# 143A: Principles of Operating Systems
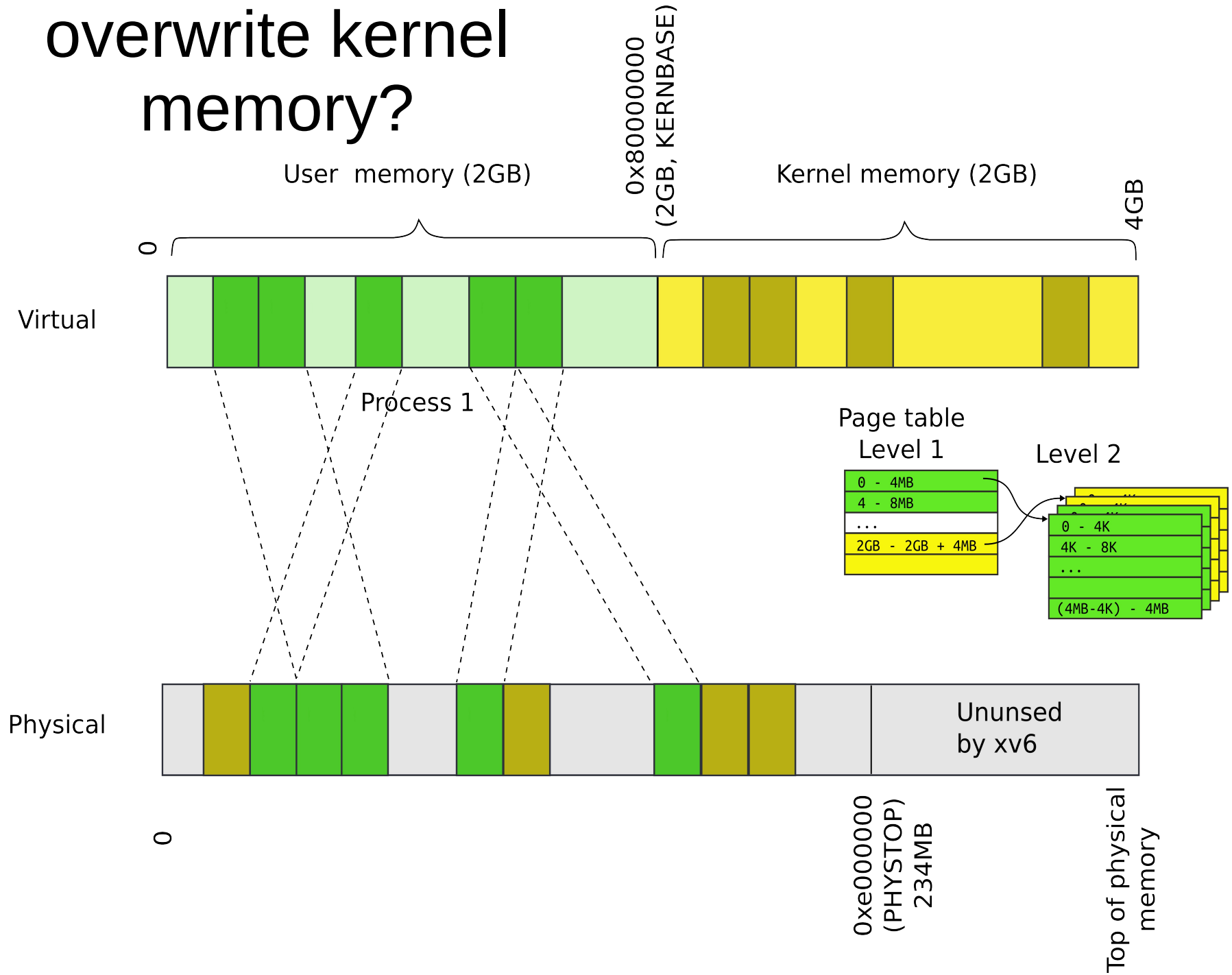
# Lecture 09: First process
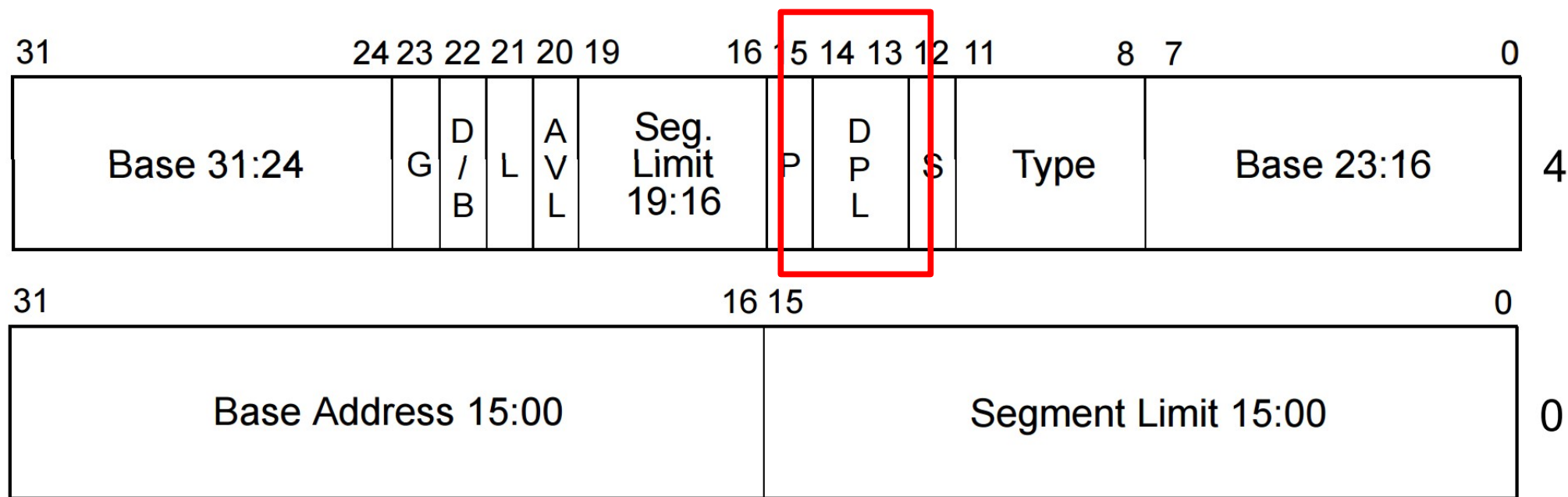
Anton Burtsev
January, 2017

# Recap: Can a process overwrite kernel memory?

User memory (2GB)

0x80000000 (2GB, KERNBASE)

Kernel memory (2GB)

4GB

Virtual

0

Process 1

Page table
Level 1

Level 2

| 0 - 4MB |
| 4 - 8MB |
| ... |
| 2GB - 2GB + 4MB |

| 0 - 4K |
| 4K - 8K |
| ... |
| (4MB-4K) - 4MB |

Physical

Ununsed
by xv6

0

0xe000000 (PHYSTOP) 234MB

Top of physical memory

# Privilege levels

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

| 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 |

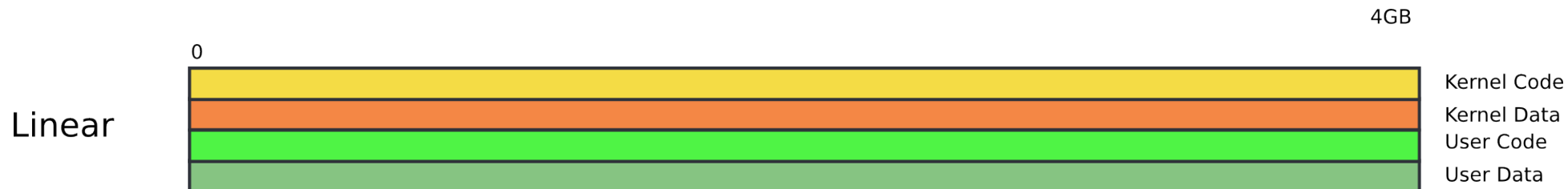| 31 | 16 15 | 0 | |
|---|---|---|---|
| Base Address 15:00 | Segment Limit 15:00 | | 0 |

