# CS5460/6460: Operating Systems
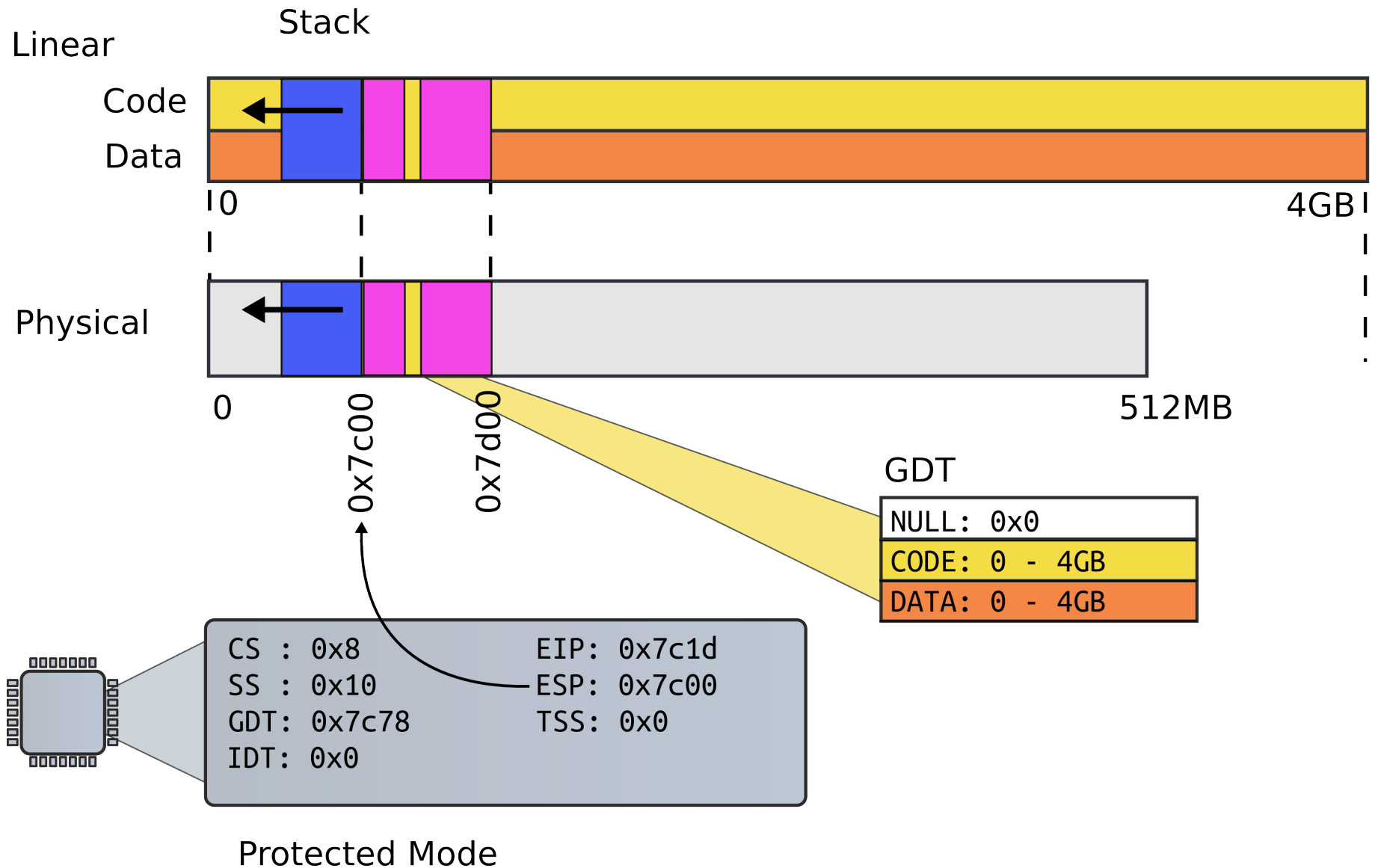
# Lecture 9: Finishing system boot, and system init

