# Cloud Operating Systems

Booting x86

**Fabian Rauscher, Daniel Gruss, Andreas Kogler**

2022-03-21

# Real Mode

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- 16 bit mode

- 16 bit mode
- Address space: 1 MB

- 16 bit mode
- Address space: 1 MB
- How is that possible?

Fabian Rauscher, Daniel Gruss, Andreas Kogler

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- 16 bit segment resgisters (CS, SS, DS, ES, FS, GS)

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- 16 bit segment resgisters (CS, SS, DS, ES, FS, GS)
- Every memory access uses a segment register and a 16 bit offset

- 16 bit segment resgisters (CS, SS, DS, ES, FS, GS)
- Every memory access uses a segment register and a 16 bit offset
- → Actual address is (segment register ≪ 4) + offset

- 16 bit segment resgisters (CS, SS, DS, ES, FS, GS)
- Every memory access uses a segment register and a 16 bit offset
- → Actual address is (segment register $\ll 4$) + offset
- Direct physical memory access

### 9.1.4    First Instruction Executed

The first instruction that is fetched and executed following a hardware reset is located at physical address FFFFFFF0H. This address is 16 bytes below the processor's uppermost physical address. The EPROM containing the software-
initialization code must be located at this address.

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Address: `0xFFFFFFF0`

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Address: `0xFFFFFFF0`
- How is that possible?

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Address: `0xFFFFFFF0`
- How is that possible?
  - CS register also has a 32-bit base address (initialized to `0xFFFF0000`)

- Address: `0xFFFFFFF0`
- How is that possible?
    - CS register also has a 32-bit base address (initialized to `0xFFFF0000`)
- What if I have $< 4\,GB$ RAM?

- Address: `0xFFFFFFF0`
- How is that possible?
  - CS register also has a 32-bit base address (initialized to `0xFFFF0000`)
- What if I have $< 4\,GB$ RAM?
  - physical address space $\neq$ RAM directly mapped

- BIOS initializes hardware platform

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- BIOS initializes hardware platform
- Select a device to boot from

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- BIOS initializes hardware platform
- Select a device to boot from
- Load MBR from device into memory (`0x7C00`)

- BIOS initializes hardware platform
- Select a device to boot from
- Load MBR from device into memory (`0x7C00`)
- Execute code from the MBR

- BIOS initializes hardware platform
- Select a device to boot from
- Load MBR from device into memory (`0x7C00`)
- Execute code from the MBR
- MBR code usually loads more data from the disk into memory

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- How can we interact with the hardware?

 Fabian Rauscher, Daniel Gruss, Andreas Kogler

- How can we interact with the hardware?
  - Where do we get the memory layout from?

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- How can we interact with the hardware?
  - Where do we get the memory layout from?
  - How can we access the disk?

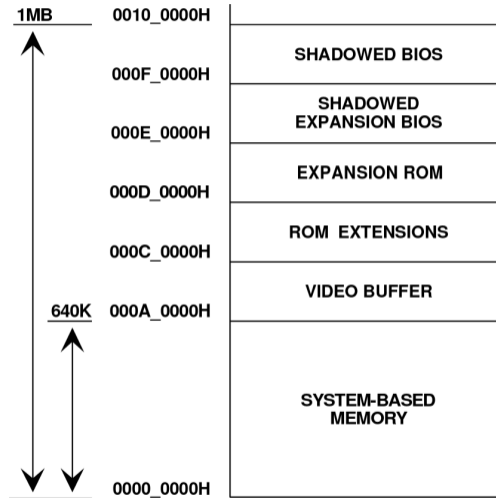Fabian Rauscher, Daniel Gruss, Andreas Kogler

- How can we interact with the hardware?
    - Where do we get the memory layout from?
    - How can we access the disk?
    - How can we configure the video output?

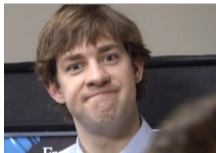Fabian Rauscher, Daniel Gruss, Andreas Kogler

- How can we interact with the hardware?
  - Where do we get the memory layout from?
  - How can we access the disk?
  - How can we configure the video output?
  - ...

- How can we interact with the hardware?
    - Where do we get the memory layout from?
    - How can we access the disk?
    - How can we configure the video output?
    - ...
$\rightarrow$ BIOS calls!

