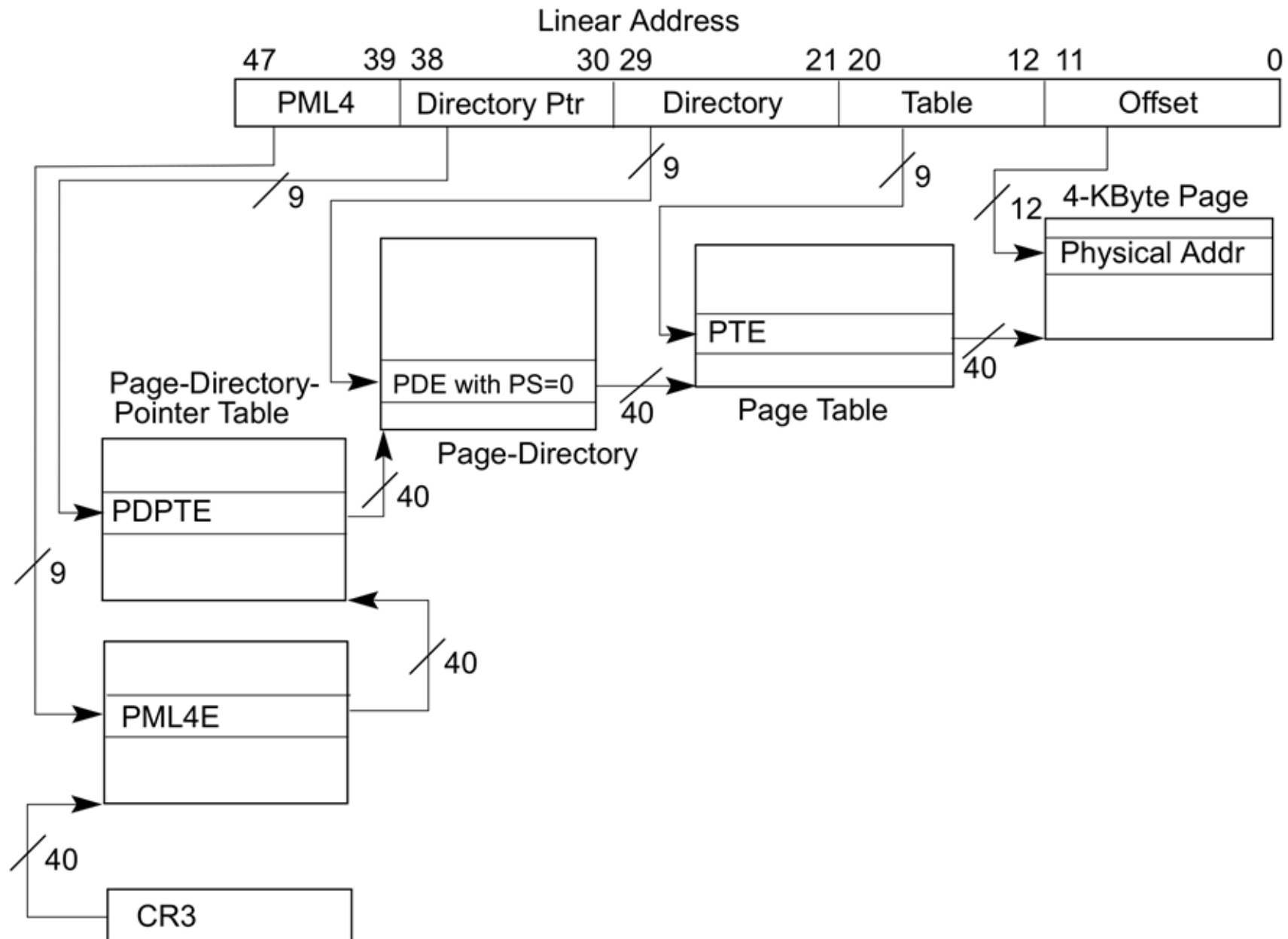
walkpml4(): walk page table

- Locate the page
  directory entry  (*pde)

```
1735 walkpglevel(uint64 *pgdir, const void *va, int alloc, int level,
1736 int *missing)
1737 {
…
1756
1757  if(*entry & PTE_P){
1758    next_pgdir = (uint64*)P2V(PTE_ADDR(*entry));
1759  } else {
1760    if(!alloc || (next_pgdir = (uint64*)kalloc()) == 0){
1761      if (missing)
1762        *missing = level;
1763      return 0;
1764    }
1765    // Make sure all those PTE_P bits are zero.
1766    memset(next_pgdir, 0, PGSIZE);
1767    // The permissions here are overly generous, but they can
1768    // be further restricted by the permissions in the page table
1769    // entries, if necessary.
1770    *entry = V2P(next_pgdir) | PTE_P | PTE_W | PTE_U;
1771  }
1772  return walkpglevel(next_pgdir, va, alloc, level - 1, missing);
1773 }
```