Anton Burtsev
January, 2014

# First stack

Linear

Stack

Code

Data

0                                                                4GB

Physical

0                                          512MB

0x7c00

0x7d00

GDT

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

CS : 0x8          EIP: 0x7c1d
SS : 0x10         ESP: 0x7c00
GDT: 0x7c78       TSS: 0x0
IDT: 0x0

Protected Mode

# Invoke first C function

9166 movl $start, %esp

9167 call bootmain

xv6/bootasm.S

# bootmain(): read kernel from disk

```
9216 void
9217 bootmain(void)
9218 {
9219     struct elfhdr *elf;
9220     struct proghdr *ph, *eph;
9221     void (*entry)(void);
9222     uchar* pa;
9223
9224     elf = (struct elfhdr*)0x10000; // scratch space
9225
9226     // Read 1st page off disk
9227     readseg((uchar*)elf, 4096, 0);
9228
9229     // Is this an ELF executable?
9230     if(elf->magic != ELF_MAGIC)
9231         return; // let bootasm.S handle error
9232
```

xv6/bootmain.c

# bootmain(): read kernel from disk

```c
9232
9233    // Load each program segment (ignores ph flags).
9234    ph = (struct proghdr*)((uchar*)elf + elf->phoff);
9235    eph = ph + elf->phnum;
9236    for(; ph < eph; ph++){
9237        pa = (uchar*)ph->paddr;
9238        readseg(pa, ph->filesz, ph->off);
9239        if(ph->memsz > ph->filesz)
9240            stosb(pa + ph->filesz, 0, ph->memsz - ph->filesz);
9241    }
9242
9243    // Call the entry point from the ELF header.
9244    // Does not return!
9245    entry = (void(*)(void))(elf->entry);
9246    entry();
9247 }
```
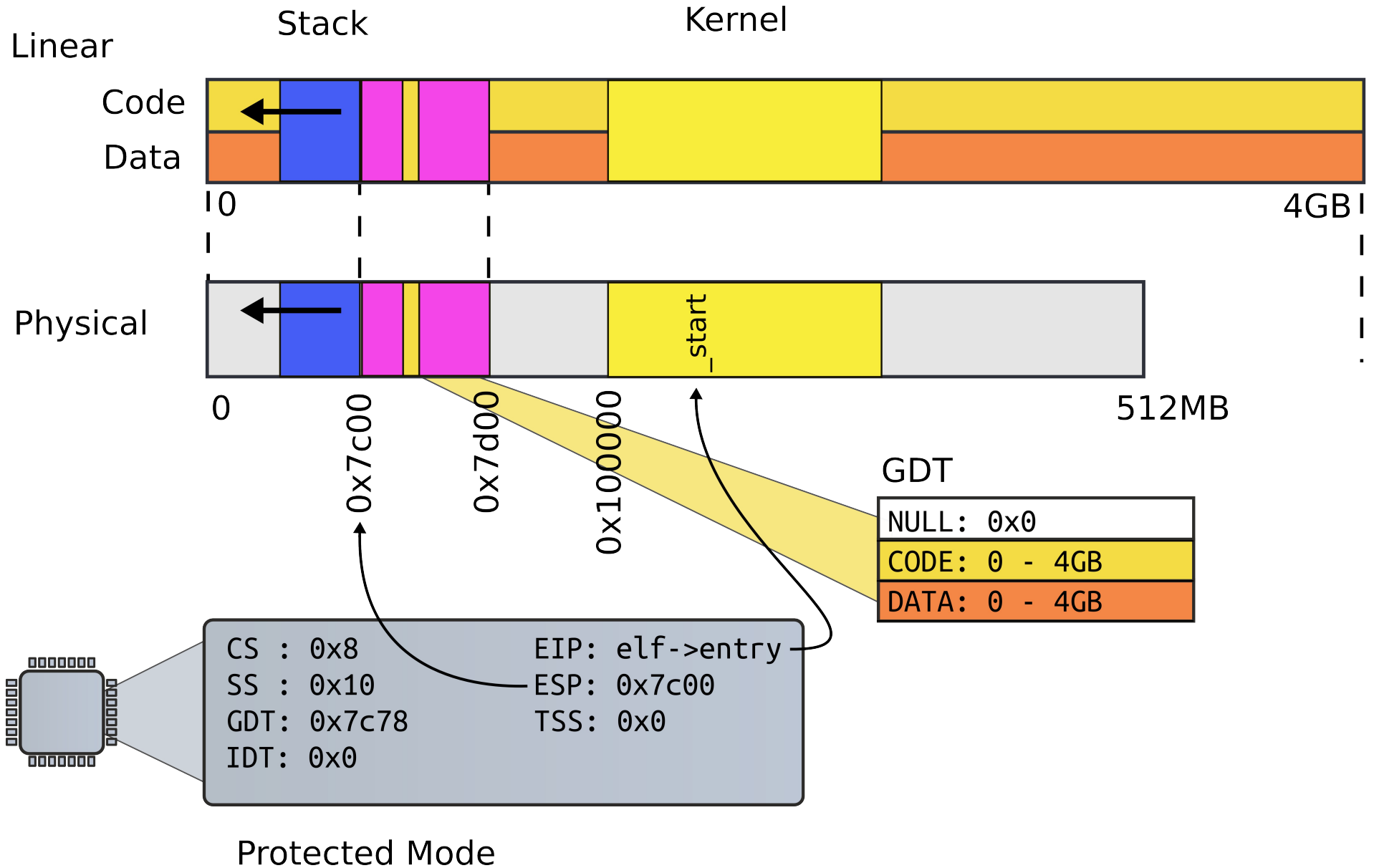
xv6/bootmain.c

