# CS5460/6460: Operating Systems
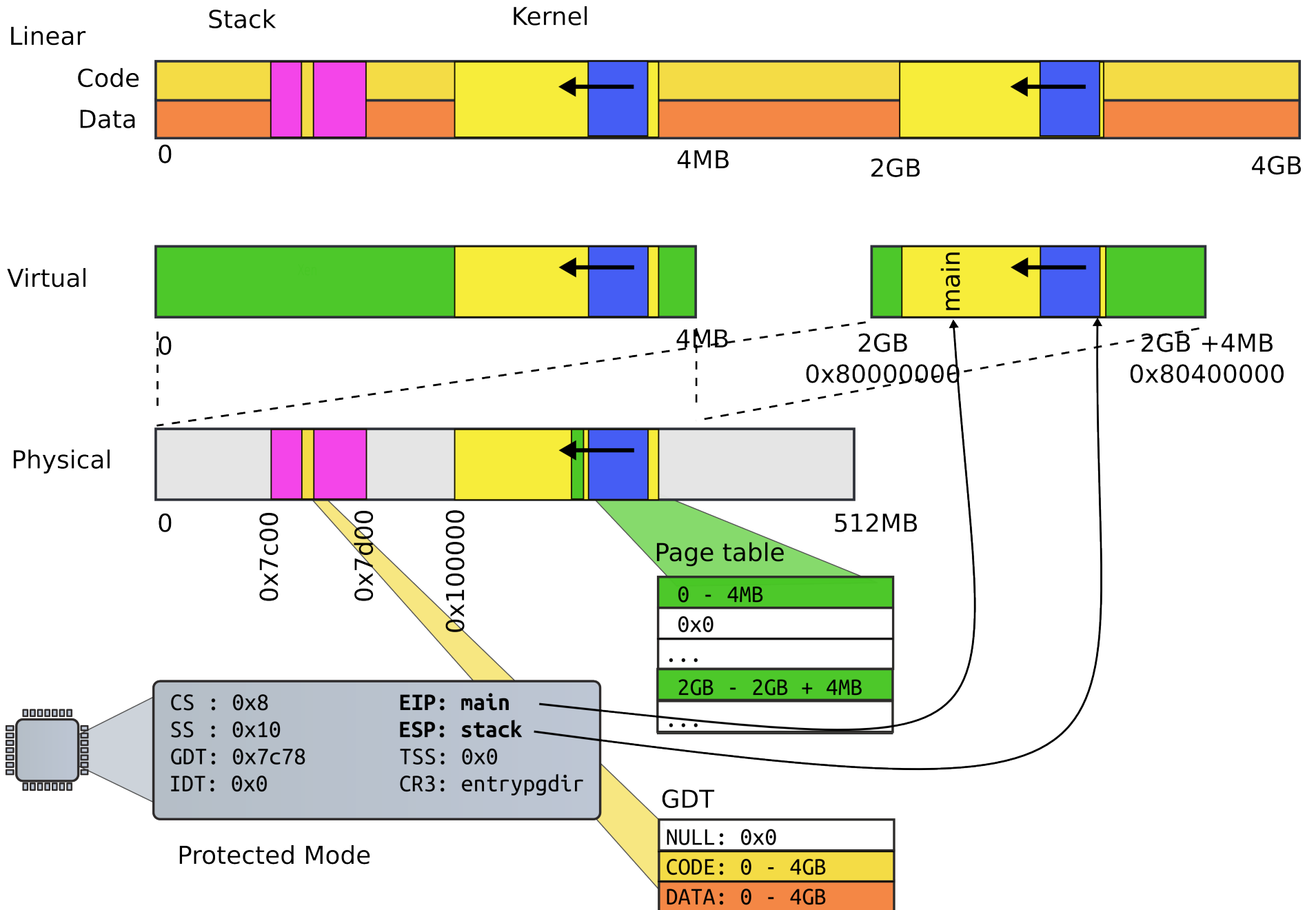
# Lecture 8: System init

Anton Burtsev
January, 2014

# Recap from last time

- Setup segments (data and code)
- Switched to protected mode
    - Loaded GDT (segmentation is on)
- Setup stack (to call C functions)
- Loaded kernel from disk
- Setup first page table
    - 2 entries [ 0 : 4MB ] and [ 2GB : (2GB + 4MB) ]
- Setup high-address stack
- Jumped to main()

```
1053 # Set up the stack pointer.
1054 movl $(stack + KSTACKSIZE), %esp
1055
1056 # Jump to main(), and switch to executing at
1057 # high addresses. The indirect call is needed because
1058 # the assembler produces a PC-relative instruction
1059 # for a direct jump.
1060 mov $main, %eax
1061 jmp *%eax
1062
1063 .comm stack, KSTACKSIZE
```

# Jumped to main()



Linear

Stack            Kernel

Code

Data

0            4MB        2GB        4GB

Virtual

main

0            4MB     2GB        2GB +4MB