walkpglevel(): walk page table

- Check if page table is allocated (present)
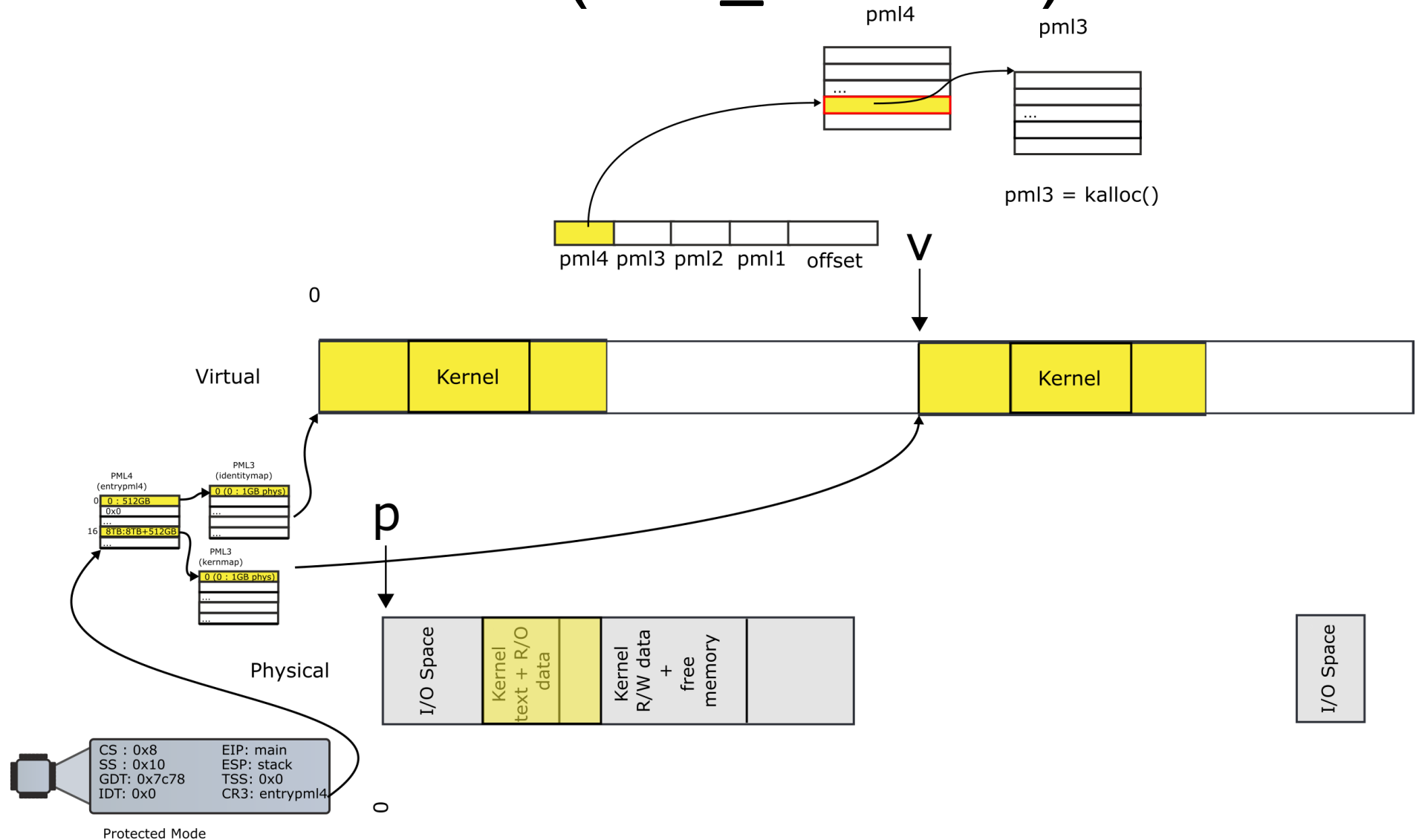
# Check if levels are allocated

```
1735 walkpglevel(uint64 *pgdir, const void *va, int alloc, int level,
1736 int *missing)
1737 {
...
1756
1757   if(*entry & PTE_P){
1758     next_pgdir = (uint64*)P2V(PTE_ADDR(*entry));
1759   } else {
1760     if(!alloc || (next_pgdir = (uint64*)kalloc()) == 0){
1761       if (missing)
1762         *missing = level;
1763       return 0;
1764     }
1765     // Make sure all those PTE_P bits are zero.
1766     memset(next_pgdir, 0, PGSIZE);
1767     // The permissions here are overly generous, but they can
1768     // be further restricted by the permissions in the page table
1769     // entries, if necessary.
1770     *entry = V2P(next_pgdir) | PTE_P | PTE_W | PTE_U;
1771   }
1772   return walkpglevel(next_pgdir, va, alloc, level - 1, missing);
1773 }
```