- How can we interact with the hardware?
  - Where do we get the memory layout from?
  - How can we access the disk?
  - How can we configure the video output?
  - ...
- → BIOS calls!
- We can trigger standardized software interrupts and let the BIOS handle it!

| | | |
|---|---|---|
| **1MB** | 0010_0000H | |
| | | SHADOWED BIOS |
| | 000F_0000H | |
| | | SHADOWED EXPANSION BIOS |
| | 000E_0000H | |
| | | EXPANSION ROM |
| | 000D_0000H | |
| | | ROM EXTENSIONS |
| | 000C_0000H | |
| | | VIDEO BUFFER |
| **640K** | 000A_0000H | |
| | | SYSTEM-BASED MEMORY |
| | 0000_0000H | |

Fabian Rauscher, Daniel Gruss, Andreas Kogler

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Memory accesses above 1MB wrap around

- Memory accesses above 1MB wrap around
- Used to fix software compatability issues
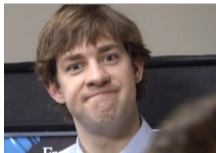
Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Memory accesses above 1MB wrap around
- Used to fix software compatability issues
- → Memory line 20 (A20) needs to be enabled by software to disable this "feature"

- Memory accesses above 1MB wrap around
- Used to fix software compatability issues
- → Memory line 20 (A20) needs to be enabled by software to disable this "feature"
- Multiple ways of doing this:

- Memory accesses above 1MB wrap around
- Used to fix software compatability issues
- → Memory line 20 (A20) needs to be enabled by software to disable this "feature"
- Multiple ways of doing this:
  - Keyboard controller

- Memory accesses above 1MB wrap around
- Used to fix software compatability issues
- → Memory line 20 (A20) needs to be enabled by software to disable this "feature"
- Multiple ways of doing this:
  - Keyboard controller
  - Fast A20 gate

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Memory accesses above 1MB wrap around
- Used to fix software compatability issues
- → Memory line 20 (A20) needs to be enabled by software to disable this "feature"
- Multiple ways of doing this:
  - Keyboard controller
  - Fast A20 gate
  - BIOS call

- Memory accesses above 1MB wrap around
- Used to fix software compatability issues
→ Memory line 20 (A20) needs to be enabled by software to disable this "feature"
- Multiple ways of doing this:
    - Keyboard controller
    - Fast A20 gate
    - BIOS call
    - ...

- Disable interrupts

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Disable interrupts
- Setup Global Descriptor Table (GDT)

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Disable interrupts
- Setup Global Descriptor Table (GDT)
- Load GDT

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Disable interrupts
- Setup Global Descriptor Table (GDT)
- Load GDT
- Enable protected mode by setting bit 0 in CR0

- Disable interrupts
- Setup Global Descriptor Table (GDT)
- Load GDT
- Enable protected mode by setting bit 0 in CR0
- Long jump to set CS

- Disable interrupts
- Setup Global Descriptor Table (GDT)
- Load GDT
- Enable protected mode by setting bit 0 in CR0
- Long jump to set CS
- Load DS, ES, FS, GS, SS, ESP

# Protected Mode

Fabian Rauscher, Daniel Gruss, Andreas Kogler

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- 32 bit mode

- 32 bit mode
- More 32 bit registers

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- 32 bit mode
- More 32 bit registers
- Access to up to 4GB of memory

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- 32 bit mode
- More 32 bit registers
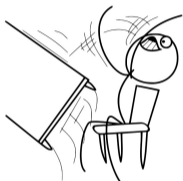- Access to up to 4GB of memory
- Segmentation uses the GDT

- 32 bit mode
- More 32 bit registers
- Access to up to 4GB of memory
- Segmentation uses the GDT
  - It is possible to address the whole address space without switching segments

- 32 bit mode
- More 32 bit registers
- Access to up to 4GB of memory
- Segmentation uses the GDT
  - It is possible to address the whole address space without switching segments
  - → Segmentation more optional!

- 32 bit mode
- More 32 bit registers
- Access to up to 4GB of memory
- Segmentation uses the GDT
  - It is possible to address the whole address space without switching segments
  - → Segmentation more optional!
- Optional paging

Fabian Rauscher, Daniel Gruss, Andreas Kogler

Fabian Rauscher, Daniel Gruss, Andreas Kogler

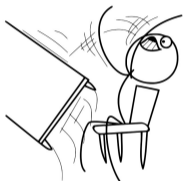Fabian Rauscher, Daniel Gruss, Andreas Kogler

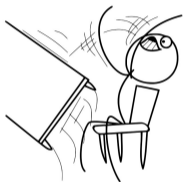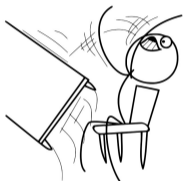Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Segment registers are offsets into the GDT

- Segment registers are offsets into the GDT
- GDT is an array of segment descriptors that each hold ...

