# KVM
## F O R U M

*Sharing IOMMU page tables with TDP in KVM*

**Lu Baolu baolu.lu@intel.com**
**Zhao Yan yan.y.zhao@intel.com**
**Tian Kevin kevin.tian@intel.com**

Sep. 2021

THE
**LINUX**
FOUNDATION

# Disclaimers

THE
**LINUX**
FOUNDATION

# Agenda

- Goal
- Sharing Advantages
- Sharing Prerequisites
- Sharing Interfaces
- Page & Page table Pinning
- Shared Page Table Root Update
- Bootup Performance
- TODOs

THE
LINUX
FOUNDATION

# Goal

# Sharing Advantages

- Reduced memory footprint

- Unified page table management
  - Dirty page tracking, page fault handling, etc.

- Probably higher performance by reducing unnecessary EPT/NPT zap

# Sharing Prerequisites

- The same address space

- Compatible page table format

- Non-conflicting page table content

THE
**LINUX**
FOUNDATION

# The Same Address Space

- Address space is GPA (L1) → HPA

- Qemu
  - KVM side
    - check TDP is enabled
    - vCPU model does not include EPT/NPT feature
  - IOMMU side
    - no vIOMMU
    - vIOMMU is not in shadow mode. (nested mode on GPA is ok)

.

