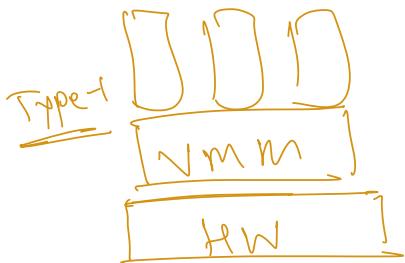


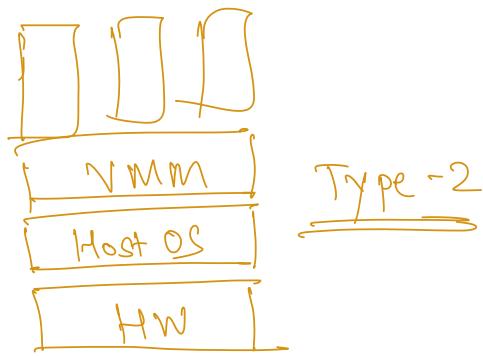
- What we are building is called VMM.

- 2 types -



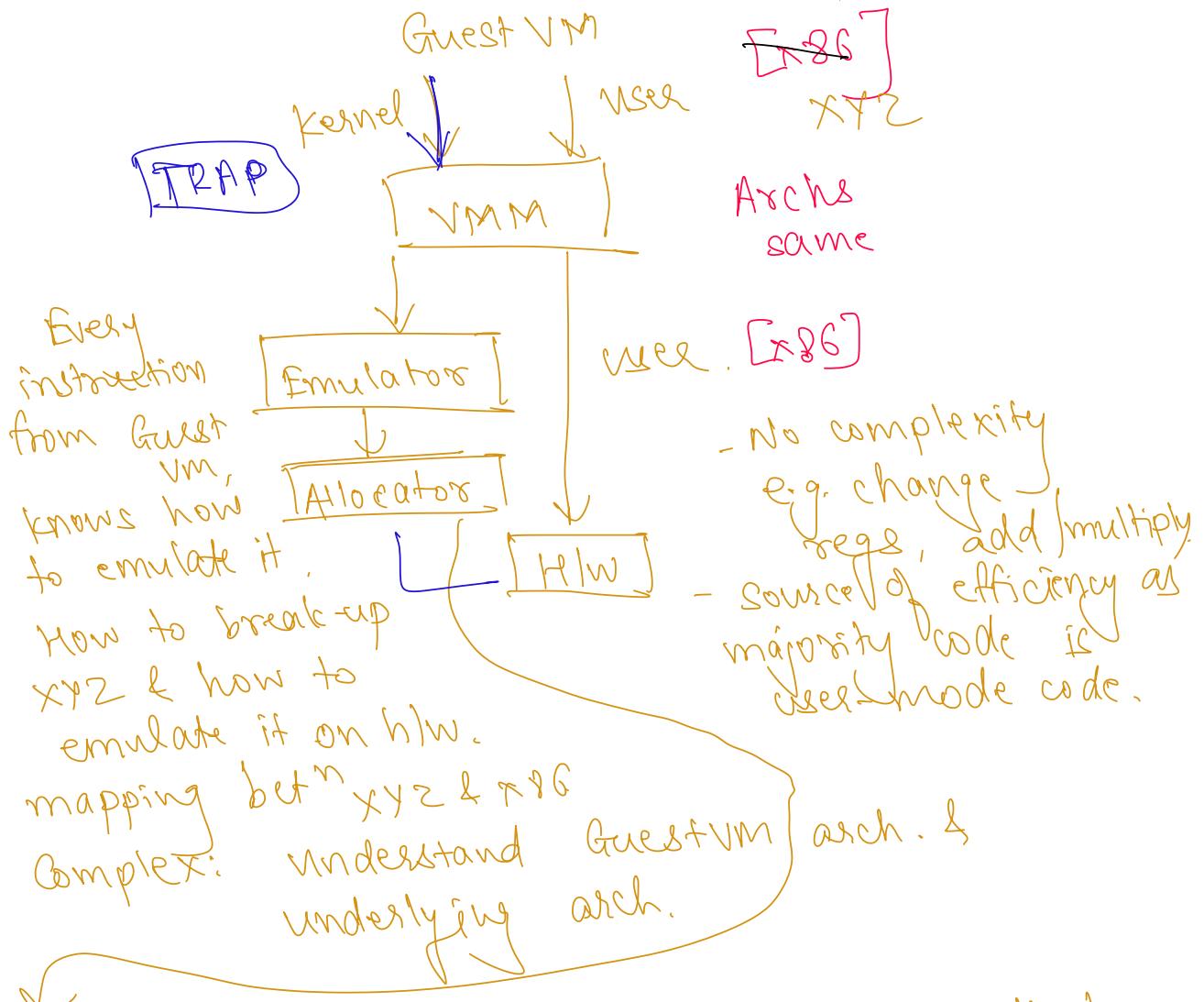
⇒ Xen
⇒ VMware workstation

- VMM directly on HW. No Host OS
- Responsibility of VMM to manage HW resources & divide among diff. VMs



- ⇒ QEMU
⇒ VirtualBox
- specific Host OS that controls HW
 - VMM works in conjunction with host OS (e.g. KVM) works with Linux to manage HW & multiplex
 - In VirtualBox, host OS thinks it is just another application
 - VBox spawns processes for VMs
 - Some cases: OS unaware
 - Some cases: Tighter Coupling

Broad View :



When guestVM: memory allocations, it thinks

- it owns whole machine -
- somebody intercept it & allocate real h/w resources so that prog can run.
- e.g. Guest VM: Create process that requires 64KB DRAM. Has to exist somewhere.
- If no Host O.S., everything done by VMM.

- If Host OS, allocator managed by Allocator. Host OS has access to hw resources.

Architectural req.: Kernel mode instruction must be intercepted.

Goldberg & Popel: 3 theorems.

Privileged inst.: Directly manipulating hw:

Allocating mem.,
creating & reg

Sensitive inst.: They read data.

e.g. read current add of interrupt table. (confidential)
User-mode inst. that dole

this w/o trap.