# Kernel

```
1039 .globl entry

1136 # By convention, the _start symbol specifies the ELF entry point.
1137 # Since we haven't set up virtual memory yet, our entry point is
1138 # the physical address of 'entry'.
1139 .globl _start
1140 _start = V2P_WO(entry)
1141
1142 # Entering xv6 on boot processor, with paging off.
1143 .globl entry
1144 entry:
1145 # Turn on page size extension for 4Mbyte pages
1146     movl %cr4, %eax
1147     orl $(CR4_PSE), %eax
1148     movl %eax, %cr4
```
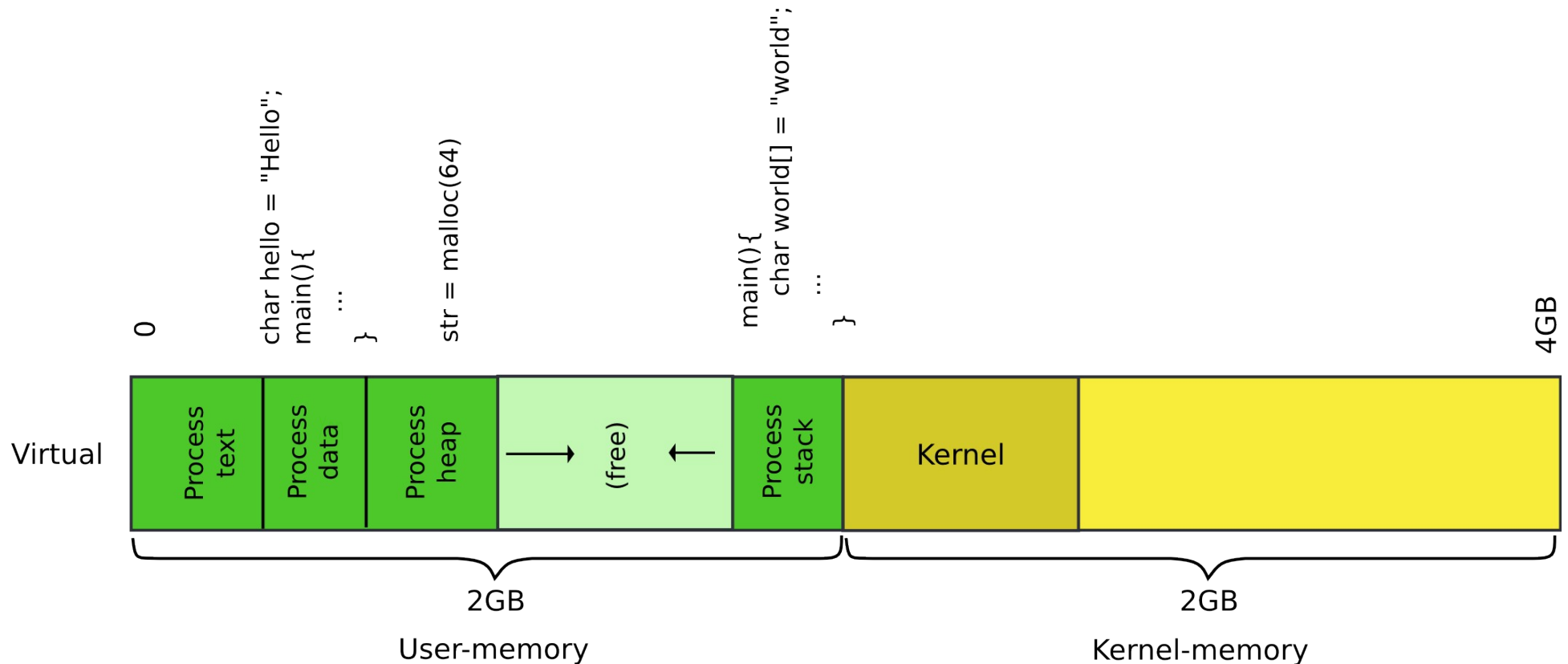
*xv6/entry.S*