# The Same Address Space (Cont.)

- Nested VM
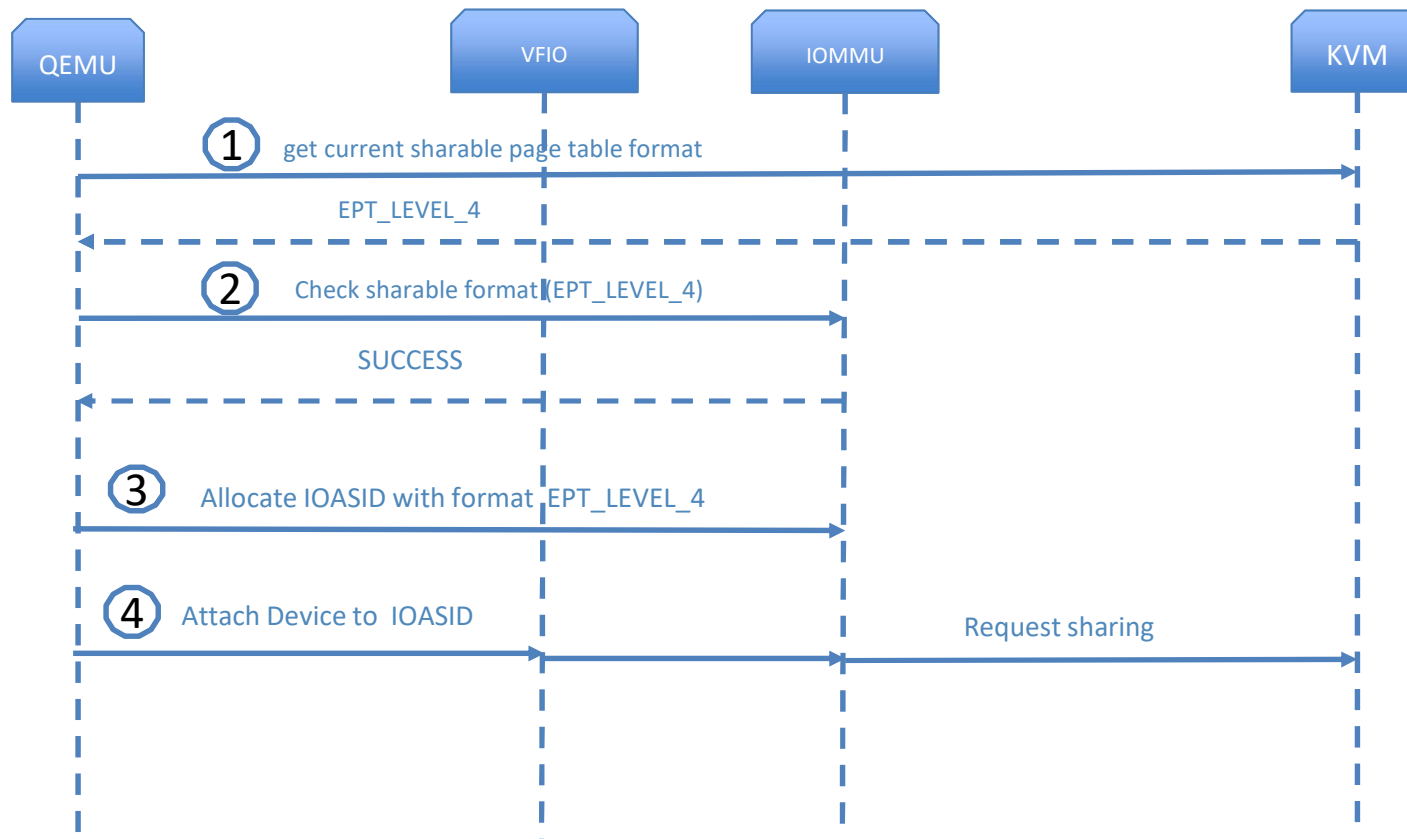  - TBD currently

- SMM in x86
  - A different address space. Cannot be shared to IOMMU.
  - Non-SMM mode EPT must be kept for sharing when vCPU is in SMM mode.

# Compatible Page Table Formats

- Unified compatible page table format definition across KVM and IOMMU

- Compatible page table formats
  - FORMAT_EPT_LEVEL_4
  - FORMAT_EPT_LEVEL_5
  - FORMAT_NPT_LEVEL_4
  - FORMAT_NPT_LEVEL_5
  - ...

# Sharing Handshake Sequence

QEMU      VFIO      IOMMU      KVM

① get current sharable page table format

EPT_LEVEL_4

② Check sharable format (EPT_LEVEL_4)

SUCCESS

③ Allocate IOASID with format EPT_LEVEL_4

④ Attach Device to IOASID      Request sharing

Note:
1. device pass-through is based on the /dev/iommu proposal, which is IOASID oriented.
2. KVM shares TDP used by vCPU 0

THE LINUX FOUNDATION
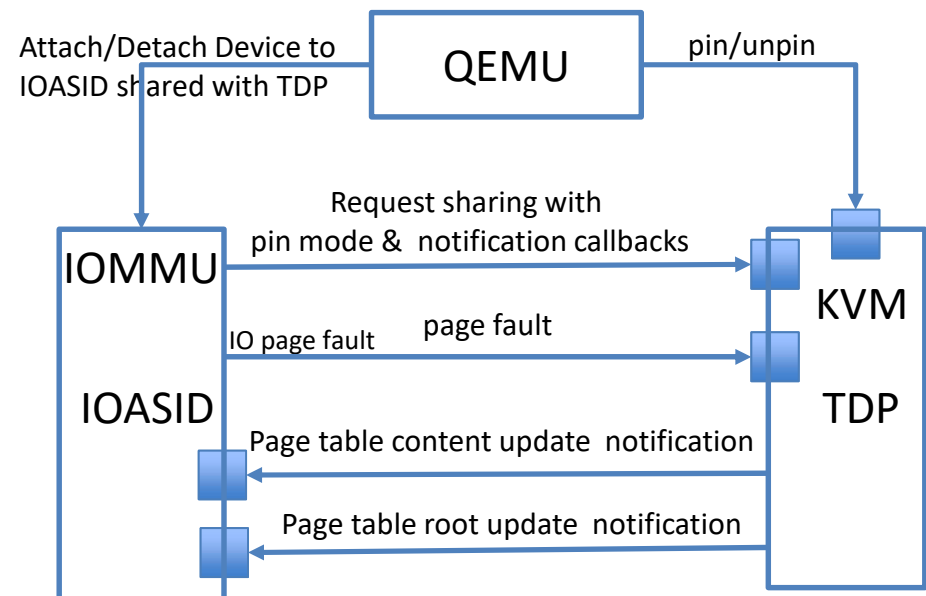
# Non-conflicting Page Table Content

- Presence of page table entry
  - For KVM user memslots
    - must be present and pinned (staying present) for DMA pages when IO page fault is not supported.
    - Can be present or zapped for non-DMA pages or when IO page fault supported

  - For KVM private memslots
    - Not present in IOPT before sharing
    - Safe to be present in IOPT after sharing

    - Local APIC
      - DMA write to 0xfeexxxxx doesn't go through DMA remapping.