# 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>
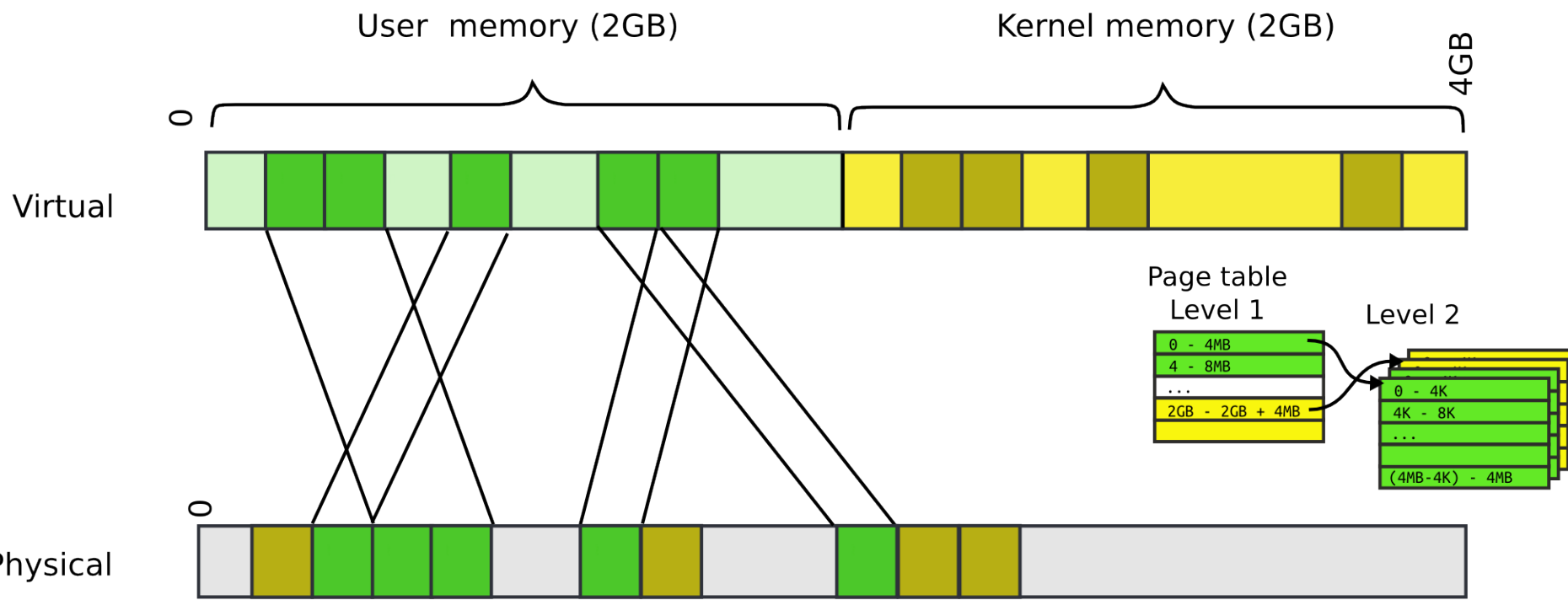            - E.g. reload a page table by changing CR3

# 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
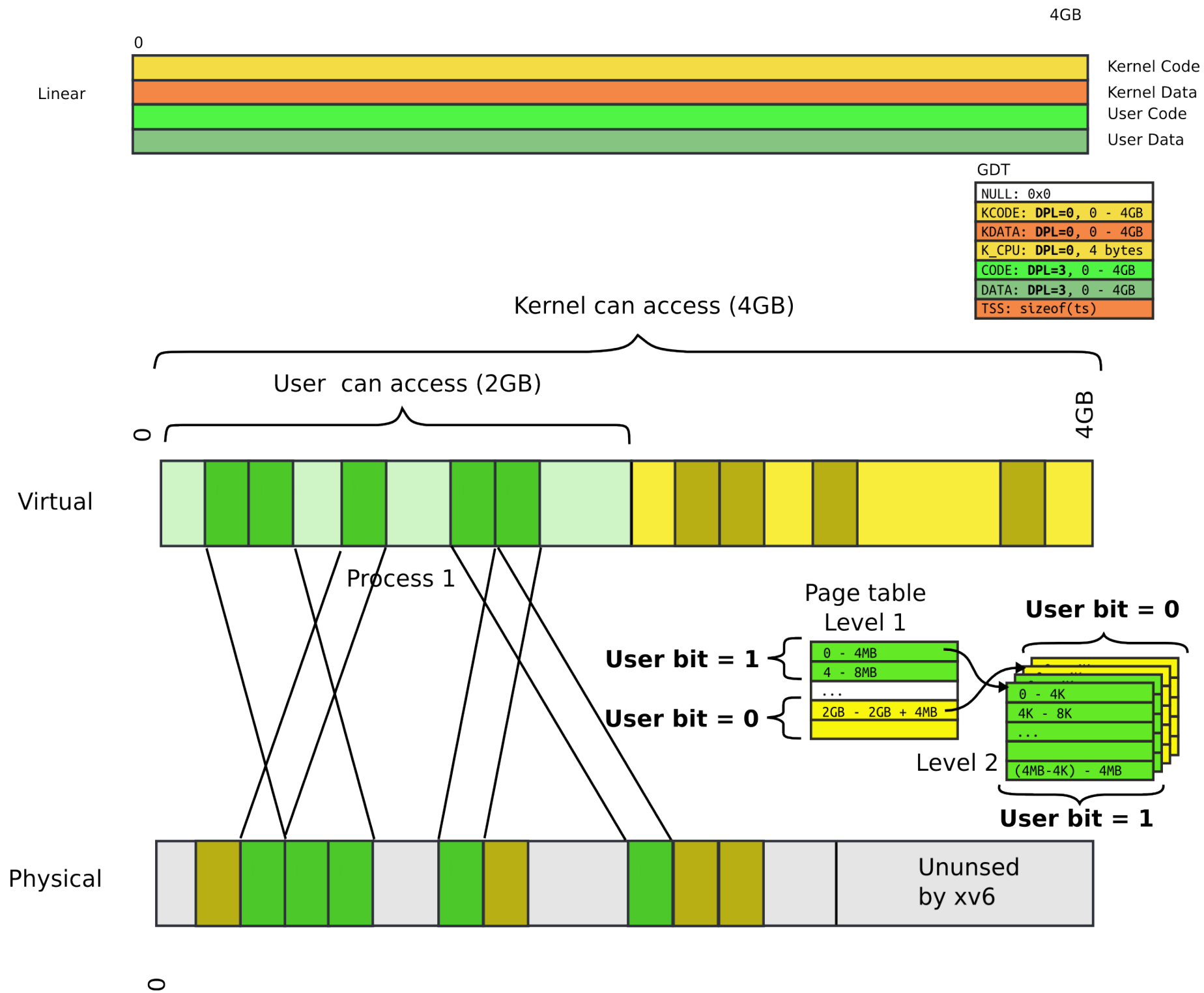- But then... how the kernel is protected?
  - Page tables