entry(): kernel ELF entry

# Set up page directory

```
1149 # Set page directory

1150 movl $(V2P_WO(entrypgdir)), %eax

1151 movl %eax, %cr3
```

xv6/entry.S

## Linear Address

| 31 | 22 | 21 | 0 |
|---|---|---|---|
| Directory | | Offset | |

/ 22

## 4-MByte Page

/ 10

## Page Directory

Physical Address

PDE with PS=1

/ 18

/ 32

CR3

# Our goal: 2GB/2GB address space

Linear

Stack

Kernel

Code

Data

0

4MB

4GB

Physical

entry

0

0x7c00

0x7d00

0x100000

512MB

Protected Mode

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

GDT

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

# First page table

- Two 4MB entries (large pages)
- Entry #0
  - 0x0 – 4MB → 0x0:0x400000
- Entry #960
  - 0x0 – 4MB → 0x8000000:0x80400000

```
1406 // The boot page table used in entry.S and entryother.S.
1407 // Page directories (and page tables) must start on page
         boundaries,
1408 // hence the __aligned__ attribute.
1409 // PTE_PS in a page directory entry enables 4Mbyte
pages.
1410
1411 __attribute__((__aligned__(PGSIZE)))
1412 pde_t entrypgdir[NPDENTRIES] = {
1413   // Map VA's [0, 4MB) to PA's [0, 4MB)
1414   [0] = (0) | PTE_P | PTE_W | PTE_PS,
1415   // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
1416   [KERNBASE>>PDXSHIFT] = (0) | PTE_P | PTE_W | PTE_PS,
1417 };
```

First page table

```
1406 // The boot page table used in entry.S and entryother.S.
1407 // Page directories (and page tables) must start on page
         boundaries,
1408 // hence the __aligned__ attribute.
1409 // PTE_PS in a page directory entry enables 4Mbyte
pages.
1410
1411 __attribute__((__aligned__(PGSIZE)))
1412 pde_t entrypgdir[NPDENTRIES] = {
1413   // Map VA's [0, 4MB) to PA's [0, 4MB)
1414   [0] = (0) | PTE_P | PTE_W | PTE_PS,
1415   // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
1416   [KERNBASE>>PDXSHIFT] = (0) | PTE_P | PTE_W | PTE_PS,
1417 };
```