walkpgdir(): walk page table

- Allocate if needed

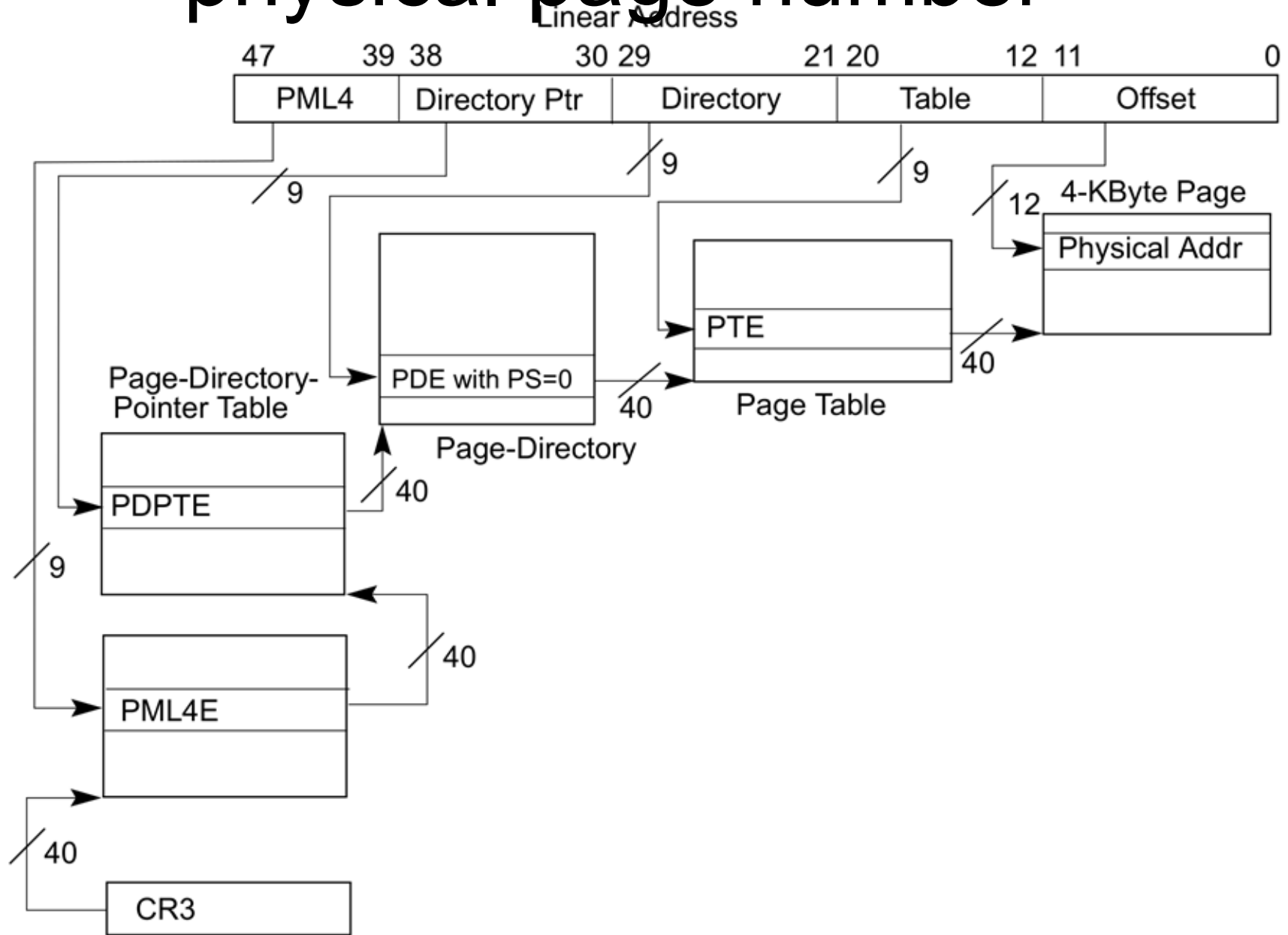# See if the next page table level exists (PTE_P is set)

```
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)
764            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 }
```

walkpgdir(): walk page table

- If exists, get the address of the next level

# PDE contains bits which represent physical page number

# Getting level 2 page

1761   pgtab = (pte_t*)P2V(PTE_ADDR(*pde));

- We need two things

- Convert from 20 bits of physical page number to physical address of the page

- PTE_ADDR(*pde)

- Convert from physical address of that page to virtual address

- P2V(...)