# Linear

0                                                          4GB

Kernel Code

Kernel Data

User Code

User Data

**GDT**

| |
|---|
| NULL: 0x0 |
| KCODE: **DPL=0**, 0 - 4GB |
| KDATA: **DPL=0**, 0 - 4GB |
| K_CPU: **DPL=0**, 4 bytes |
| CODE: **DPL=3**, 0 - 4GB |
| DATA: **DPL=3**, 0 - 4GB |
| TSS: sizeof(ts) |

User memory (2GB)                  Kernel memory (2GB)

0                                                               4GB

**Virtual**

Page table
Level 1            Level 2

| |
|---|
| 0 - 4MB |
| 4 - 8MB |
| ... |
| 2GB - 2GB + 4MB |

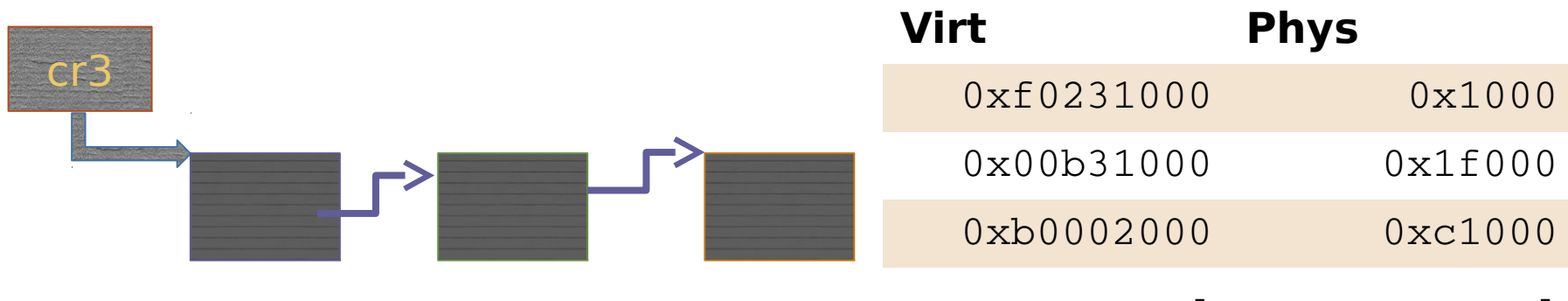| |
|---|
| 0 - 4K |
| 4K - 8K |
| ... |
| (4MB-4K) - 4MB |

0

**Physical**

# 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

4GB

0

Linear

Kernel Code

Kernel Data

User Code

User Data

GDT

NULL: 0x0

KCODE: **DPL=0**, 0 - 4GB

KDATA: **DPL=0**, 0 - 4GB

K_CPU: **DPL=0**, 4 bytes

CODE: **DPL=3**, 0 - 4GB

DATA: **DPL=3**, 0 - 4GB

TSS: sizeof(ts)

Kernel can access (4GB)

User can access (2GB)

0

4GB

Virtual

Process 1

Page table
Level 1

**User bit = 0**

**User bit = 1**

0 - 4MB

4 - 8MB

...

2GB - 2GB + 4MB

**User bit = 0**

0 - 4K

4K - 8K

...

(4MB-4K) - 4MB

Level 2

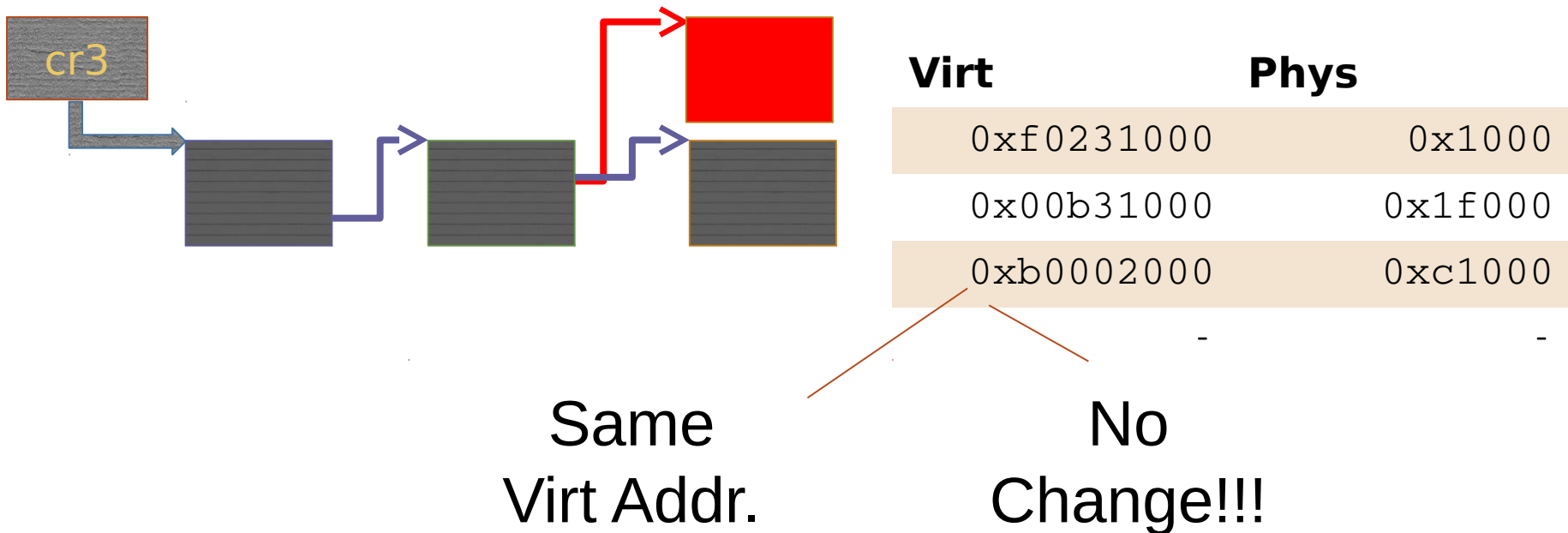**User bit = 1**

Physical

Ununsed
by xv6

0

# 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

cr3

| Virt | Phys |
|---|---|
| 0xf0231000 | 0x1000 |
| 0x00b31000 | 0x1f000 |
| 0xb0002000 | 0xc1000 |
| - | - |