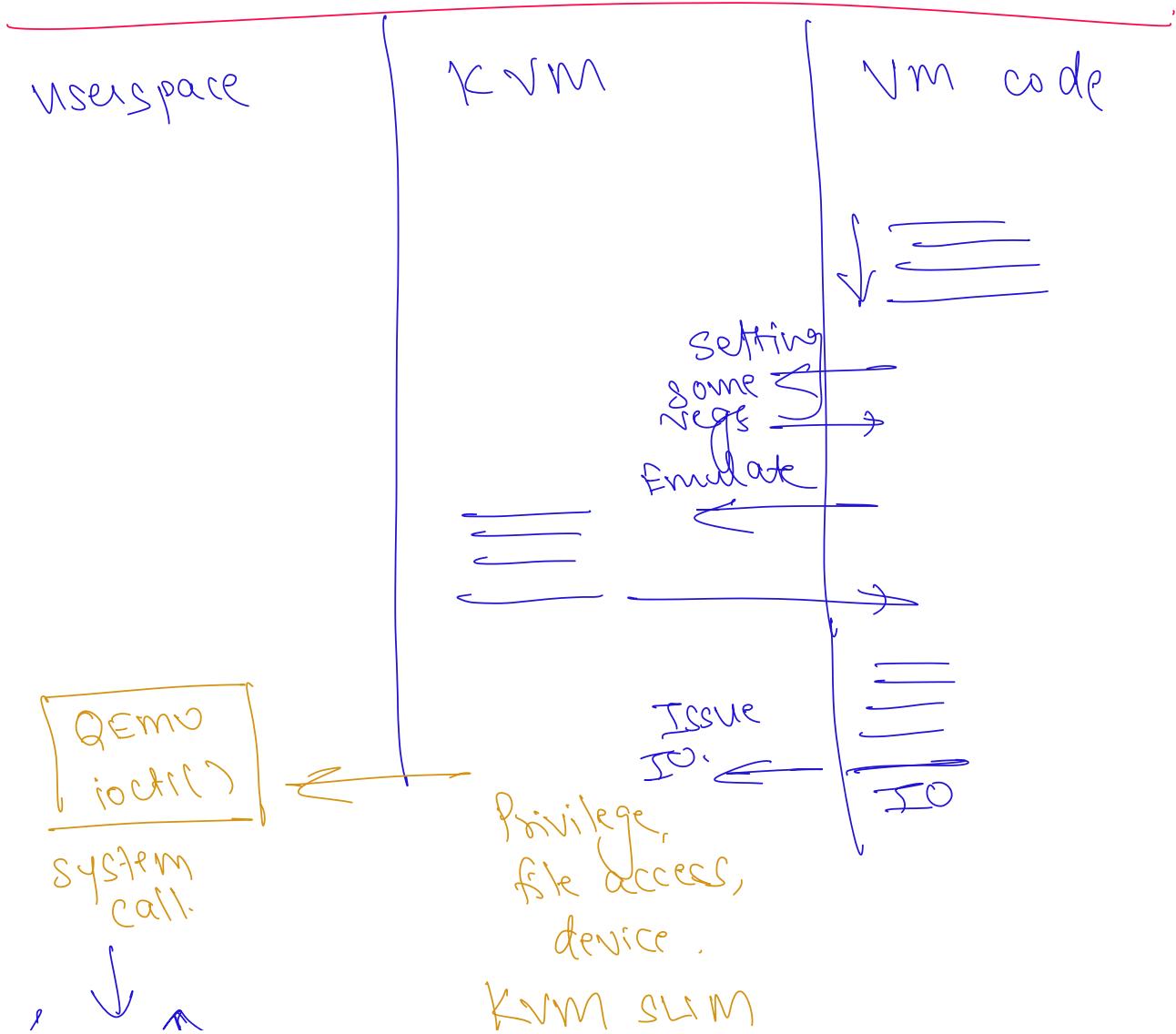
must cause a TRAP

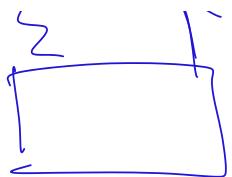
QEMU: Full system Simulator.

- can simulate x86, PowerPC.
- w/o any support is fully capable of simulating it.
- Does the loop: Every instruction:
emulate using hw

- no hardware support needed.
- Not efficient. Slow
- KVM: Accelerate this. use hw for some things.
- Helps get some machine
- KVM representing host O.S.

WHAT IS THE USE OF KVM IN QEMU?