  - We can't access physical addresses directly

  - We can only access virtual addresses

  - Registers, mov instructions, etc. contain virtual addresses

  - Physical address have to be mapped by the current page table
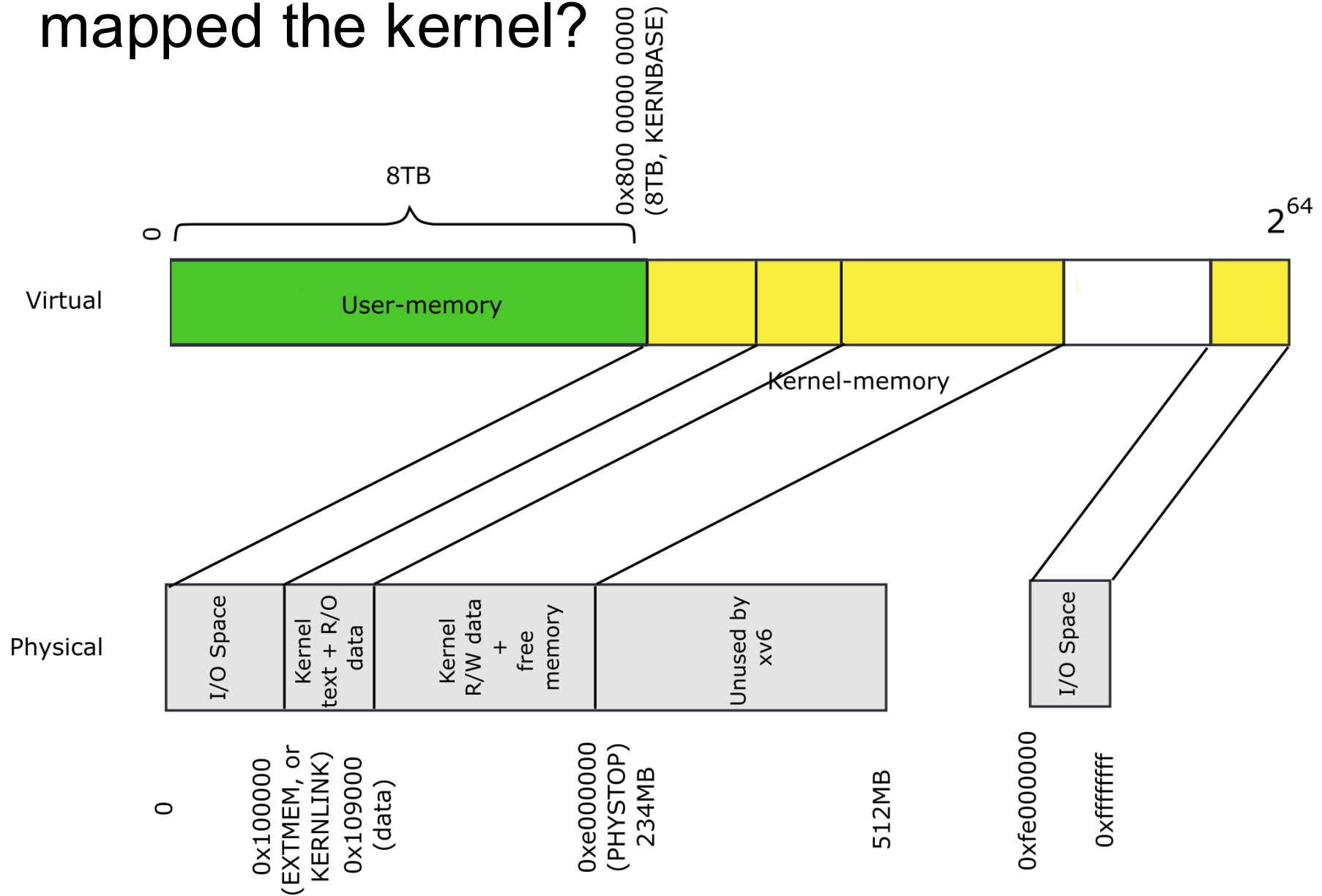
# Step 1