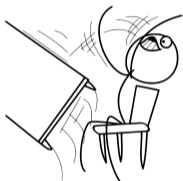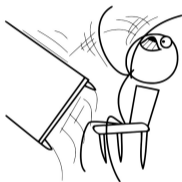Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Segment registers are offsets into the GDT
- GDT is an array of segment descriptors that each hold ...
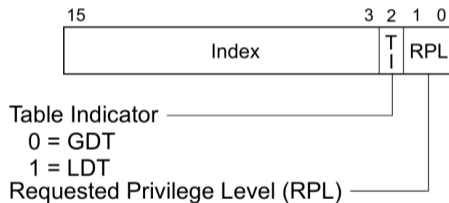  - Base address

- Segment registers are offsets into the GDT
- GDT is an array of segment descriptors that each hold ...
    - Base address
    - Limit

- Segment registers are offsets into the GDT
- GDT is an array of segment descriptors that each hold ...
  - Base address
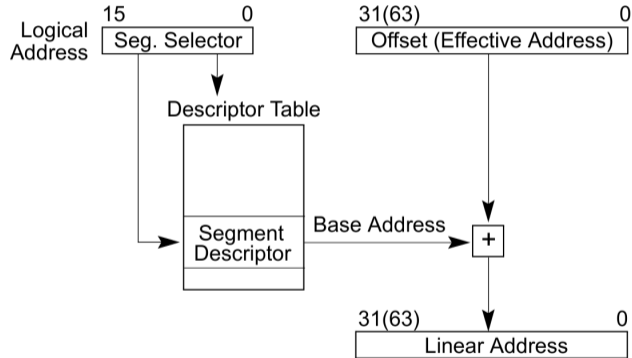  - Limit
  - Access rights

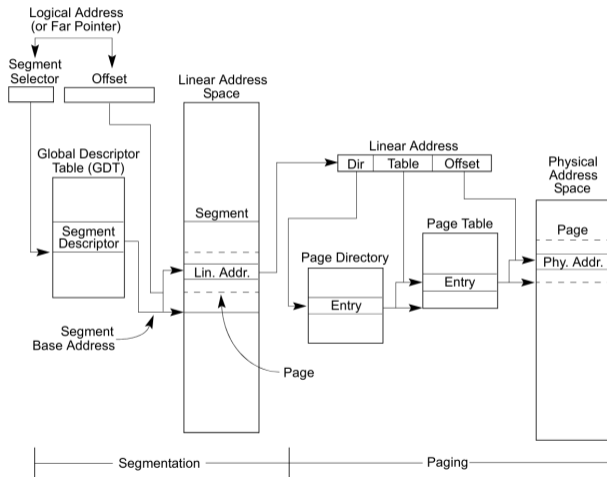- Segment registers are offsets into the GDT
- GDT is an array of segment descriptors that each hold ...
    - Base address
    - Limit
    - Access rights
    - Flags

Table Indicator
0 = GDT
1 = LDT
Requested Privilege Level (RPL)

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Disable interrupts

- Disable interrupts
- Setup 64 bit segmend descriptors

 Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Disable interrupts
- Setup 64 bit segmend descriptors
- Setup paging structure and set CR3

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Disable interrupts
- Setup 64 bit segmend descriptors
- Setup paging structure and set CR3
- Enable long Mode via bit 8 of the EFER MSR

- Disable interrupts
- Setup 64 bit segmend descriptors
- Setup paging structure and set CR3
- Enable long Mode via bit 8 of the EFER MSR
- Enable PAE via bit 5 of CR4

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Disable interrupts
- Setup 64 bit segmend descriptors
- Setup paging structure and set CR3
- Enable long Mode via bit 8 of the EFER MSR
- Enable PAE via bit 5 of CR4
- Enable paging via bit 31 in CR0

- Disable interrupts
- Setup 64 bit segmend descriptors
- Setup paging structure and set CR3
- Enable long Mode via bit 8 of the EFER MSR
- Enable PAE via bit 5 of CR4
- Enable paging via bit 31 in CR0
- Reload segment selectors to enter long mode proper

Fabian Rauscher, Daniel Gruss, Andreas Kogler

# Long Mode

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Only supported with paging enabled