I/O Device

→ very heavy
orders of
magnitude
slower than
instructions



Inst. 0.5 ns

HDD: 10 ms.

Handle I/O out of
critical path.

Instantiate I/O.
With QEMU, KVM
slimmer.

VT-X (2005)

→ Hardware support for virt. Intel 2005.

→ VMware = binary translation

Xen: paravirtualization.

most popular approach

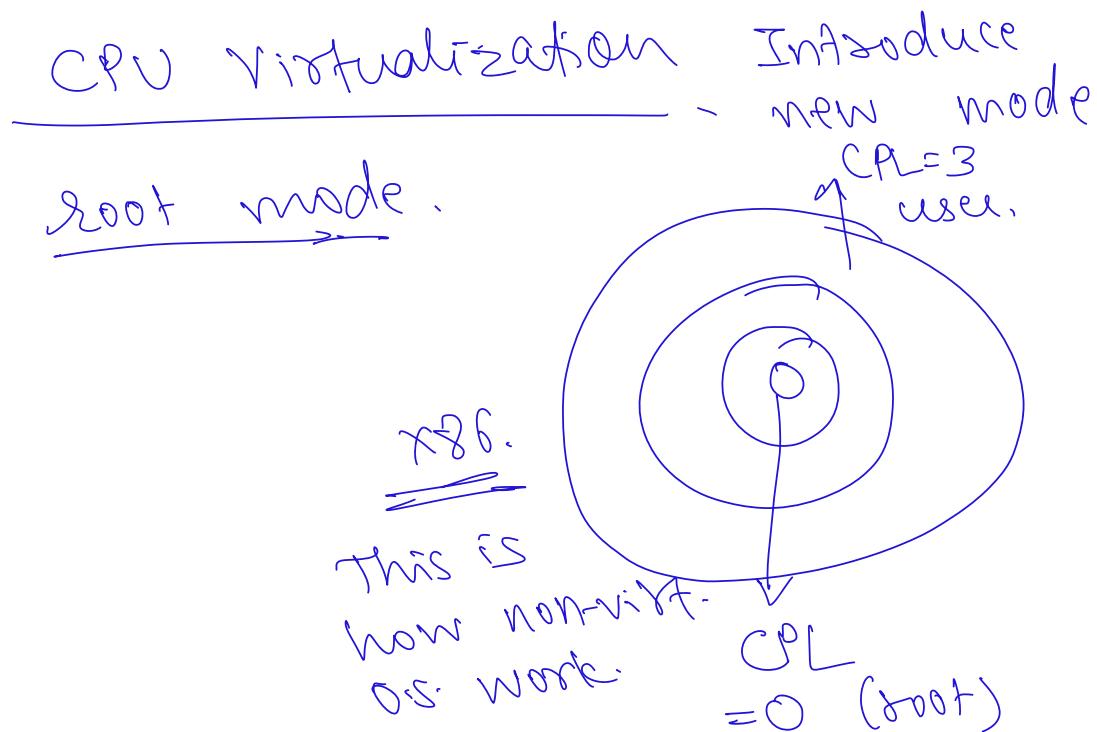
Goal: hardware support for virt.
such that unmodified x86-based
OS directly on HW.