- Convert from bits of physical page number to physical address of the page

- PTE_ADDR(*pde)

- This is trivial

# Step 2

- Convert from physical address of that page to virtual address
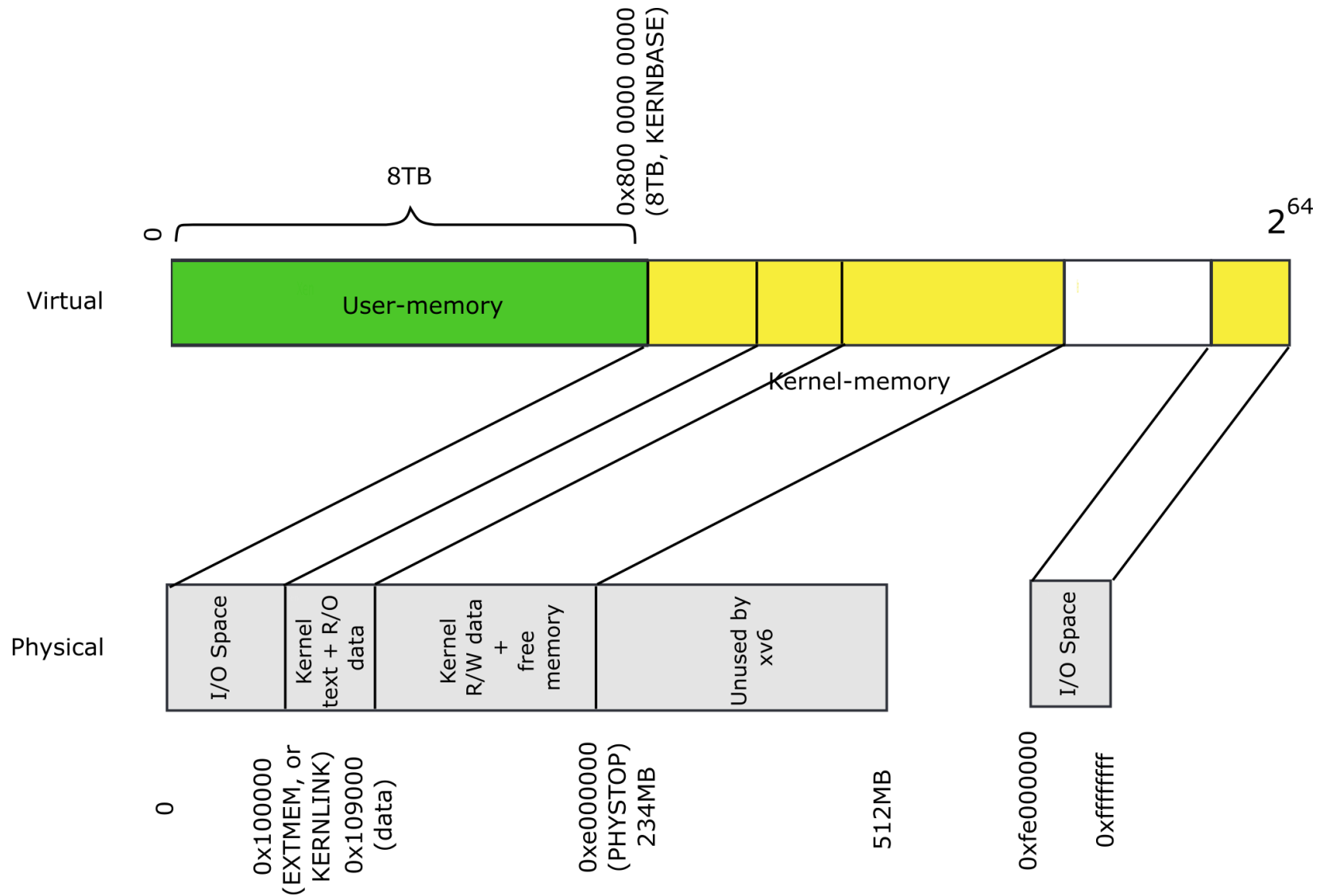
- P2V(...)