First page table

```
1406 // The boot page table used in entry.S and entryother.S.
1407 // Page directories (and page tables) must start on page
         boundaries,
1408 // hence the __aligned__ attribute.
1409 // PTE_PS in a page directory entry enables 4Mbyte
pages.
1410
1411 __attribute__((__aligned__(PGSIZE)))
1412 pde_t entrypgdir[NPDENTRIES] = {
1413   // Map VA's [0, 4MB) to PA's [0, 4MB)
1414   [0] = (0) | PTE_P | PTE_W | PTE_PS,
1415   // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
1416   [KERNBASE>>PDXSHIFT] = (0) | PTE_P | PTE_W | PTE_PS,
1417 };
```

First page table

```
1406 // The boot page table used in entry.S and entryother.S.
1407 // Page directories (and page tables) must start on page
        boundaries,
1408 // hence the __aligned__ attribute.
1409 // PTE_PS in a page directory entry enables 4Mbyte
pages.
1410
1411 __attribute__((__aligned__(PGSIZE)))
1412 pde_t entrypgdir[NPDENTRIES] = {
1413   // Map VA's [0, 4MB) to PA's [0, 4MB)
1414   [0] = (0) | PTE_P | PTE_W | PTE_PS,
1415   // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
1416   [KERNBASE>>PDXSHIFT] = (0) | PTE_P | PTE_W | PTE_PS,
1417 };
```

First page table

```
1406 // The boot page table used in entry.S and entryother.S.
1407 // Page directories (and page tables) must start on page
          boundaries,
1408 // hence the __aligned__ attribute.
1409 // PTE_PS in a page directory entry enables 4Mbyte
pages.
1410
1411 __attribute__((__aligned__(PGSIZE)))
1412 pde_t entrypgdir[NPDENTRIES] = {
1413   // Map VA's [0, 4MB) to PA's [0, 4MB)
1414   [0] = (0) | PTE_P | PTE_W | PTE_PS,
1415   // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
1416   [KERNBASE>>PDXSHIFT] = (0) | PTE_P | PTE_W | PTE_PS,
1417 };
```
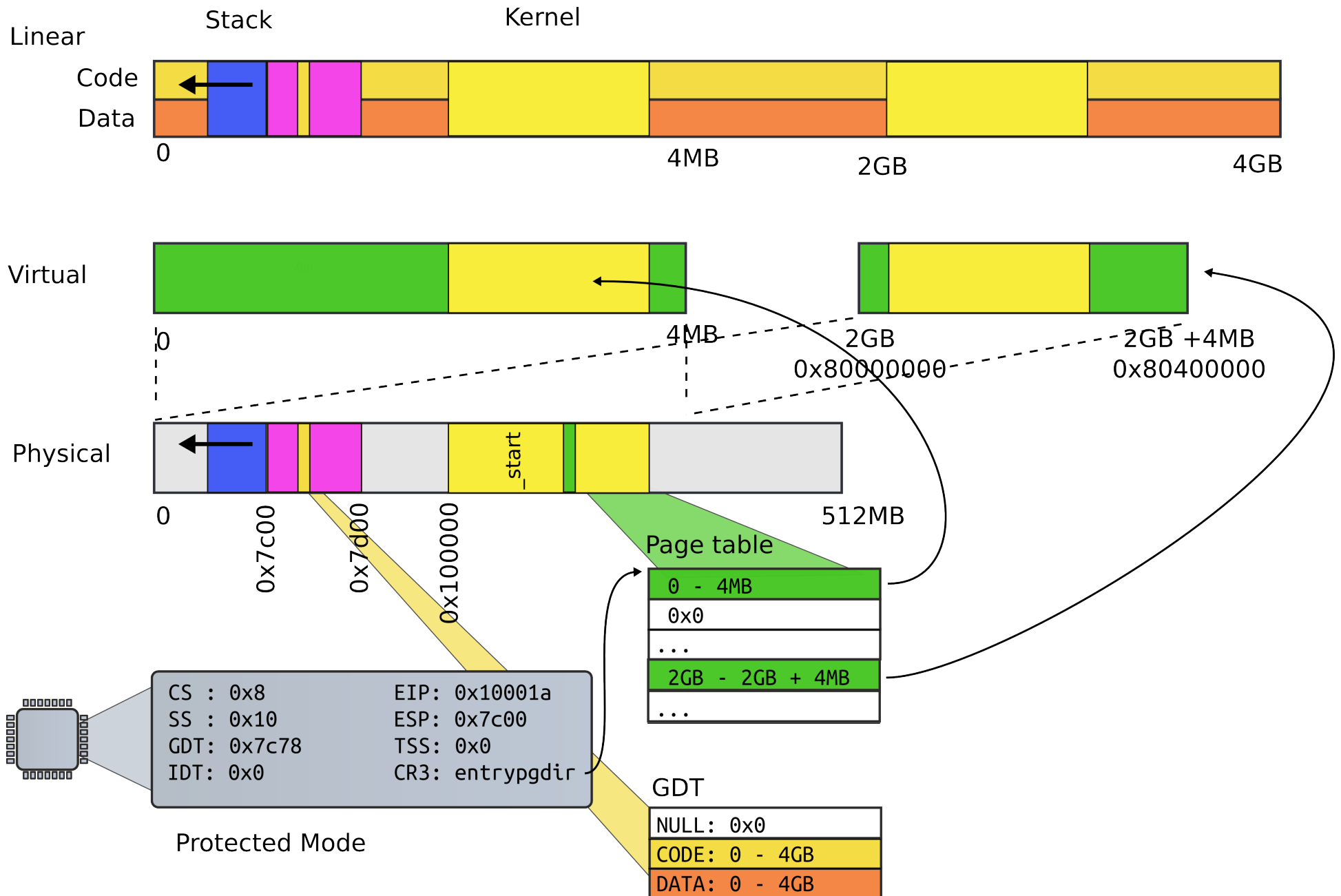
First page table