- first version 2005: focus not perf.
Just get it functional

- over time: faster.

- In 2005, perf. was inferior than

- Since 2005: NTx faster other approaches instructions that NTx introduces



NTx introduces new mode called root mode

HW → Root mode → 4 rings
[ARCH STATE]

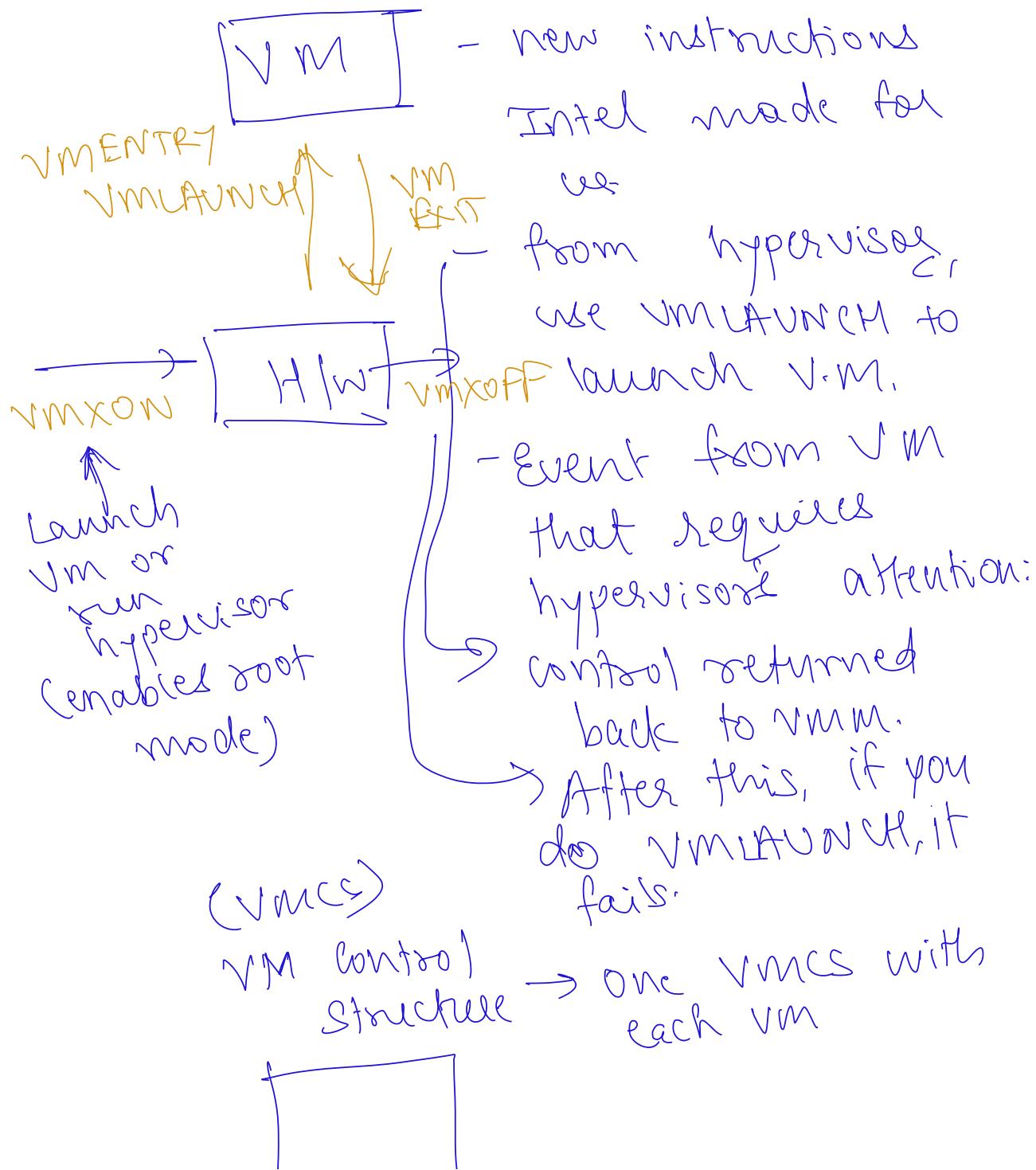
HW → non-root mode → 4 rings.
[ARCH STATE]

