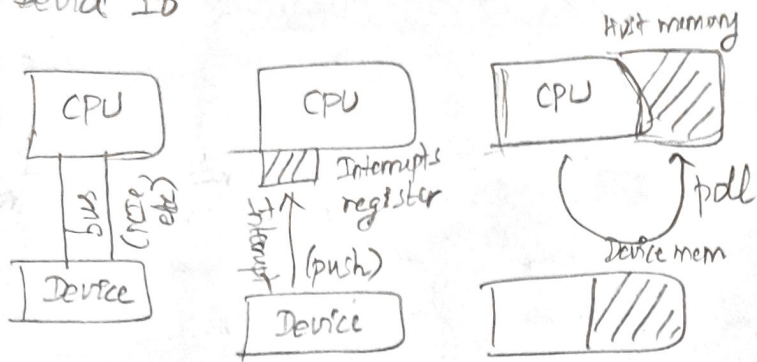


① Device IO



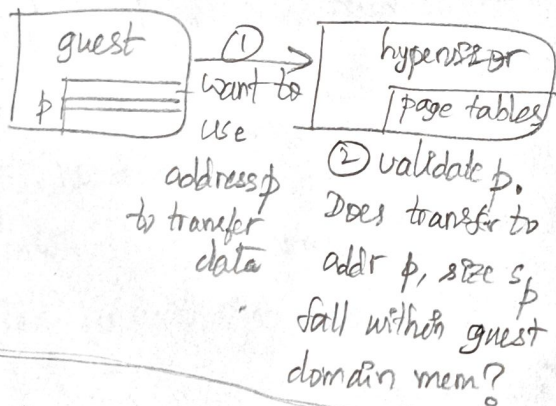
- ① Mix Interrupts with pulling
- ② network: interrupt → interrupt disable → poll for timeout / until no data → reusable interrupts

③ how should IO be virtualized?

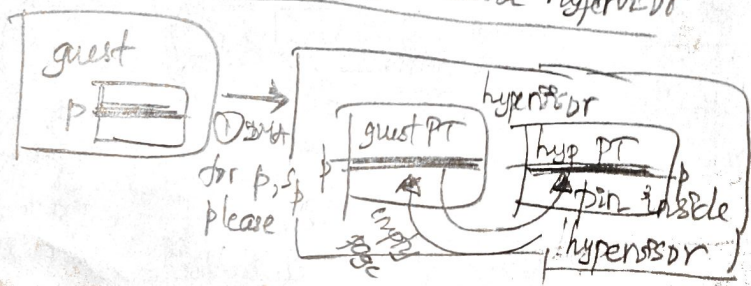
- ① guest can't ^(directly) program sensitive device registers
- ② device I/O on behalf of guest must be vetted by the hypervisor - addresses part of guest domain.

⑤ Approach 2 Validation of guest memory in DMA descriptors

- ① Don't copy. Validate guest memory.

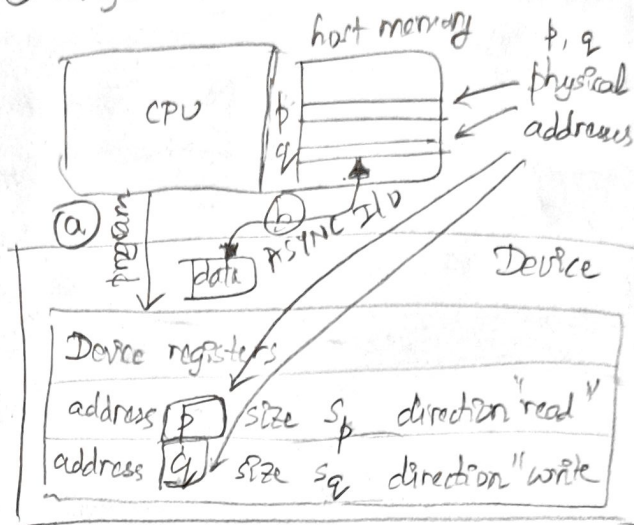


⑥ Approach 3 Pin pages inside hypervisor

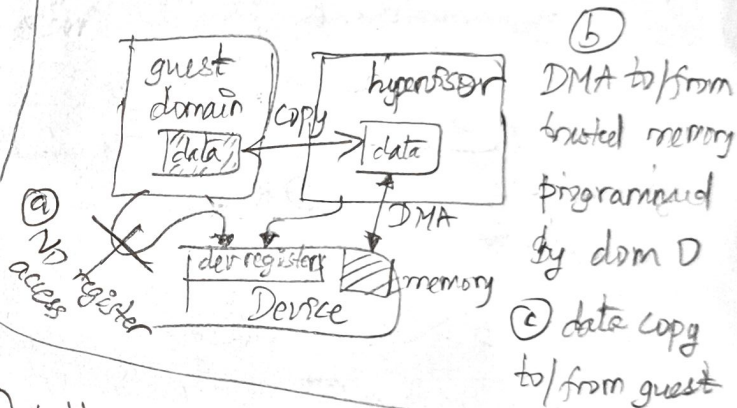


② Direct Memory Access (DMA)