# First page table (cont)

```
0870 // Page directory and page table constants.

0871 #define NPDENTRIES 1024
```

# First page table

Linear

Stack        Kernel

Code  ←————

Data

0            4MB        2GB            4GB

Virtual

0            4MB        2GB            2GB +4MB
                        0x80000000     0x80400000

Physical  ←————

_start

0    0x7c00  0x7d00  0x100000    512MB

Page table

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

```
CS : 0x8          EIP: 0x10001a
SS : 0x10         ESP: 0x7c00
GDT: 0x7c78       TSS: 0x0
IDT: 0x0          CR3: entrypgdir
```

Protected Mode

GDT

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

# Turn on paging

```
1152 # Turn on paging.
1153 movl %cr0, %eax
1154 orl $(CR0_PG|CR0_WP), %eax
1155 movl %eax, %cr0
```

# High address stack (4K)

```
1157 # Set up the stack pointer.
1158 movl $(stack + KSTACKSIZE), %esp
1159

...

1167 .comm stack, KSTACKSIZE


0151 #define KSTACKSIZE 4096 // size of

                        per-process kernel stack
```
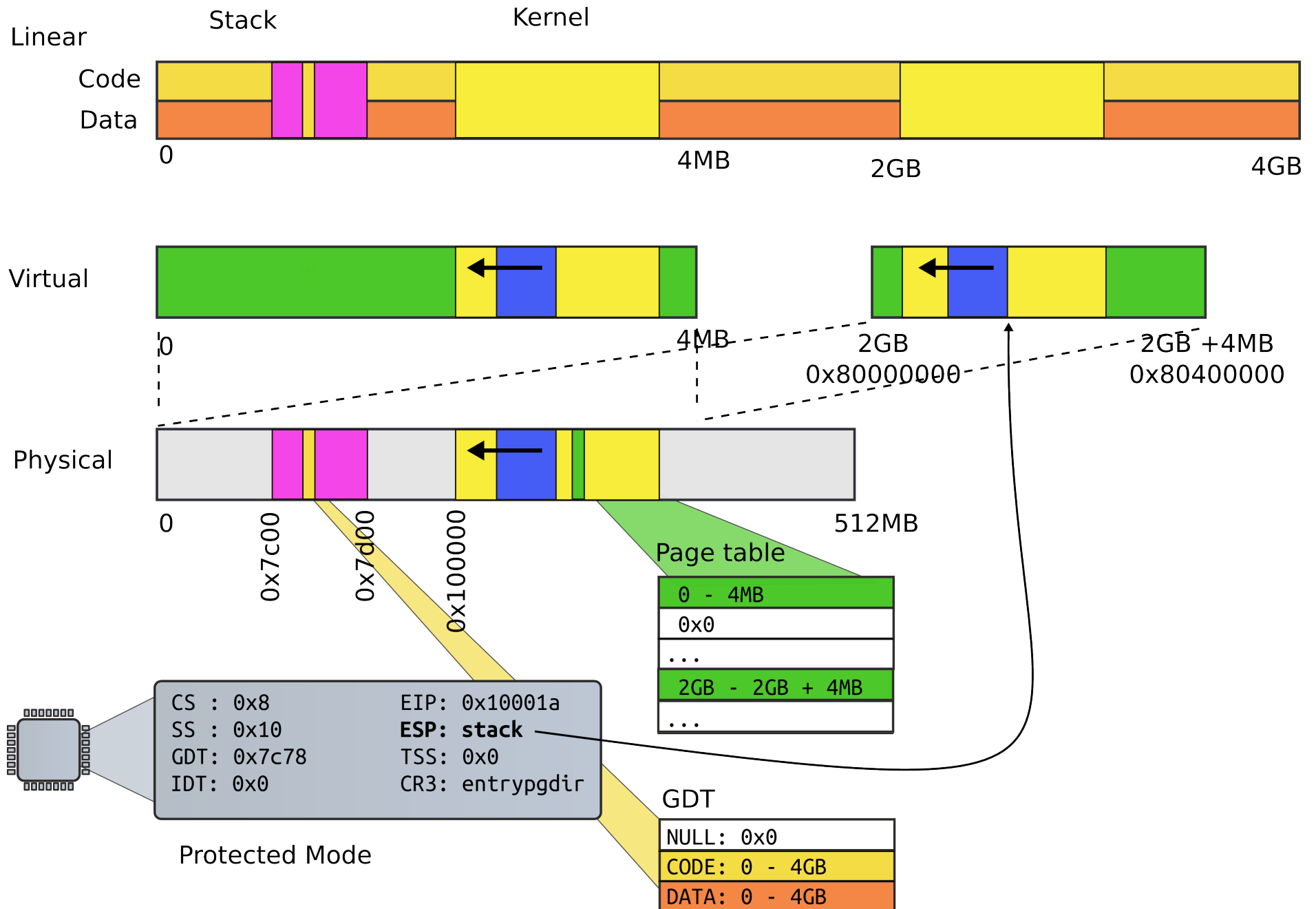
# High address stack (4K)

**Linear**

Stack　　　　　　　Kernel

Code

Data

0　　　　　　　　　　　　　　　　4MB　　　　2GB　　　　　　　　4GB

**Virtual**

0　　　　　　　　　　　　4MB　　　2GB　　　　　　2GB +4MB
　　　　　　　　　　　　　　　0x80000000　　0x80400000

**Physical**

0　0x7c00　0x7d00　0x100000　　　　512MB

Page table

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

```
CS : 0x8        EIP: 0x10001a
SS : 0x10       ESP: stack
GDT: 0x7c78     TSS: 0x0
IDT: 0x0        CR3: entrypgdir
```

Protected Mode

GDT

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

# Jump to main()

```
1160 # Jump to main(), and switch to executing at
1161 # high addresses. The indirect call is
         needed because
1162 # the assembler produces a PC-relative
         instruction
1163 # for a direct jump.
1164 mov $main, %eax
1165 jmp *%eax
1166
```

# 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 }
```

# Recap of the boot sequence

- 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()

# Conclusion

- We've booted
  - We're running in main()

- Next time:
  - 2-level page tables
  - Process and kernel address space

# Thank you!

# Thank you!