- Before: physical machine had registers which VM has access to but only written by hypervisor

2 sets of resources:

- One set just for root mode
- Another set just for non-root mode.
- Makes it easy: Even if we have sensitive inst. that is going to read state, it is going to be non-root state
[Never see root mode state]
- Fact of root & non-root mode orthogonal to inst. that Intel already had.
- Whenever you read reg.: You had no idea whether root or non-root mode
 - Assumed that it is the only privileged reg. in system, but actually 2 copies.
 - unless in root mode, cannot access state associated with root mode.

→ Doubled resources on the chip
Kept one for hypervisor, other
for V.M.

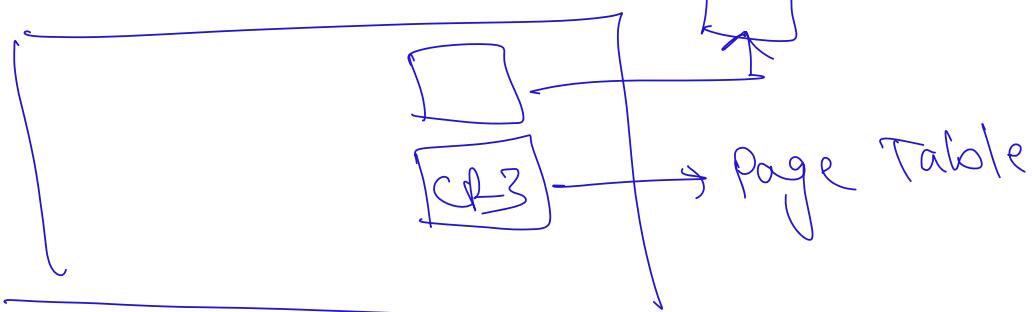


F H

- whenever there is VM EXIT for eg., VMCS populated & stored in mem.
- on VMLAUNCH & VMENTRY, data from VMCS used to populate reg in h/w.
 - stored in physical DRAM.

Reg. points to VMCS structure.