# 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

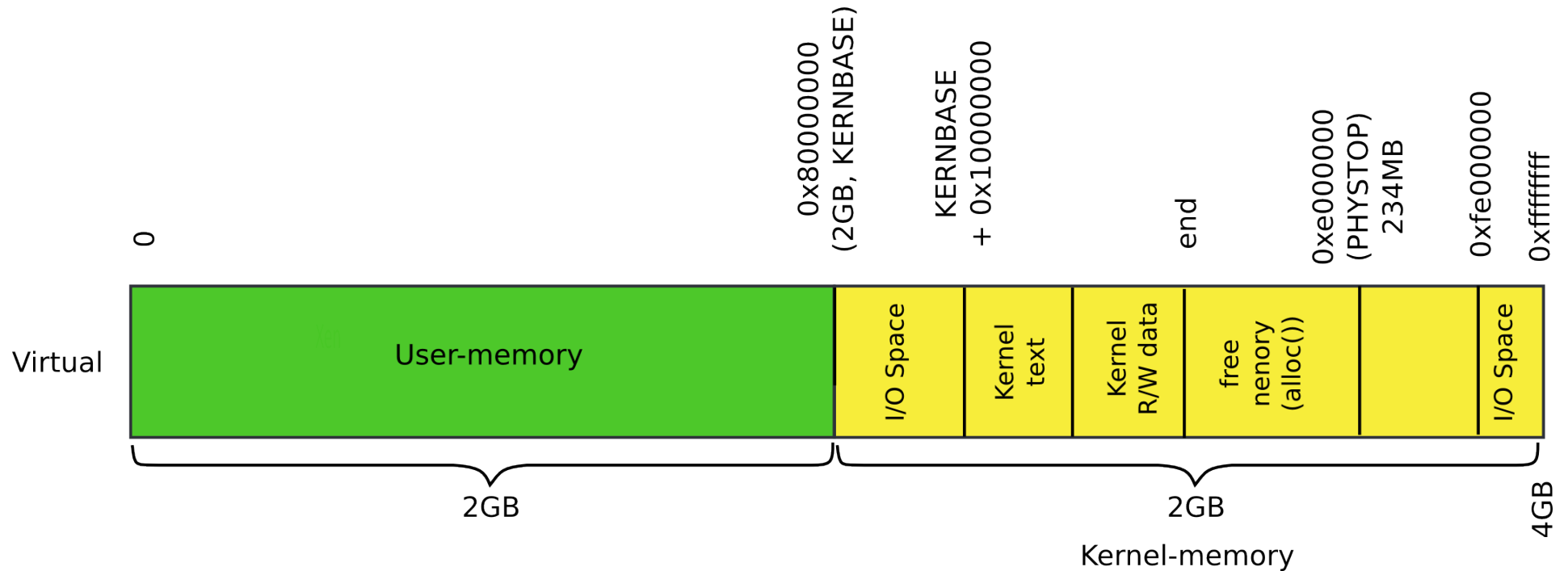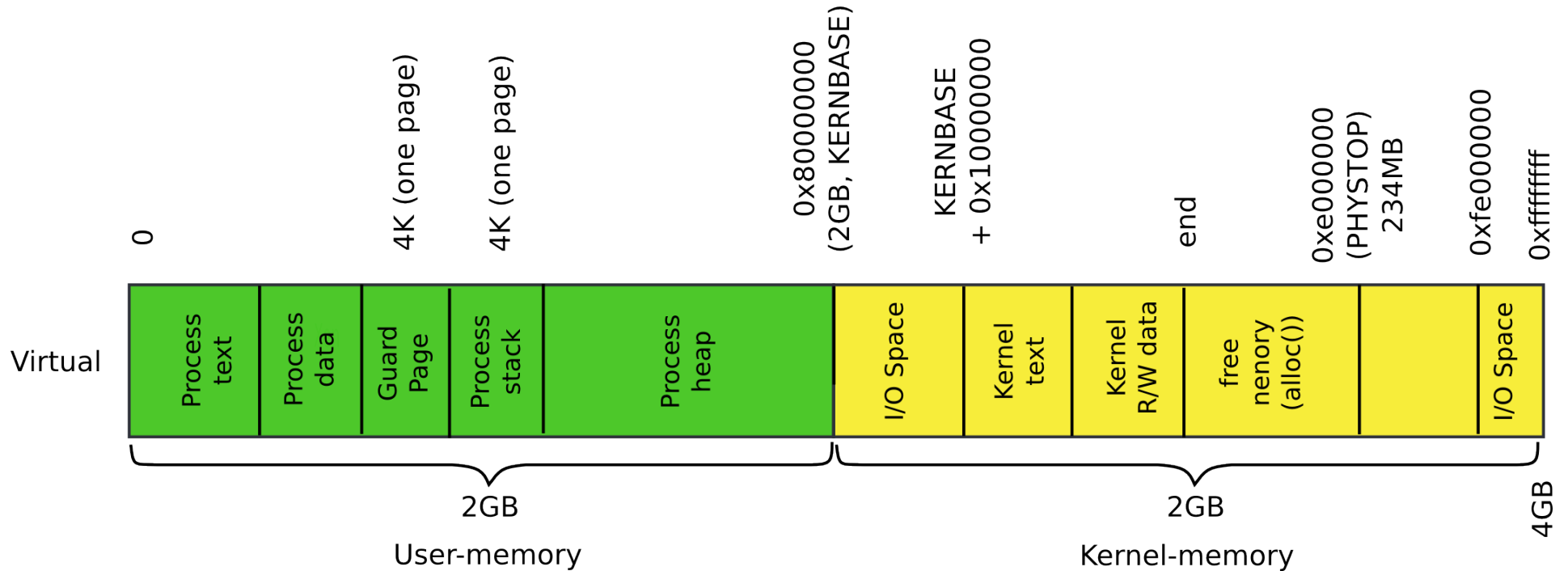| Virt | Phys |
|---|---|
| 0xf0231000 | 0x1000 |
| 0x00b31000 | 0x1f000 |
| 0xb0002000 | 0xc1000 |
| - | - |

Same
Virt Addr.

No
Change!!!

# 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)

# Creating Processes

# Recap: kernel memory layout



Virtual

| User-memory | I/O Space | Kernel text | Kernel R/W data | free memory (alloc()) | | I/O Space |

0

0x80000000 (2GB, KERNBASE)