0x80000000        0x80400000

Physical

0     0x7c00     0x7d00     0x100000        512MB

Page table

| 0 - 4MB |
| --- |
| 0x0 |
| ... |
| 2GB - 2GB + 4MB |
| ... |

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

Protected Mode

GDT

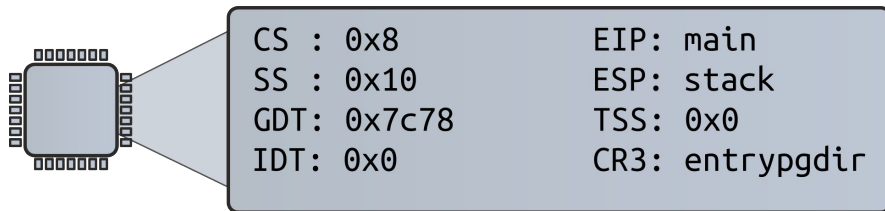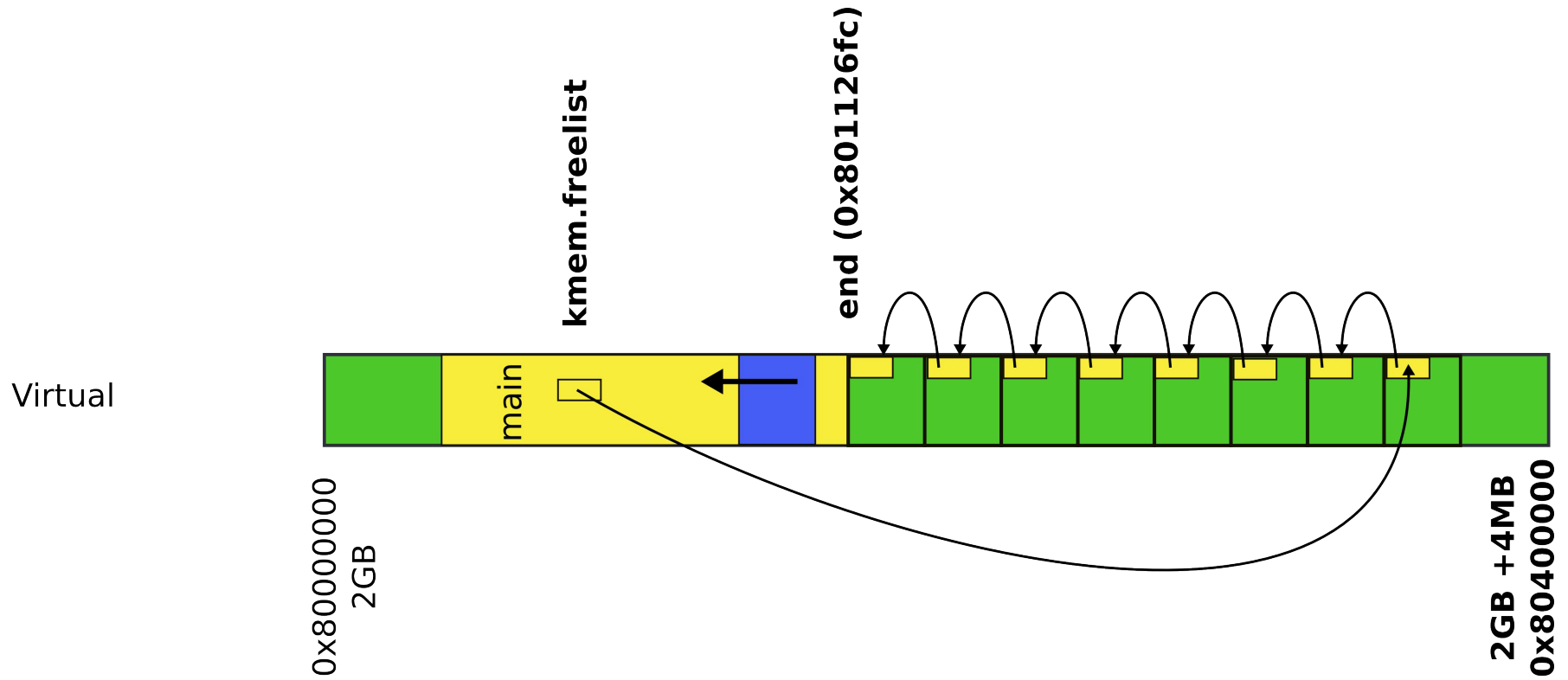| NULL: 0x0 |
| --- |
| CODE: 0 - 4GB |
| DATA: 0 - 4GB |

```
1217 main(void)
1218 {
1219   kinit1(end, P2V(4*1024*1024)); // phys page allocator
1220   kvmalloc(); // kernel page table
1221   mpinit(); // collect info about this machine
1222   lapicinit();
1223   seginit(); // set up segments
1224   cprintf("\ncpu%d: starting xv6\n\n", cpu->id);
1225   picinit(); // interrupt controller
1226   ioapicinit(); // another interrupt controller
1227   consoleinit(); // I/O devices & their interrupts
1228   uartinit(); // serial port
1229   pinit(); // process table
1230   tvinit(); // trap vectors
```