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Only supported with paging enabled
- Register extensions and new 64 bit registers

- Only supported with paging enabled
- Register extensions and new 64 bit registers
  - $\rightarrow$ RAX, RBX, RCX RDX, ..., R8-R15

- Only supported with paging enabled
- Register extensions and new 64 bit registers
  - → RAX, RBX, RCX RDX, ..., R8-R15
- 64 bit pointers

- Only supported with paging enabled
- Register extensions and new 64 bit registers
  - → RAX, RBX, RCX RDX, ..., R8-R15
- 64 bit pointers
- 48 bit virtual addresses (57 bit with 5-level paging)

- Only supported with paging enabled
- Register extensions and new 64 bit registers
    - → RAX, RBX, RCX RDX, ..., R8-R15
- 64 bit pointers
- 48 bit virtual addresses (57 bit with 5-level paging)
- Segmentation is simplified

- Only supported with paging enabled
- Register extensions and new 64 bit registers
  - → RAX, RBX, RCX RDX, ..., R8-R15
- 64 bit pointers
- 48 bit virtual addresses (57 bit with 5-level paging)
- Segmentation is simplified
  - All segment limits are ignored

- Only supported with paging enabled
- Register extensions and new 64 bit registers
    - → RAX, RBX, RCX RDX, ..., R8-R15
- 64 bit pointers
- 48 bit virtual addresses (57 bit with 5-level paging)
- Segmentation is simplified
    - All segment limits are ignored
    - CS, SS, ES, DS have a base of 0

- Only supported with paging enabled
- Register extensions and new 64 bit registers
  - → RAX, RBX, RCX RDX, ..., R8-R15
- 64 bit pointers
- 48 bit virtual addresses (57 bit with 5-level paging)
- Segmentation is simplified
  - All segment limits are ignored
  - CS, SS, ES, DS have a base of 0
  - Bases for FS and GS can be set via MSRs

Fabian Rauscher, Daniel Gruss, Andreas Kogler

# Booting SWEB

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Grand Unified Bootloader (GRUB)

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Grand Unified Bootloader (GRUB)
- Runs in real mode and protected mode

- Grand Unified Bootloader (GRUB)
- Runs in real mode and protected mode
  - → Protected mode: most of the time

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Grand Unified Bootloader (GRUB)
- Runs in real mode and protected mode
    - $\rightarrow$ Protected mode: most of the time
    - $\rightarrow$ Real mode: entry and BIOS calls

- Grand Unified Bootloader (GRUB)
- Runs in real mode and protected mode
  - → Protected mode: most of the time
  - → Real mode: entry and BIOS calls
- Basic hardware detection and framebuffer setup

- Grand Unified Bootloader (GRUB)
- Runs in real mode and protected mode
  - → Protected mode: most of the time
  - → Real mode: entry and BIOS calls
- Basic hardware detection and framebuffer setup
- Loads SWEB into memory

- Grand Unified Bootloader (GRUB)
- Runs in real mode and protected mode
  - → Protected mode: most of the time
  - → Real mode: entry and BIOS calls
- Basic hardware detection and framebuffer setup
- Loads SWEB into memory
- Provides hardware information according to the Multiboot spec

- Grand Unified Bootloader (GRUB)
- Runs in real mode and protected mode
  - → Protected mode: most of the time
  - → Real mode: entry and BIOS calls
- Basic hardware detection and framebuffer setup
- Loads SWEB into memory
- Provides hardware information according to the Multiboot spec
- Runs SWEB

- Provides standardized format for passing information from the Bootloader to the OS

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Provides standardized format for passing information from the Bootloader to the OS
  - Memory areas and their type

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Provides standardized format for passing information from the Bootloader to the OS
    - Memory areas and their type
    - Location of the framebuffer and the configured video mode

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Provides standardized format for passing information from the Bootloader to the OS
  - Memory areas and their type
  - Location of the framebuffer and the configured video mode
  - Drive infos

- Provides standardized format for passing information from the Bootloader to the OS
  - Memory areas and their type
  - Location of the framebuffer and the configured video mode
  - Drive infos
  - ...

Fabian Rauscher, Daniel Gruss, Andreas Kogler

- Provides standardized format for passing information from the Bootloader to the OS
  - Memory areas and their type
  - Location of the framebuffer and the configured video mode
  - Drive infos
  - ...
- Pointer to the multiboot header passed via EBX to the kernel entry