KERNBASE + 0x10000000

end

0xe000000 (PHYSTOP) 234MB

0xfe000000

0xffffffff

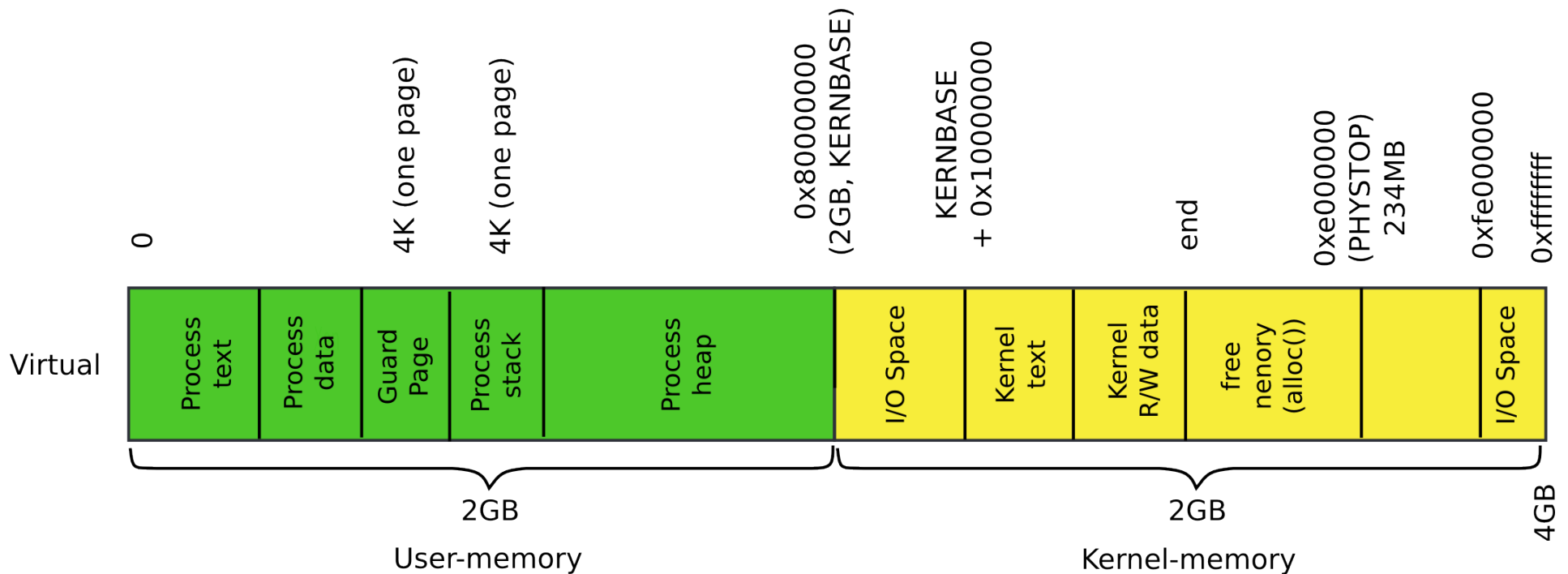2GB

2GB

Kernel-memory

4GB

# Process memory layout

# How does kernel creates new processes?

# How does kernel creates new processes?
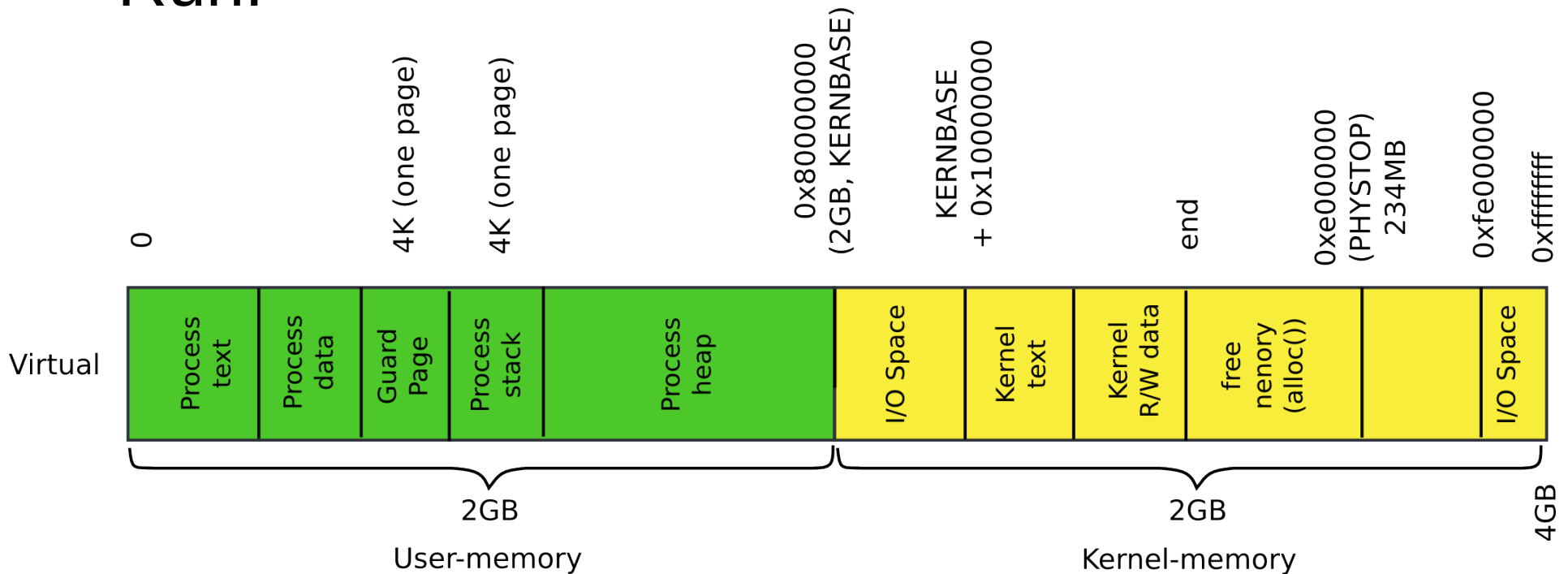
- Exec
    - exec("/bin/ls", argv);

# exec(): high-level outline

- We want to create the following layout
- What shall we do?

# exec(): high-level outline

- Load program from disk
- Create user-stack
- Run!

# exec(): locate inode

```
6309 int
6310 exec(char *path, char **argv)
6311 {
...
6321   if((ip = namei(path)) == 0){
6322     end_op();
6323     return -1;
6324   }
6328   // Check ELF header
6329   if(readi(ip, (char*)&elf, 0, sizeof(elf)) <
                                      sizeof(elf))
6330     goto bad;
6331   if(elf.magic != ELF_MAGIC)
6332     goto bad;
```