    - TSS and IDENTITY_PAGETABLE
      - for !enable_unrestricted_guest, E820 Reserved

# Non-conflicting Page Table Content (Cont.)

- Read/Write/Execute bit
  - RO for RO memslots
  - RW for other memslots
  - Execute bit
    - currently ignored in IOMMU and no device uses it.
  - Write protection for live migration
    - Allowed when IO page fault is supported
    - Must be disabled otherwise
      - All pinned ranges are dirty or
      - traversal for Dirty bit

THE
LINUX
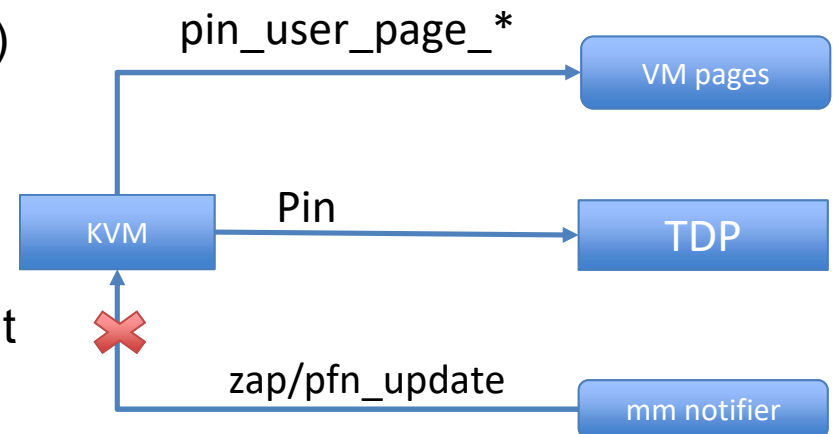FOUNDATION

# Sharing Interfaces

- Request/stop sharing

- Page/page table pinning for DMAs without IO page fault

- Page fault for IO page fault support

- Notification
  - Page table content update notification
  - Page table root update notification