# Physical page allocator

- Goal:
  - List of free physical pages
  - To allocate page tables, stacks, data structures, etc.
  - Remember current page table is only 1! page

- Where to get memory to keep the list itself?
  - 1 level, only 4MB entries
    - You don't even have space to keep the second level page tables

# Physical page allocator



**kmem.freelist**

**end (0x801126fc)**

Virtual

0x80000000
2GB

**2GB +4MB
0x80400000**

main

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

Protected Mode

```
2780 kinit1(void *vstart, void *vend)
2781 {
...
2784   freerange(vstart, vend);
2785 }


2801 freerange(void *vstart, void *vend)
2802 {
2803   char *p;
2804   p = (char*)PGROUNDUP((uint)vstart);
2805   for(; p + PGSIZE <= (char*)vend; p += PGSIZE)
2806     kfree(p);
2807 }
```

```
2815 kfree(char *v)
2816 {
2817   struct run *r;

...

2827   r = (struct run*)v;
2828   r->next = kmem.freelist;
2829   kmem.freelist = r;

...

2832 }
```

# Kernel page table

```
1217 main(void)

1218 {

1219   kinit1(end, P2V(4*1024*1024)); // phys page allocator

1220   kvmalloc(); // kernel page table

1221   mpinit(); // collect info about this machine

1222   lapicinit();

1223   seginit(); // set up segments
```