# exec(): check ELF header

```
6309 int
6310 exec(char *path, char **argv)
6311 {
...
6321   if((ip = namei(path)) == 0){
6322     end_op();
6323     return -1;
6324   }
6328   // Check ELF header
6329   if(readi(ip, (char*)&elf, 0, sizeof(elf)) <
                              sizeof(elf))
6330     goto bad;
6331   if(elf.magic != ELF_MAGIC)
6332     goto bad;
```
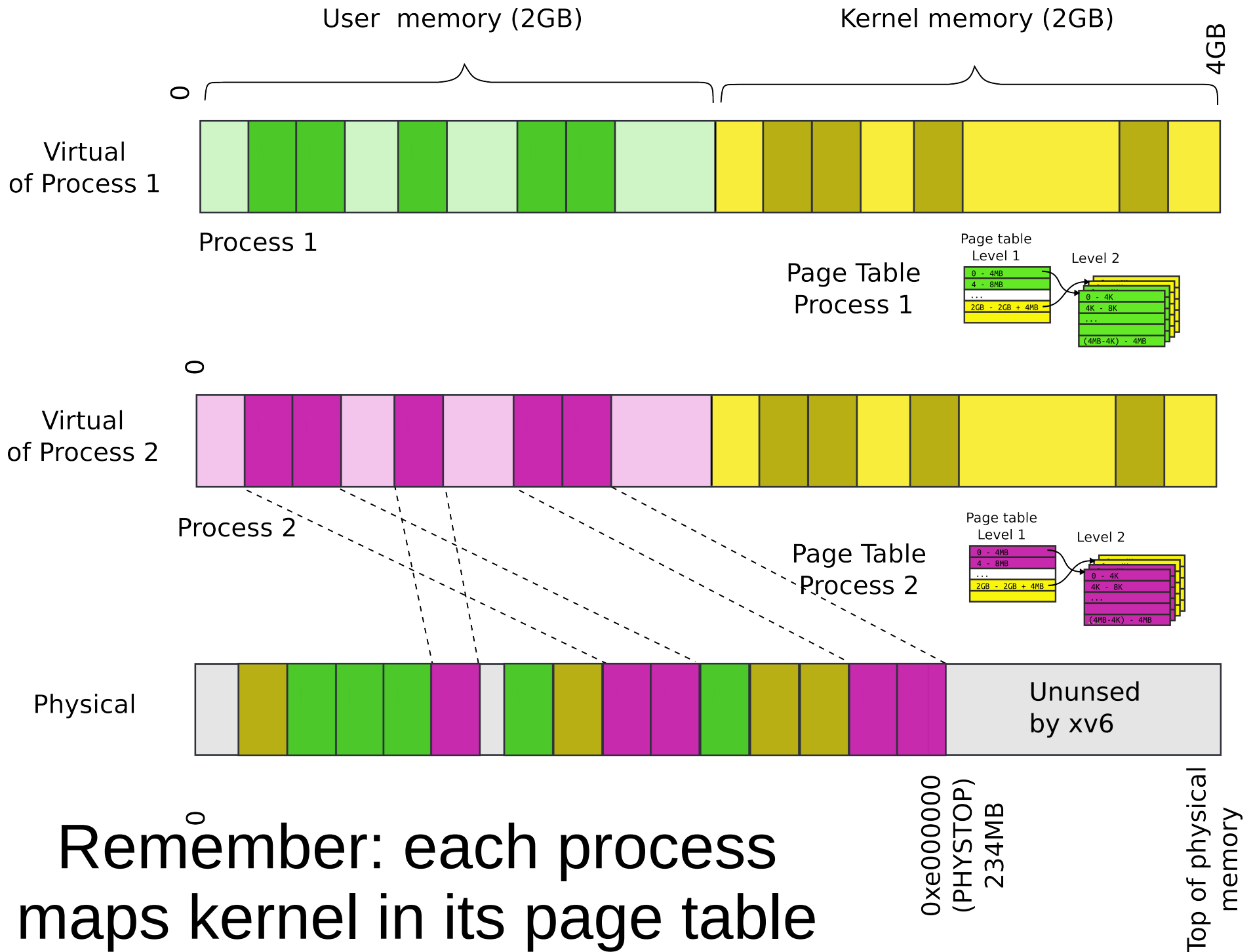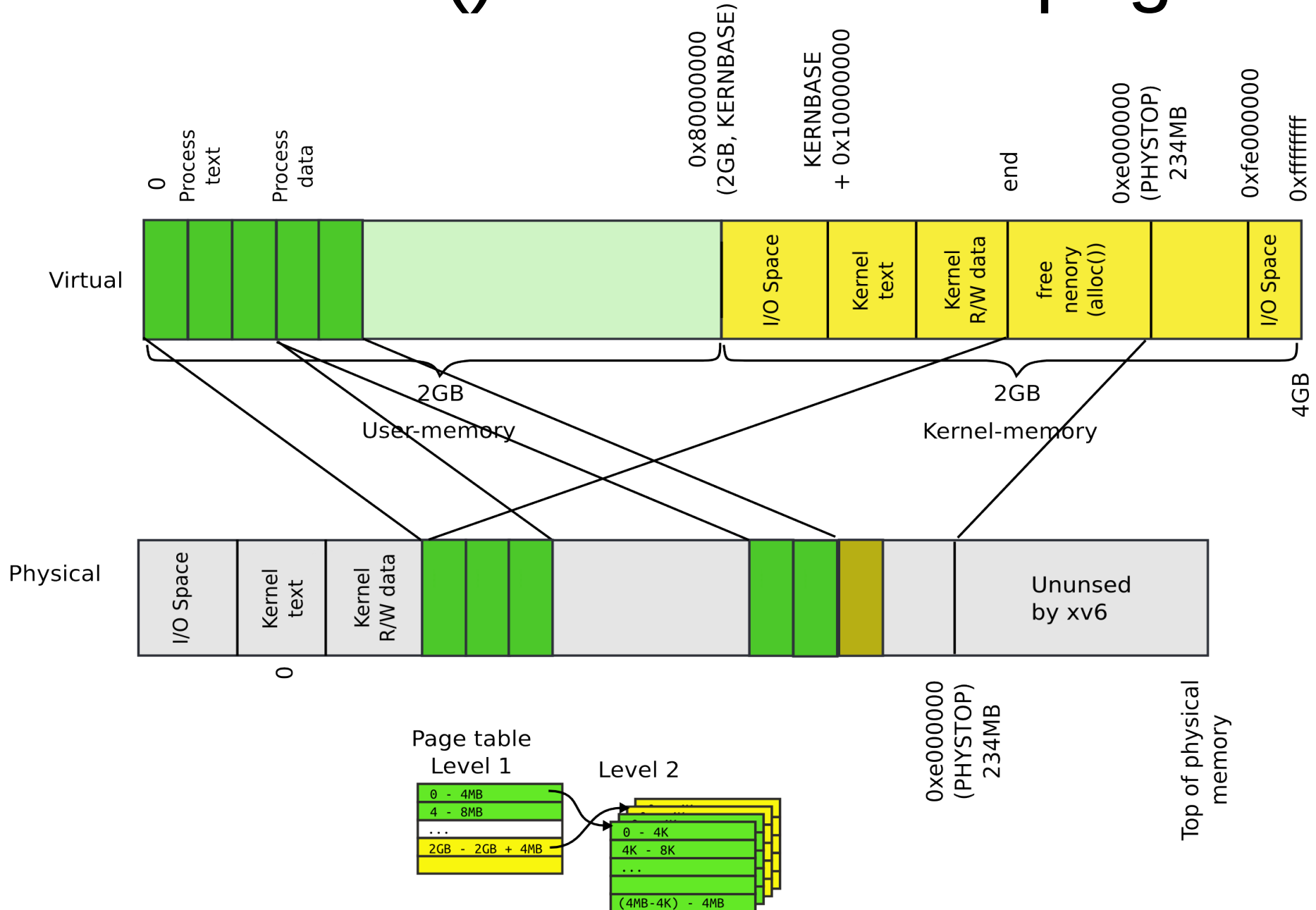
# Setup kernel address space()

6333

6334    if((pgdir = **setupkvm**()) == 0)

6335      goto bad;

6336

User memory (2GB)    Kernel memory (2GB)