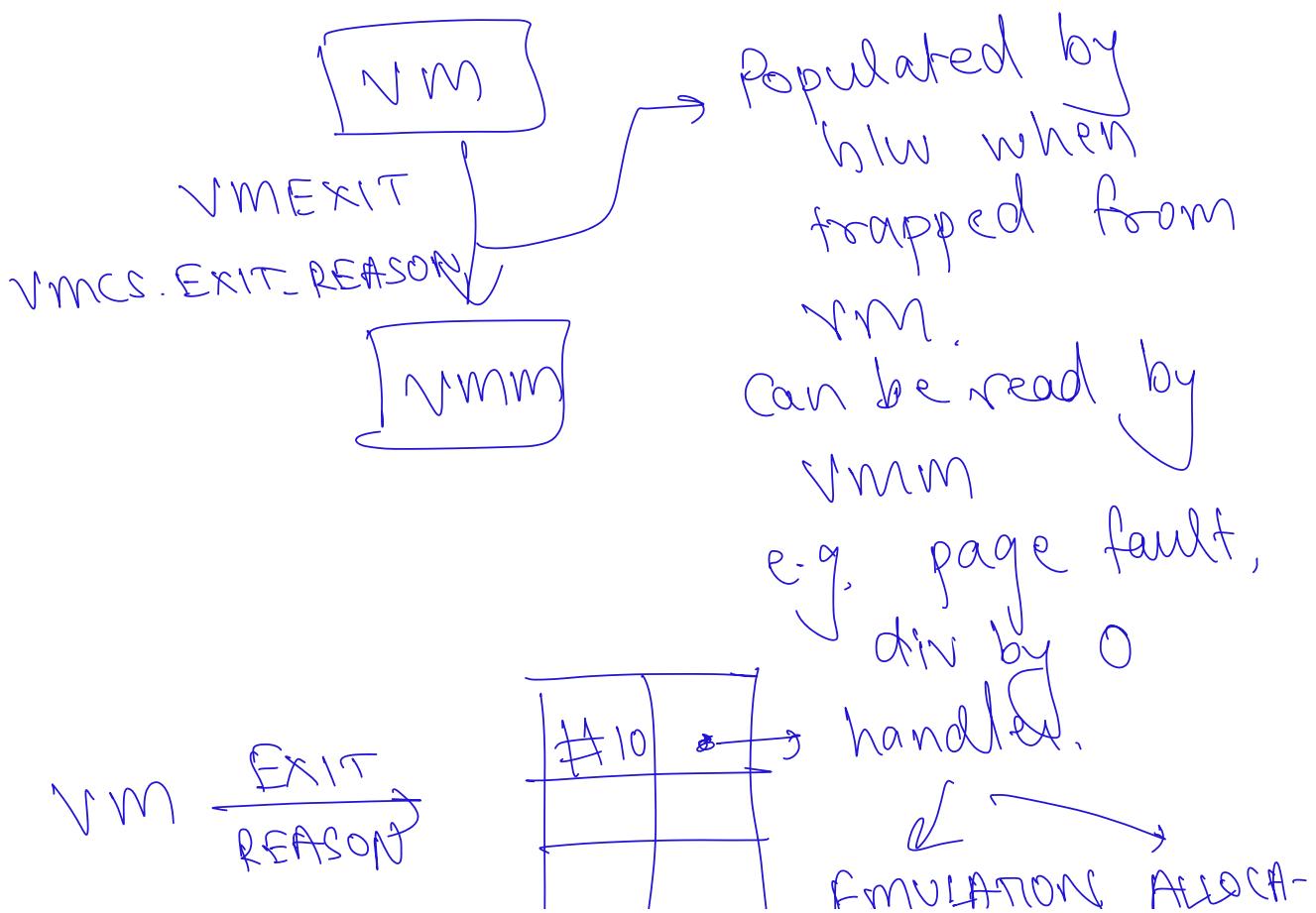
VMCS.



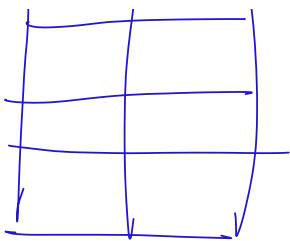
Every time context switch b/w V.M,
reg. points to diff. VMCS structure.

- VMCS^{struct} can only be read using VMREAD (from VMM)
- write: VMWRITE (from VMM)

- Structure of VMCS not specified by INTEL.
- Only readable & writable by VMREAD & VMWRITE.
- Allows Intel to change how VMCS is stored.
 - Not any guarantee provided by Intel.
- Whenever VMEXIT -



Total 55
diff. reasons
for TRAP.
e.g. exception, pte fault.