- This seems a bit tricky

# Remember how we mapped the kernel?

0207 #define KERNBASE 0x80000000 // First kernel virtual address

0210 #define V2P(a) (((uint) (a)) − KERNBASE)

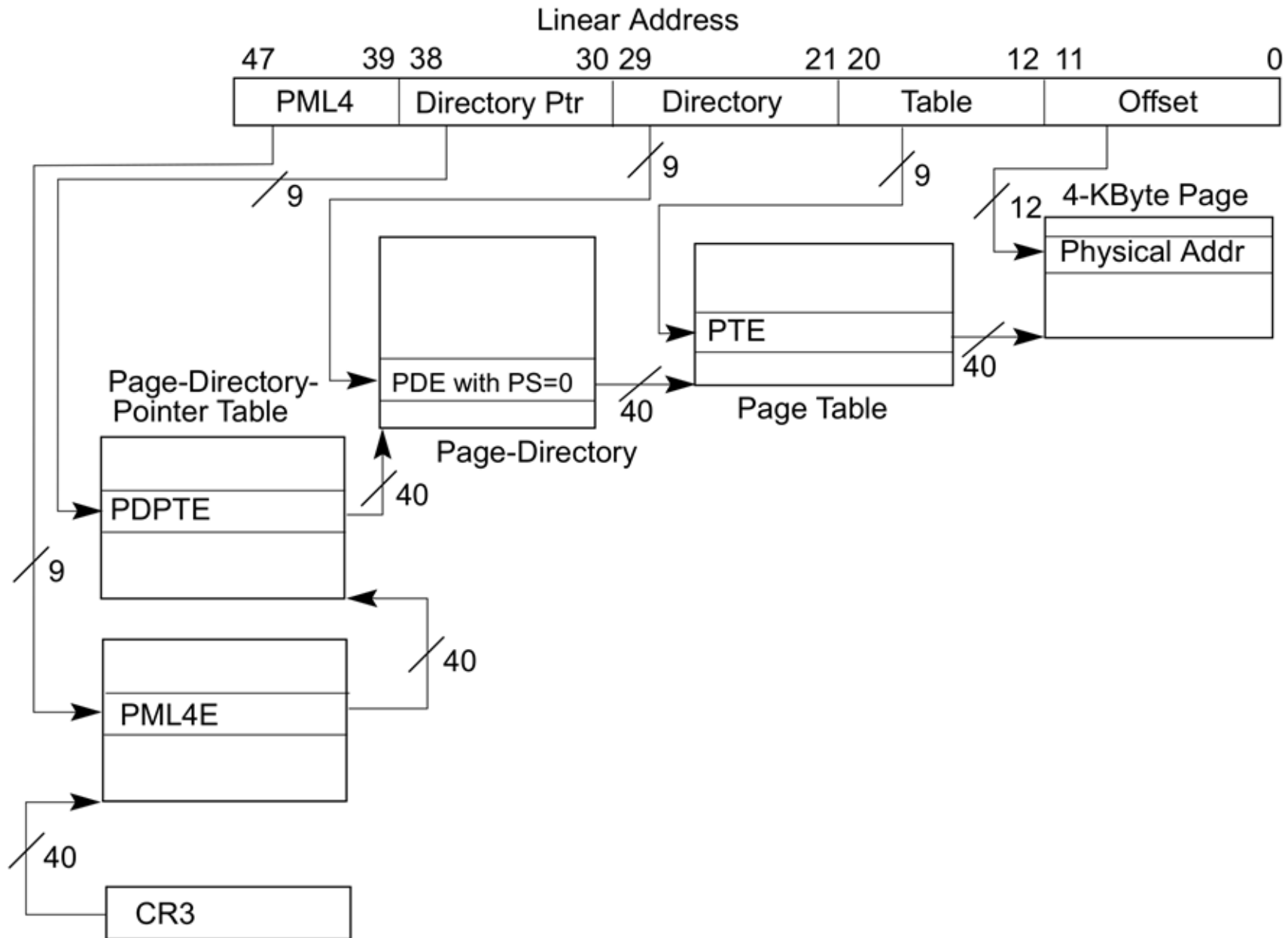0211 #define P2V(a) (((void *) (a)) + KERNBASE)

```
1735 walkpglevel(uint64 *pgdir, const void *va, int alloc, int level,
1736 int *missing)
1737 {
...
1756
1757   if(*entry & PTE_P){
1758     next_pgdir = (uint64*)P2V(PTE_ADDR(*entry));
1759   } else {
1760     if(!alloc || (next_pgdir = (uint64*)kalloc()) == 0){
1761       if (missing)
1762         *missing = level;
1763       return 0;
1764     }
1765     // Make sure all those PTE_P bits are zero.
1766     memset(next_pgdir, 0, PGSIZE);
1767     // The permissions here are overly generous, but they can
1768     // be further restricted by the permissions in the page table
1769     // entries, if necessary.
1770     *entry = V2P(next_pgdir) | PTE_P | PTE_W | PTE_U;
1771   }
1772   return walkpglevel(next_pgdir, va, alloc, level - 1, missing);
1773 }
```
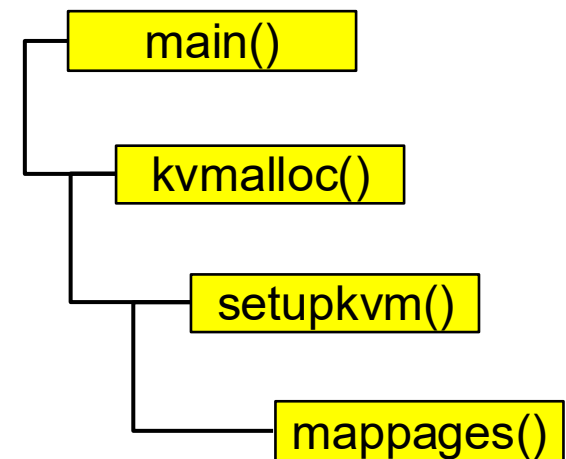
walkpgdir(): walk page table