0                                    4GB

Virtual of Process 1

Process 1

Page Table Process 1

Page table Level 1
0 - 4MB
4 - 8MB
...
2GB - 2GB + 4MB

Level 2
0 - 4K
4K - 8K
...
(4MB-4K) - 4MB

0

Virtual of Process 2

Process 2

Page Table Process 2

Page table Level 1
0 - 4MB
4 - 8MB
...
2GB - 2GB + 4MB

Level 2
0 - 4K
4K - 8K
...
(4MB-4K) - 4MB

Physical

Ununsed by xv6

0

0xe000000 (PHYSTOP) 234MB

Top of physical memory

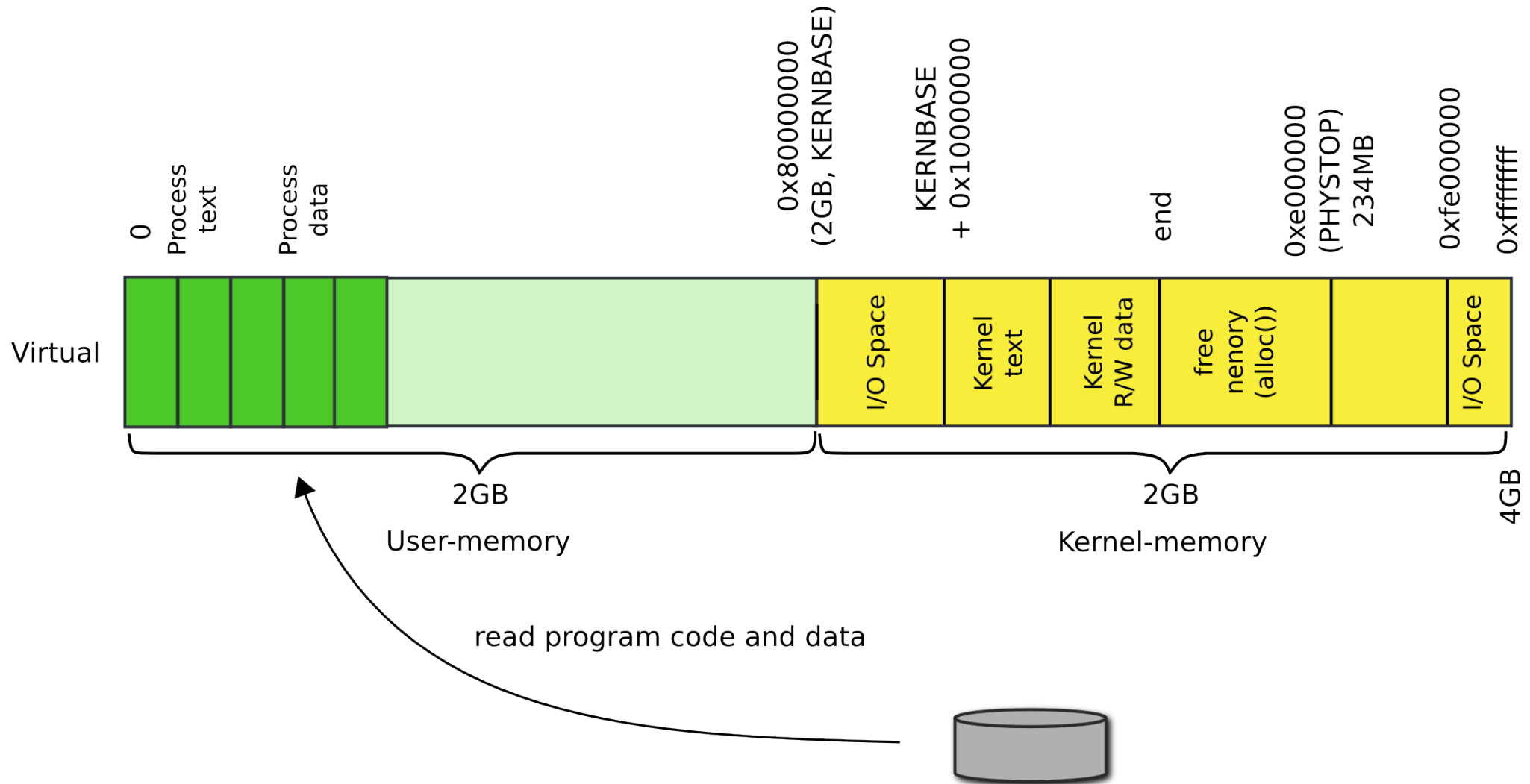# Remember: each process maps kernel in its page table

# allocuvm(): allocate user pages
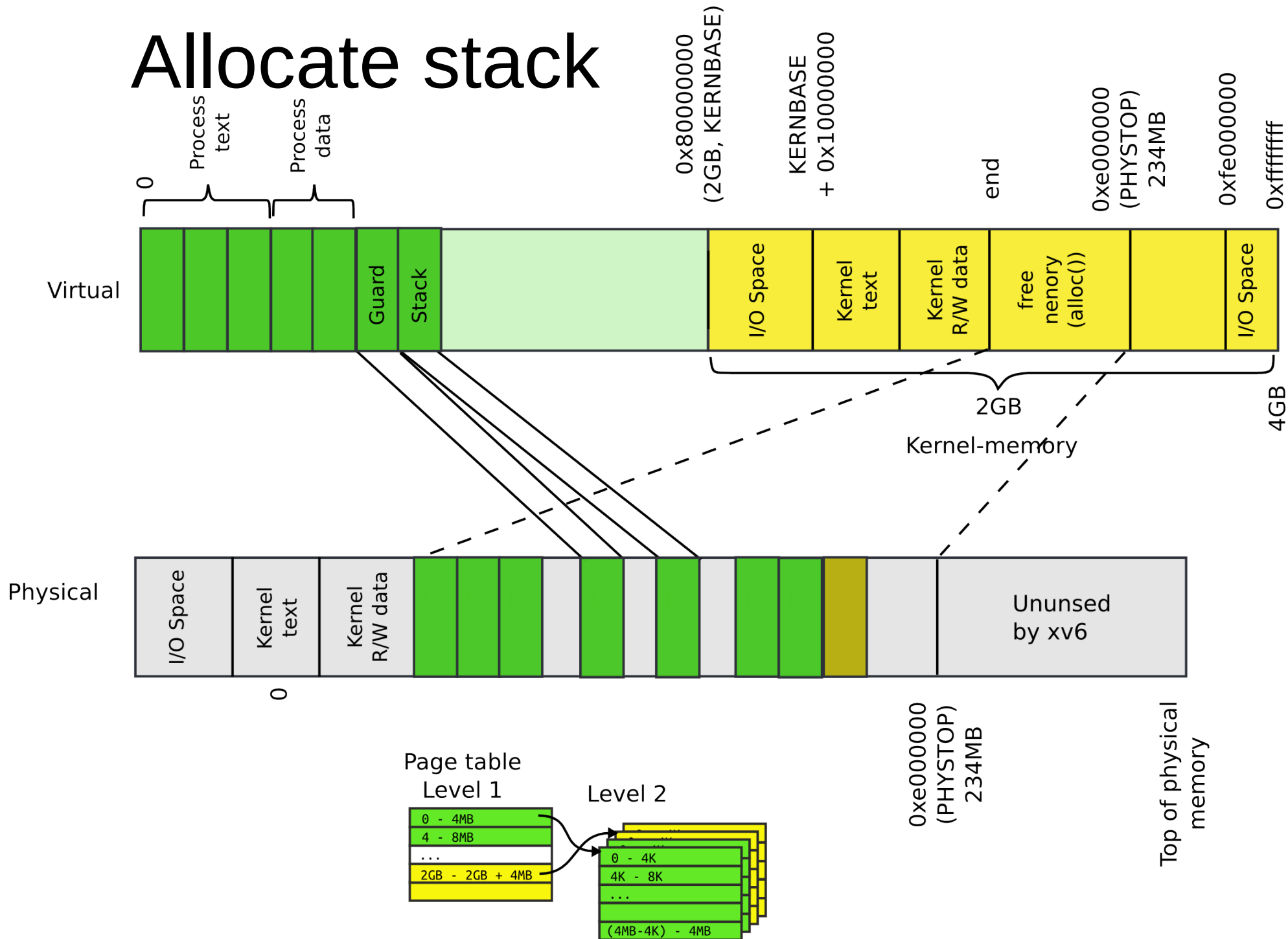
# loaduvm(): read program from disk

# exec(): allocate process' stack

- Allocate two pages
  - One will be stack
  - Mark another one as inaccessible

```
6361   sz = PGROUNDUP(sz);

6362   if((sz = allocuvm(pgdir, sz, sz + 2*PGSIZE)) == 0)

6363     goto bad;

6364   clearpteu(pgdir, (char*)(sz - 2*PGSIZE));

6365   sp = sz;
```

# Allocate stack



Virtual

0 · Process text · Process data · Guard · Stack · 0x80000000 (2GB, KERNBASE) · KERNBASE + 0x10000000 · I/O Space · Kernel text · Kernel R/W data · end · free nenory (alloc()) · 0xe000000 (PHYSTOP) 234MB · 0xfe000000 · 0xffffffff · I/O Space · 2GB · 4GB · Kernel-memory

Physical

I/O Space · Kernel text · Kernel R/W data · 0 · Ununsed by xv6 · 0xe000000 (PHYSTOP) 234MB · Top of physical memory