# kvmalloc()

```
1757 kvmalloc(void)
1758 {
1759    kpgdir = setupkvm();
1760    switchkvm();
1761 }
```

```
1736 pde_t*
1737 setupkvm(void)
1738 {
1739   pde_t *pgdir;
1740   struct kmap *k;
1741
1742   if((pgdir = (pde_t*)kalloc()) == 0)
1743     return 0;
1744   memset(pgdir, 0, PGSIZE);
...
1747   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1748     if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
1749                 (uint)k->phys_start, k->perm) < 0)
1750       return 0;
1751   return pgdir;
1752 }
```
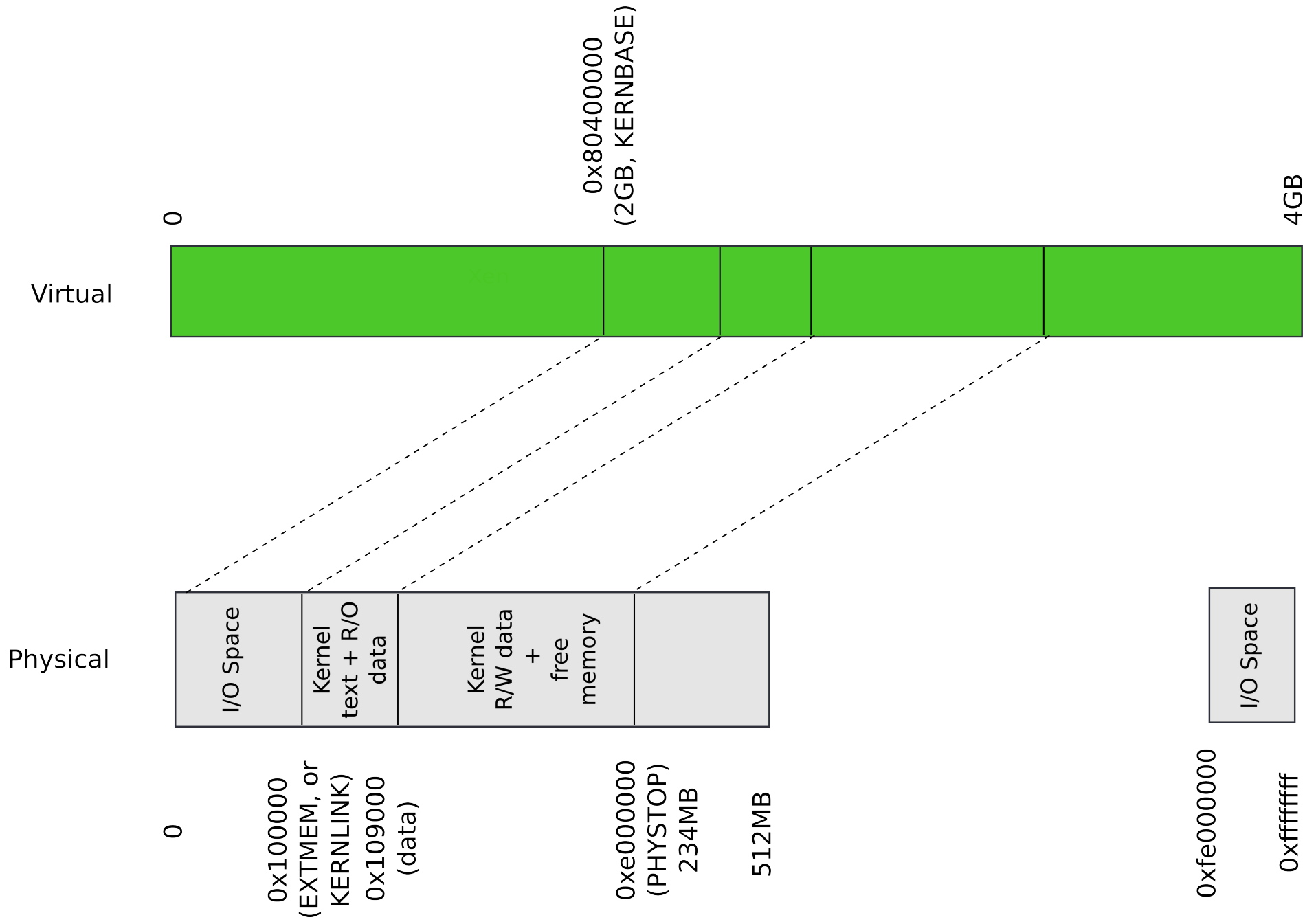
# Kalloc() - kernel allocator

```
2837 char*

2838 kalloc(void)

2839 {

2840   struct run *r;

...

2844   r = kmem.freelist;

2845   if(r)

2846     kmem.freelist = r->next;

…

2849   return (char*)r;

2850 }
```

```c
1736 pde_t*
1737 setupkvm(void)
1738 {
1739   pde_t *pgdir;
1740   struct kmap *k;
1741
1742   if((pgdir = (pde_t*)kalloc()) == 0)
1743     return 0;
1744   memset(pgdir, 0, PGSIZE);
...
1747   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1748     if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
1749                 (uint)k->phys_start, k->perm) < 0)
1750       return 0;
1751   return pgdir;
1752 }
```