ITION

- EXIT REASON: what happened,
= handle invoked & handle it
e.g. if req emulation, can emulate,
if requires allocation resources.
- subtle thing in how we handle
exit from VMM

↳ Retry instruction.

e.g. READ 2000

Re-execute
some
instruction,
not next

Fault
↓
VMM

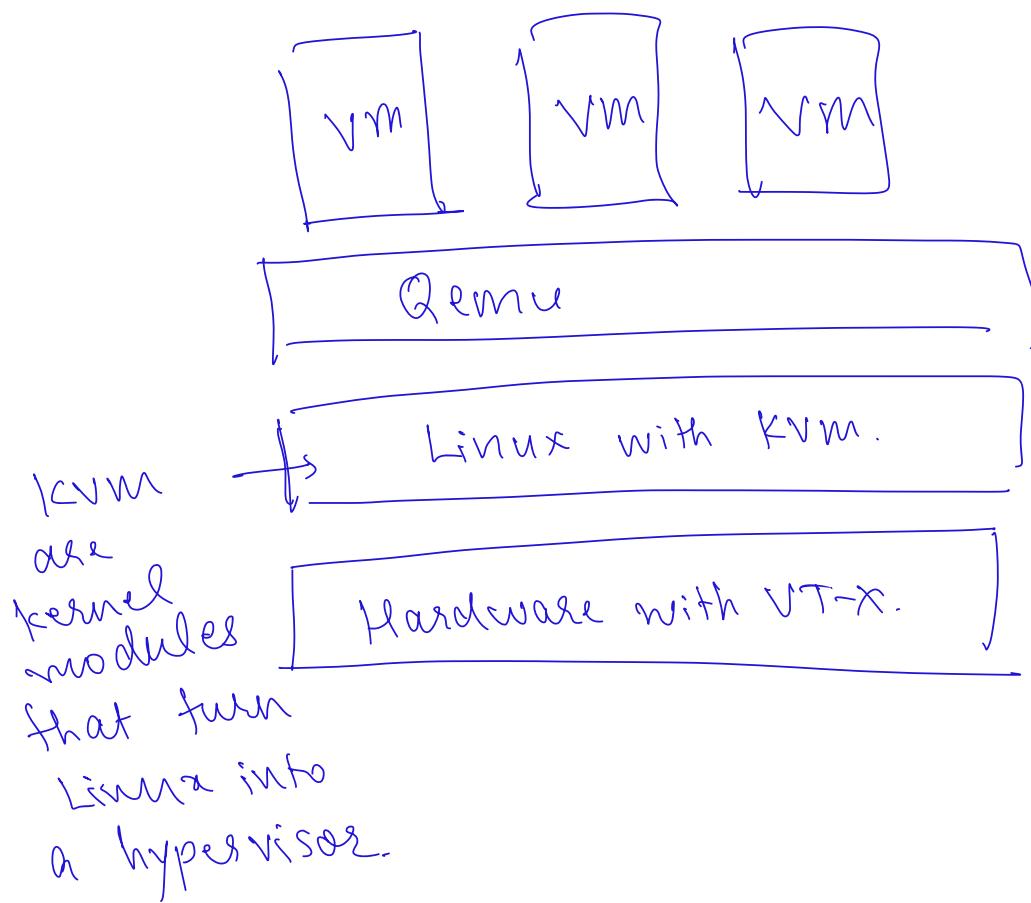
Do not increment instruction pth.

e.g. I-pointee.

↳ Increment i.p.
e.g. hw doesn't support inst.,
emulate it & increment

ptr by l.

- Confusion: if trapped always next instruction.



2 main modules:

kvm.ko → core module

kvm-intel.ko → machine-specific module

kvm-amd.ko.

KVM: Launches VM_LAUNCH | VM_RESUME
set up VMCS, deal with VMEXIT.

QEMU talks to KVM using system calls.
ioctls calls.

QEMU:

User-space component. Can run Guest using
powerpc & host using x86

2-modes: Emulator & Virtualizer



Fake
hardware;

binary translation.

No KVM.

Gets help
from KVM

Because QEMU supports emulation, it can
emulate I/O:

e.g. assume that VM on QEMU emulates
certain serial port that doesn't exist

on host.

When VM tries I/O on port, VM exits to KVM. KVM looks at reason & passes to QEMU.

QEMU emulates the I/O device & passes control to KVM.

KVM executes VMRESUME to let VM proceed.