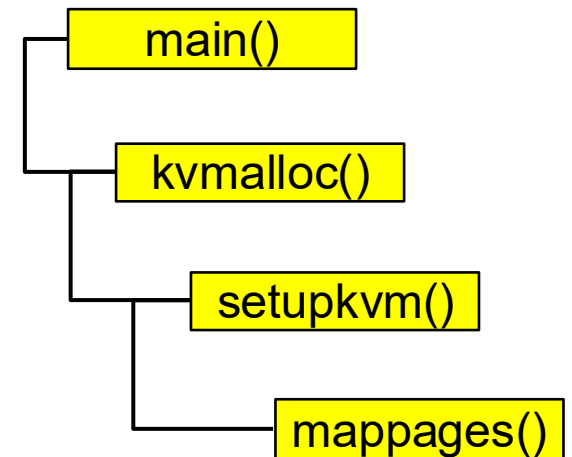- Go to the next level

# Return a pointer to PTE

Back to `mappages()` function that maps a region of virtual memory into continuous region of physical memory

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```



main()
kvmalloc()
setupkvm()
mappages()

Remember we just
discussed walkpml4()

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```



Page present
(PTE_P) – panic

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```

main()

kvmalloc()

setupkvm()

mappages()

- Update page table entry
- Where does *pte point?
  - pa – physical address of the page

```
1804 mappages(pml4e_t *pml4, void *va, uint size, uint64 pa, int perm)
1805 {
1806   char *a, *last;
1807   pte_t *pte;
1808
1809   a = (char*)PGROUNDDOWN((uint64)va);
1810   last = (char*)PGROUNDDOWN(((uint64)va) + size - 1);
1811   for(;;){
1812     if((pte = walkpml4(pml4, a, 1, 0)) == 0)
1813       return -1;
1814     if(*pte & PTE_P)
1815       panic("remap");
1816     *pte = pa | perm | PTE_P;
1817     if(a == last)
1818       break;
1819     a += PGSIZE;
1820     pa += PGSIZE;
1821   }
1822   return 0;
1823 }
```
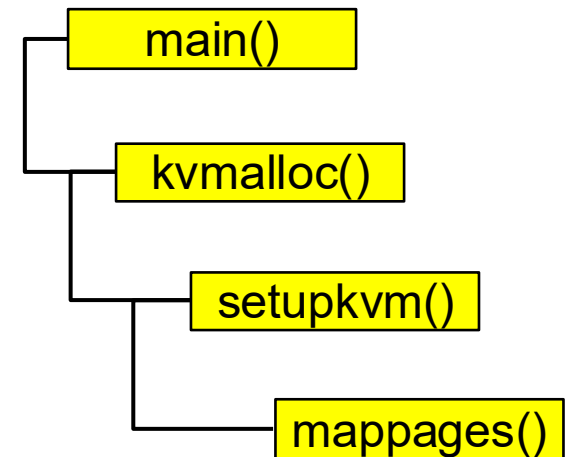


- Move to the next page

# kvmalloc()

1906 // Allocate one page table for the machine for the kernel address

1907 // space for scheduler processes.

1908 void

1909 kvmalloc(void)

1910 {

1911   kpml4 = setupkvm();

1912   switchkvm();

1913 }

main()

kvmalloc()

# Switch to the new page table

1915 // Switch h/w page table register to the kernel-only page table,

1916 // for when no process is running.