# Kmap – kernel map

```
1723 static struct kmap {
1724 void *virt;
1725 uint phys_start;
1726 uint phys_end;
1727 int perm;
1728 } kmap[] = {
1729 { (void*)KERNBASE, 0, EXTMEM, PTE_W}, // I/O space
1730 { (void*)KERNLINK, V2P(KERNLINK), V2P(data), 0}, // kern
text+rodata
1731 { (void*)data, V2P(data), PHYSTOP, PTE_W}, // kern data+memory
1732 { (void*)DEVSPACE, DEVSPACE, 0, PTE_W}, // more devices
1733 };
```

Virtual

0      0x80400000 (2GB, KERNBASE)      4GB

Physical

| I/O Space | Kernel text + R/O data | Kernel R/W data + free memory | | I/O Space |

0

0x100000 (EXTMEM, or KERNLINK)

0x109000 (data)

0xe000000 (PHYSTOP) 234MB

512MB

0xfe000000

0xffffffff

```c
1679 mappages(pde_t *pgdir, void *va, uint size, uint pa, int perm)
1680 {
1681   char *a, *last;
1682   pte_t *pte;
1683
1684   a = (char*)PGROUNDDOWN((uint)va);
1685   last = (char*)PGROUNDDOWN(((uint)va) + size - 1);
1686   for(;;){
1687     if((pte = walkpgdir(pgdir, a, 1)) == 0)
1688       return -1;
1689     if(*pte & PTE_P)
1690       panic("remap");
1691     *pte = pa | perm | PTE_P;
1692     if(a == last)
1693       break;
1694     a += PGSIZE;
1695     pa += PGSIZE;
1696   }
1697   return 0;
1698 }
```

# PDX()

```
0805 // +--------10------+------10-------+--------12---------+
0806 // | Page Directory |   Page Table  | Offset within Page |
0807 // |    Index       |     Index     |                    |
0808 // +----------------+---------------+--------------------+
0809 // \--- PDX(va) --/  \--- PTX(va) --/
0810
0811 // page directory index
0812 #define PDX(va) (((uint)(va) >> PDXSHIFT) & 0x3FF)
...
0827 #define PDXSHIFT 22 // offset of PDX in a linear address
```

```c
1654 walkpgdir(pde_t *pgdir, const void *va, int alloc)
1655 {
1656   pde_t *pde;
1657   pte_t *pgtab;
1658
1659   pde = &pgdir[PDX(va)];
1660   if(*pde & PTE_P){
1661     pgtab = (pte_t*)p2v(PTE_ADDR(*pde));
1662   } else {
1663     if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
1664       return 0;
1665     // Make sure all those PTE_P bits are zero.
1666     memset(pgtab, 0, PGSIZE);
...
1670     *pde = v2p(pgtab) | PTE_P | PTE_W | PTE_U;
1671   }
1672   return &pgtab[PTX(va)];
1673 }
```

```
1654 walkpgdir(pde_t *pgdir, const void *va, int alloc)
1655 {
1656   pde_t *pde;
1657   pte_t *pgtab;
1658
1659   pde = &pgdir[PDX(va)];
1660   if(*pde & PTE_P){
1661     pgtab = (pte_t*)p2v(PTE_ADDR(*pde));
1662   } else {
1663     if(!alloc || (pgtab = (pte_t*)kalloc()) == 0)
1664       return 0;
1665     // Make sure all those PTE_P bits are zero.
1666     memset(pgtab, 0, PGSIZE);
…
1670     *pde = v2p(pgtab) | PTE_P | PTE_W | PTE_U;
1671   }
1672   return &pgtab[PTX(va)];
1673 }
```

```
1736 pde_t*
1737 setupkvm(void)
1738 {
1739   pde_t *pgdir;
1740   struct kmap *k;
1741
1742   if((pgdir = (pde_t*)kalloc()) == 0)
1743     return 0;
1744   memset(pgdir, 0, PGSIZE);
...
1747   for(k = kmap; k < &kmap[NELEM(kmap)]; k++)
1748     if(mappages(pgdir, k->virt, k->phys_end - k->phys_start,
1749               (uint)k->phys_start, k->perm) < 0)
1750       return 0;
1751   return pgdir;
1752 }
```