Page table
Level 1          Level 2

| 0 - 4MB |
| 4 - 8MB |
| ... |
| 2GB - 2GB + 4MB |

| 0 - 4K |
| 4K - 8K |
| ... |
| (4MB-4K) - 4MB |

# Switch page tables

- Switch page tables
- Deallocate old page table

```
6398    switchuvm(proc);

6399    freevm(oldpgdir);

6400    return 0;
```
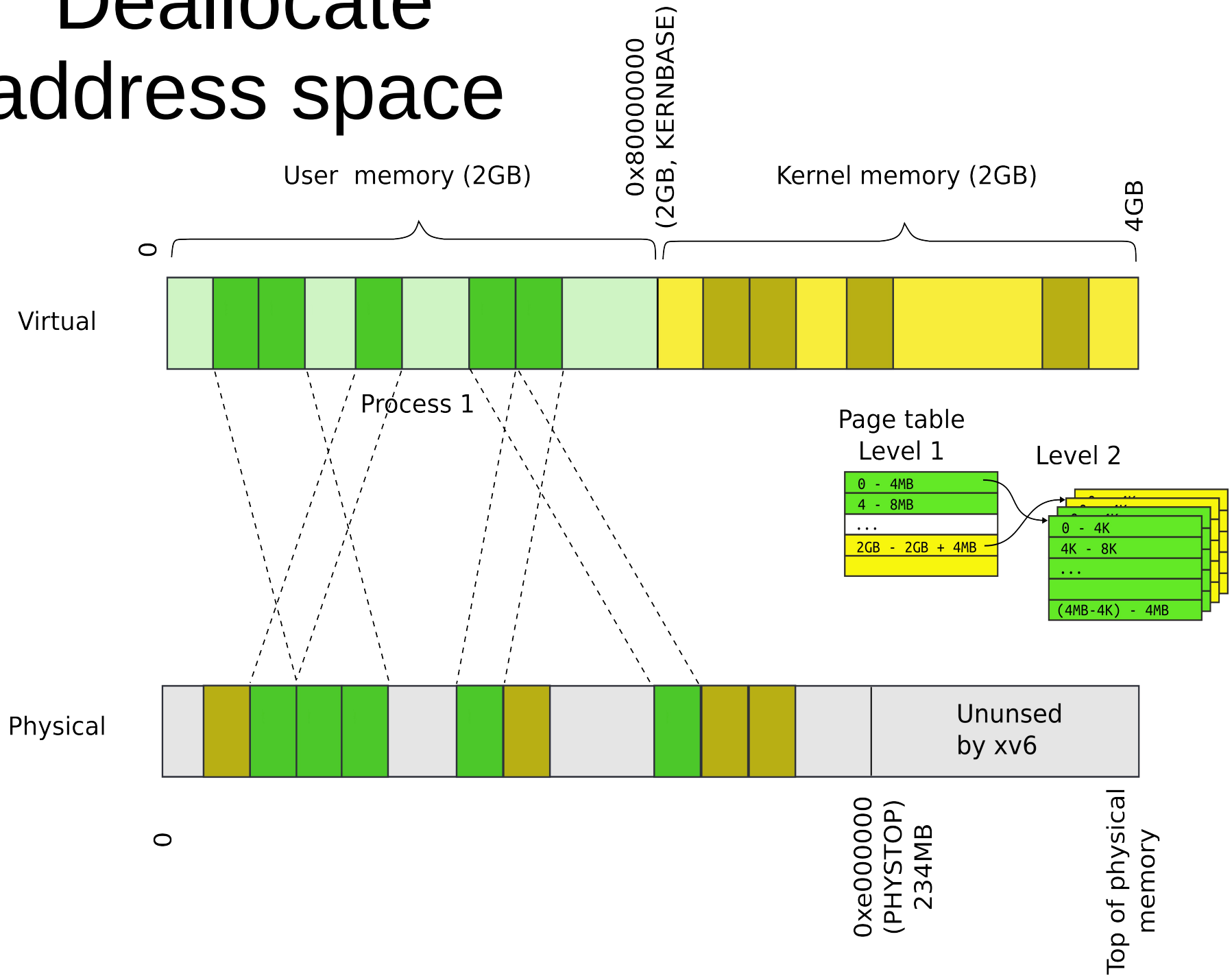
# Wait... which page table we are deallocating?

# Wait... which page table we are deallocating?

- Remember exec() replaces content of an already existing process
  - That process had a page table
  - We have to deallocate it

# Deallocate address space

User memory (2GB)

0x80000000 (2GB, KERNBASE)

Kernel memory (2GB)

4GB

0

Virtual

Process 1

Page table
Level 1

Level 2

| 0 - 4MB |
| 4 - 8MB |
| ... |
| 2GB - 2GB + 4MB |

| 0 - 4K |
| 4K - 8K |
| ... |
| (4MB-4K) - 4MB |

Physical

0

0xe000000 (PHYSTOP) 234MB

Ununsed by xv6

Top of physical memory

# Outline: deallocate process address space

- Walk the page table

  - Deallocate all pages mapped by the page table

- Deallocate pages that contain Level 2 of the page-table

- Deallocate page directory

```
2015 freevm(pde_t *pgdir)
2016 {
2017   uint i;
2018
2019   if(pgdir == 0)
2020     panic("freevm: no pgdir");
2021   deallocuvm(pgdir, KERNBASE, 0);
2022   for(i = 0; i < NPDENTRIES; i++){
2023     if(pgdir[i] & PTE_P){
2024       char * v = P2V(PTE_ADDR(pgdir[i]));
2025       kfree(v);
2026     }
2027   }
2028   kfree((char*)pgdir);
2029 }
```

Deallocate user address space

```
1987 deallocuvm(pde_t *pgdir, uint oldsz, uint newsz)
1988 {
...
1995   a = PGROUNDUP(newsz);
1996   for(; a < oldsz; a += PGSIZE){
1997     pte = walkpgdir(pgdir, (char*)a, 0);
1998     if(!pte)
1999       a += (NPTENTRIES - 1) * PGSIZE;
2000     else if((*pte & PTE_P) != 0){
2001       pa = PTE_ADDR(*pte);
2002       if(pa == 0)
2003         panic("kfree");
2004       char *v = P2V(pa);
2005       kfree(v);
2006       *pte = 0;
2007     }
2008   }
2009   return newsz;
2010 }
```

Walk page table and get pte

```
1987 deallocuvm(pde_t *pgdir, uint oldsz, uint newsz)
1988 {
...
1995   a = PGROUNDUP(newsz);
1996   for(; a < oldsz; a += PGSIZE){
1997     pte = walkpgdir(pgdir, (char*)a, 0);
1998     if(!pte)
1999       a += (NPTENTRIES - 1) * PGSIZE;
2000     else if((*pte & PTE_P) != 0){
2001       pa = PTE_ADDR(*pte);
2002       if(pa == 0)
2003         panic("kfree");
2004       char *v = P2V(pa);
2005       kfree(v);
2006       *pte = 0;
2007     }
2008   }
2009   return newsz;
2010 }
```

Deallocate a page

# Deallocate Level 2

```
2015 freevm(pde_t *pgdir)
2016 {
2017   uint i;
2018
2019   if(pgdir == 0)
2020     panic("freevm: no pgdir");
2021   deallocuvm(pgdir, KERNBASE, 0);
2022   for(i = 0; i < NPDENTRIES; i++){
2023     if(pgdir[i] & PTE_P){
2024       char * v = P2V(PTE_ADDR(pgdir[i]));
2025       kfree(v);
2026     }
2027   }
2028   kfree((char*)pgdir);
2029 }
```

# Recap

- We know how exec works!
- We can create new processes

# Questions?