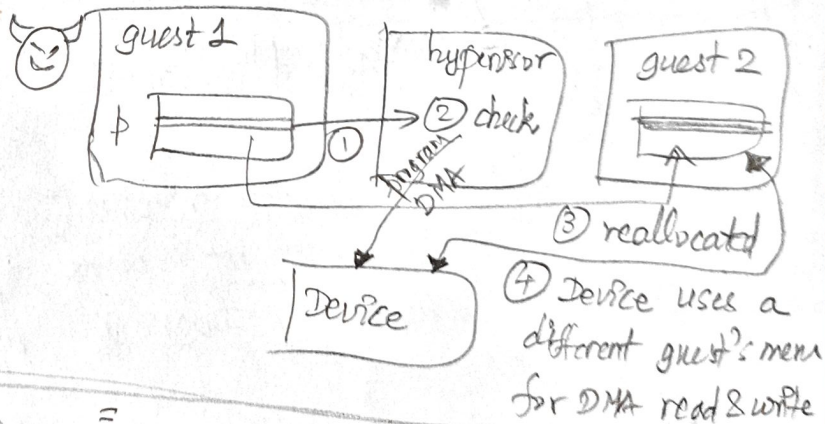
use CPU cycles when data moved
 ① large data ② scattered data



④ Approach 1 data copy



① problem: time of check to time of use



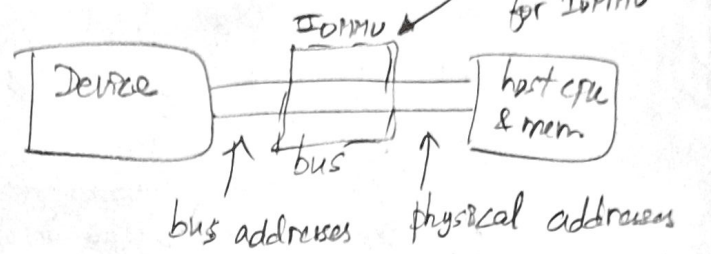
① upon DMA request validation, pin page inside hypervisor's own page tables. Cannot be reallocated to anyone else until DMA complete.

⑦ Need for hardware support

- validation per DMA data transfer
- 30% of native network throughput

⑧ IO Memory Management unit (IOMMU)

Hardware to translate & validate devices accessing host memory.



⑧a translation: useful when device can only access part of host memory (e.g.: 32 bit device versus 64 bit cpu)

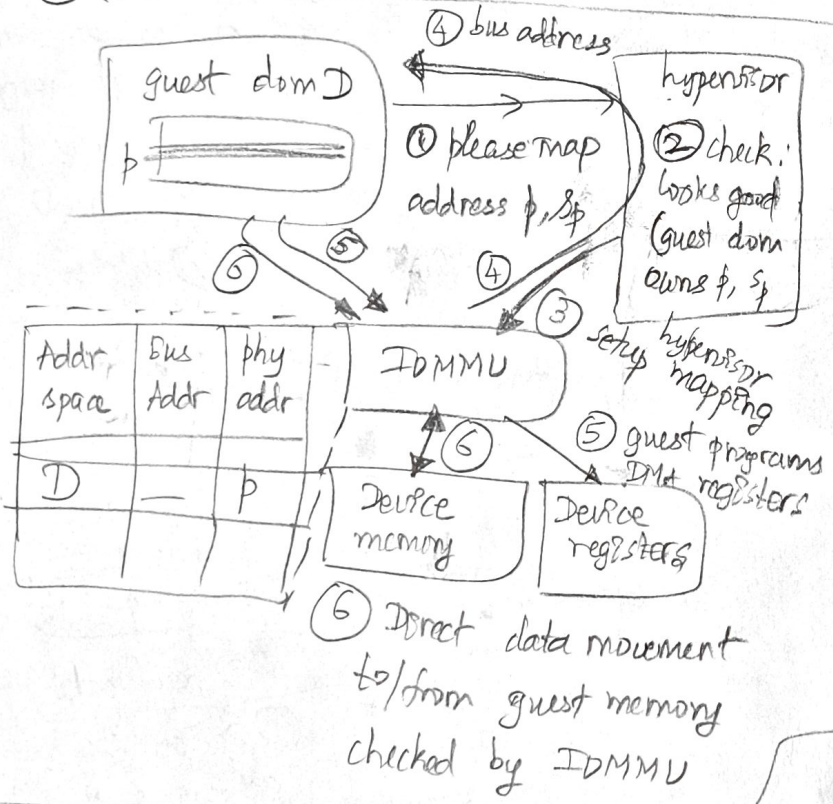
⑧b Isolation IOMMU has a page table. Entries are only mapped/unmapped by privileged entity (hypervisor) after validation.

⑧c page table structure

address space ID	bus address	physical address

Key → value

⑨ Data movement with IOMMU & address space isolation



- ① separate address spaces per guest domain
- ② validation at map/unmap time; guest gets a bus address
- ③ direct DMA access from guest for all mapped bus addresses (checking directly in hardware IOMMU at the time of use)
- ④ share most other resources

⑩ Single Root IO virtualization (SR-IOV)

- ① expose single device as multiple devices (PF/VF/UNC)
- ② dedicated DMA data movement registers per VF
- ③ sensitive registers programmed by privileged entity