Attach/Detach Device to IOASID shared with TDP

QEMU

pin/unpin

IOMMU

Request sharing with pin mode & notification callbacks

KVM

IOASID

IO page fault    page fault

TDP

Page table content update notification

Page table root update notification

THE LINUX FOUNDATION

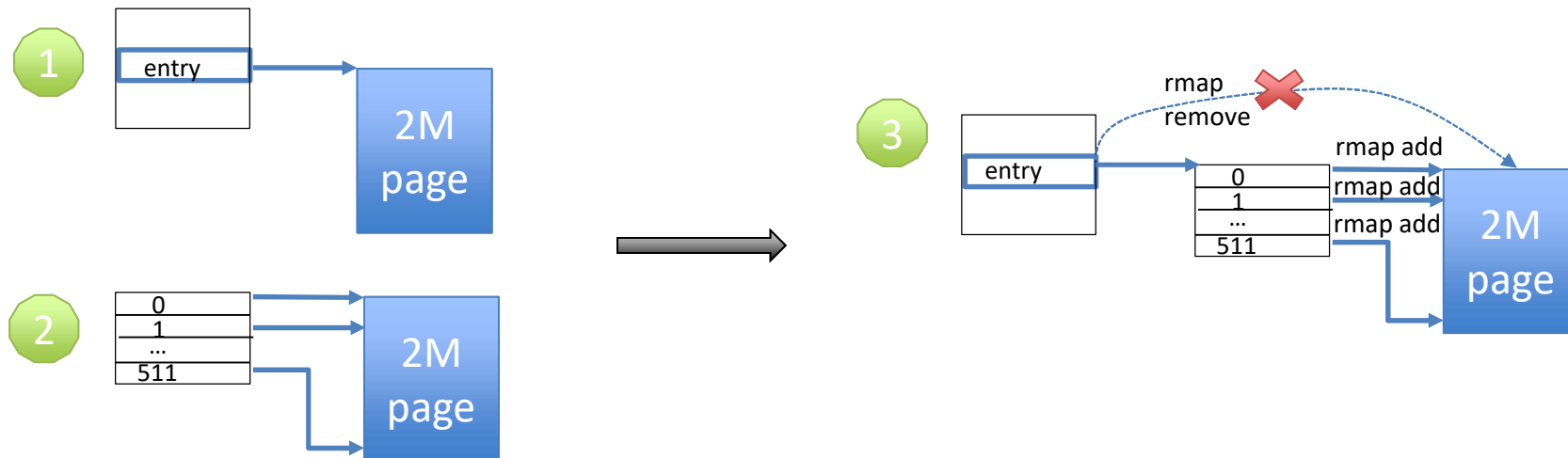# Page & Page Table Pinning

For sharing without IO page fault,

- Pinning of VM pages
  - pin_user_pages_* (FOLL_LONGTERM)

- Pinning of TDP entries
  - Pre-population of pinned ranges
  - No zap/pfn update
  - No reclaiming of mmu pages with parent linked
  - Atomic update of TDP entries when permission or page size change

pin_user_page_*

VM pages

KVM

Pin

TDP

zap/pfn_update

mm notifier

# Atomic Update for TDP Entries

Atomic update is required for TDP entries for pinned ranges, when

- Splitting huge pages
- Updating of PTE permission



TDP entry being atomically updated from non-zero value to another non-zero value.
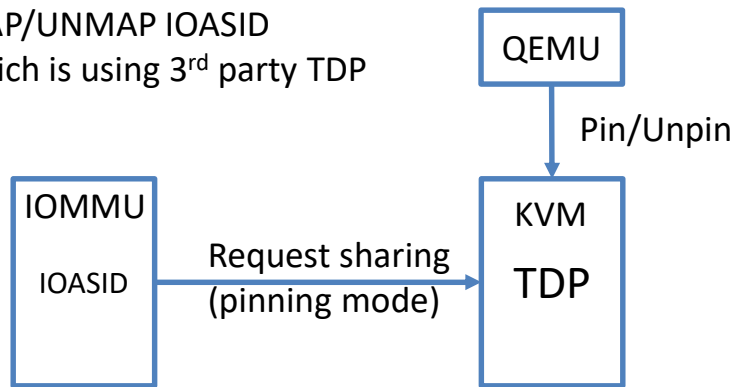
# Page & Page Table Pinning Interfaces