1917 void

1918 switchkvm(void)

1919 {

1920   lcr3(V2P(kpml4)); // switch to the kernel page table

1921 }

```
main()
kvmalloc()
setupkvm()
switchkvm()
```

# Recap

- Kernel has a memory allocator

- Kernel has a its own address space

- It uses 4KB page tables

- It is ready to create processes

```
1366 int
1367 main(void)
1368 {
1369   void *kinit1_end = (1024*1024*1024 > PHYSTOP) ? P2V(PHYSTOP)
1370                                   : P2V(1024*1024*1024);
1371   kinit1(end, kinit1_end); // phys page allocator
1372   kvmalloc(); // kernel page table
1373   mpinit(); // detect other processors
1374   lapicinit(); // interrupt controller
1375   seginit(); // segment descriptors
1376   picinit(); // disable pic
1377   ioapicinit(); // another interrupt controller
...
```

main()

# Initialize GDT

```
1712 // Set up CPU's kernel segment descriptors.

1713 // Run once on entry on each CPU.

1714 void

1715 seginit(void)

1716 {

1717   struct cpu *c;

1718

1719   // Map "logical" addresses to virtual addresses using identity map.

1720   // Cannot share a CODE descriptor for both kernel and user

1721   // because it would have to have DPL_USR, but the CPU forbids

1722   // an interrupt from CPL=0 to DPL=3.

1723   c = &cpus[cpuid()];

1724   c->gdt[SEG_KCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, 0);

1725   c->gdt[SEG_KDATA] = SEG(STA_W, 0, 0xffffffff, 0);

1726   c->gdt[SEG_UCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, DPL_USER);

1727   c->gdt[SEG_UDATA] = SEG(STA_W, 0, 0xffffffff, DPL_USER);

1728   lgdt(c->gdt, sizeof(c->gdt));

1729 }
```

main()

seginit()

# Struct CPU

2300 // Per–CPU state

2301 struct cpu {

2302   uchar apicid;            // Local APIC ID

2303   struct context *scheduler;   // swtch() here to enter scheduler

2304   struct taskstate ts;        // Used by x86 to find stack for interrupt

2305   struct segdesc gdt[NSEGS];    // x86 global descriptor table

2306   volatile uint started;      // Has the CPU started?

2307   int ncli;              // Depth of pushcli nesting.

2308   int intena;             // Were interrupts enabled before pushcli?

2309   struct proc *proc;         // The process running on this cpu or null

2310 };

2311

2312 extern struct cpu cpus[NCPU];

# Segment descriptor

| 31 | | 24 | 23 | 22 | 21 | 20 | 19 | | 16 | 15 | 14 | 13 | 12 | 11 | | 8 | 7 | | 0 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Base 31:24 | | | G | D/B | L | AVL | Seg. Limit 19:16 | | | P | DPL | | S | Type | | | Base 23:16 | | | 4 |

| 31 | | 16 | 15 | | 0 | |
|---|---|---|---|---|---|---|
| Base Address 15:00 | | | Segment Limit 15:00 | | | 0 |

L      — 64-bit code segment (IA-32e mode only)
AVL    — Available for use by system software
BASE — Segment base address
D/B    — Default operation size (0 = 16-bit segment; 1 = 32-bit segment)
DPL    — Descriptor privilege level
G      — Granularity
LIMIT — Segment Limit
P      — Segment present
S      — Descriptor type (0 = system; 1 = code or data)
TYPE — Segment type

# Segment Descriptor

```
0724 // Segment Descriptor
0725 struct segdesc {
0726   uint lim_15_0 : 16;  // Low bits of segment limit
0727   uint base_15_0 : 16; // Low bits of segment base address
0728   uint base_23_16 : 8; // Middle bits of segment base address
0729   uint type : 4;     // Segment type (see STS_ constants)
0730   uint s : 1;        // 0 = system, 1 = application
0731   uint dpl : 2;      // Descriptor Privilege Level
0732   uint p : 1;        // Present
0733   uint lim_19_16 : 4;  // High bits of segment limit
0734   uint avl : 1;      // Unused (available for software use)
0735   uint rsv1 : 1;     // Reserved
0736   uint db : 1;       // 0 = 16-bit segment, 1 = 32-bit segment
0737   uint g : 1;        // Granularity: limit scaled by 4K when set
0738   uint base_31_24 : 8; // High bits of segment base address
0739 };
```

# Thank you!

## (Next time: interrupts!)