# kvmalloc()

```
1757 kvmalloc(void)
1758 {
1759    kpgdir = setupkvm();
1760    switchkvm();
1761 }
```

# Switch to the new page table

```
1765 void

1766 switchkvm(void)

1767 {

1768    lcr3(v2p(kpgdir));

1769 }
```

```
1217 main(void)
1218 {
1219   kinit1(end, P2V(4*1024*1024)); // phys page allocator
1220   kvmalloc(); // kernel page table
1221   mpinit(); // collect info about this machine
1222   lapicinit();
1223   seginit(); // set up segments
1224   cprintf("\ncpu%d: starting xv6\n\n", cpu->id);
1225   picinit(); // interrupt controller
1226   ioapicinit(); // another interrupt controller
1227   consoleinit(); // I/O devices & their interrupts
1228   uartinit(); // serial port
1229   pinit(); // process table
1230   tvinit(); // trap vectors
```

```
1616 seginit(void)
1617 {
1618   struct cpu *c;
...
1624   c = &cpus[cpunum()];
1625   c->gdt[SEG_KCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, 0);
1626   c->gdt[SEG_KDATA] = SEG(STA_W, 0, 0xffffffff, 0);
1627   c->gdt[SEG_UCODE] = SEG(STA_X|STA_R, 0, 0xffffffff, DPL_USER);
1628   c->gdt[SEG_UDATA] = SEG(STA_W, 0, 0xffffffff, DPL_USER);
1629
1630   // Map cpu, and curproc
1631   c->gdt[SEG_KCPU] = SEG(STA_W, &c->cpu, 8, 0);
1632
1633   lgdt(c->gdt, sizeof(c->gdt));
1634   loadgs(SEG_KCPU << 3);
1635
1636   // Initialize cpu-local storage.
1637   cpu = c;
1638   proc = 0;
1639 }
```

```
1217 main(void)
1218 {
1219   kinit1(end, P2V(4*1024*1024)); // phys page allocator
1220   kvmalloc(); // kernel page table
1221   mpinit(); // collect info about this machine
1222   lapicinit();
1223   seginit(); // set up segments
1224   cprintf("\ncpu%d: starting xv6\n\n", cpu->id);
1225   picinit(); // interrupt controller
1226   ioapicinit(); // another interrupt controller
1227   consoleinit(); // I/O devices & their interrupts
1228   uartinit(); // serial port
1229   pinit(); // process table
1230   tvinit(); // trap vectors
```

# Interrupt descriptor table

```
3066 void

3067 tvinit(void)

3068 {

3069   int i;

3070

3071   for(i = 0; i < 256; i++)

3072     SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);

3073   SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3,

                vectors[T_SYSCALL], DPL_USER);
...

3076 }
```

# Thank you!