- For sharing without IO page fault,
  - Pinning of all ranges in user memslots: memslot add
  - Pinning a specific range: extra interface
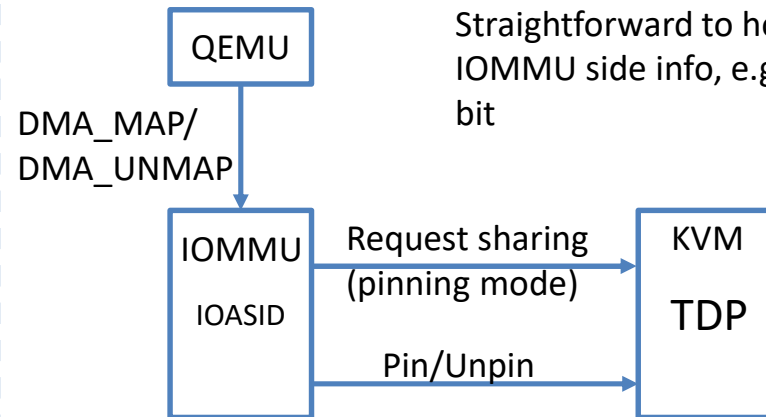
**Pin/Unpin from Qemu**
Pros:
QEMU doesn't need to
MAP/UNMAP IOASID
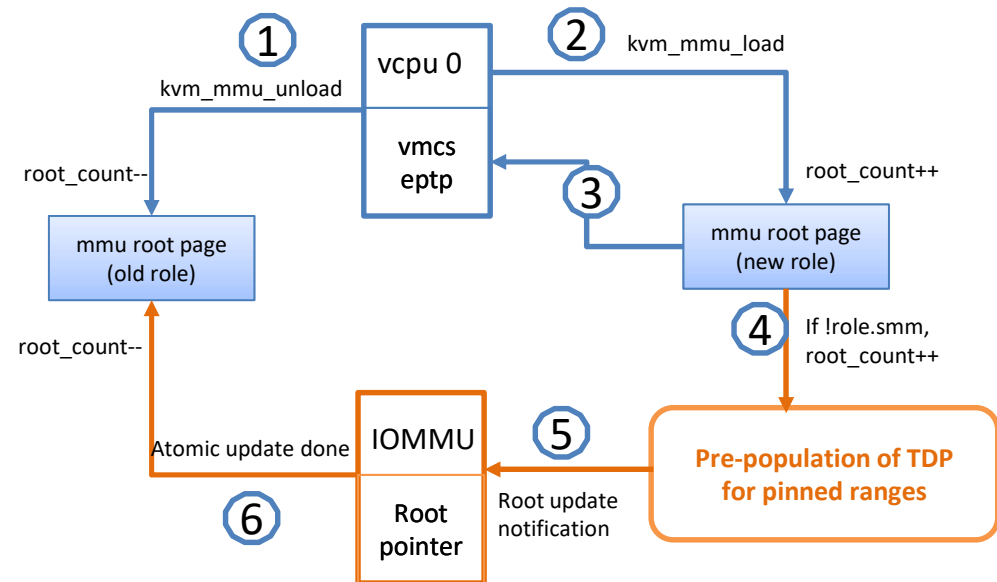which is using 3rd party TDP
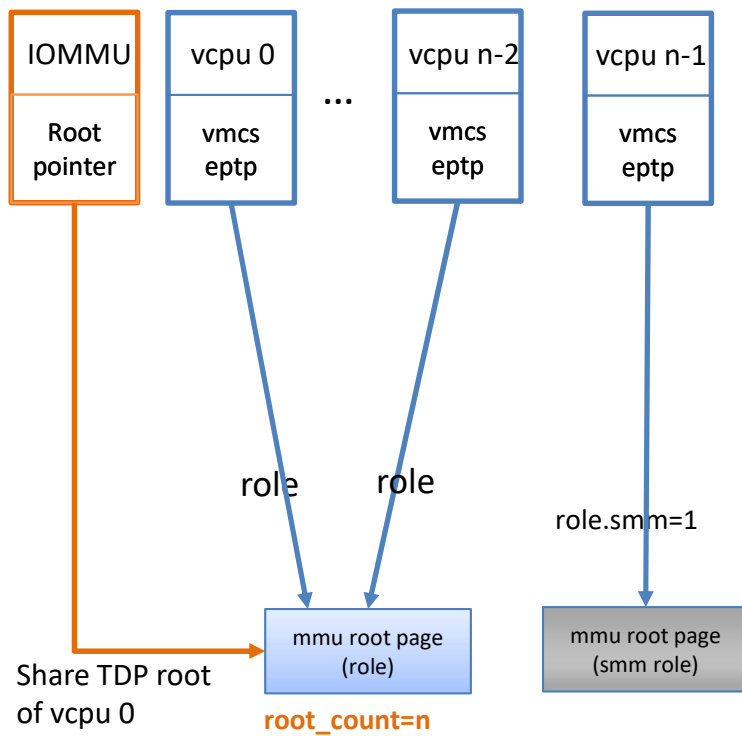
**Pin/Unpin from IOMMU**
Pros:
Straightforward to hold more
IOMMU side info, e.g. snoop
bit

# Shared Page Table Root Update

# Bootup Performance

Rough performance data without any optimization yet.

| 8G memory | Bootup Time | Pre-population count |
|---|---|---|
| Base (no sharing) | 29s | 0 |
| Sharing (huge page enabled) | 32s | 132 |
| Sharing (huge page disabled) | 63s | 132 |

- All VM pages were pinned/unpinned on user memslots creation/deletion.

- TDP was pre-populated on page table changes (when switching to new root, memslot add, and huge page splitting)

- IOTLB was flushed on page table root/content update notification (~1s)

- Quite a lot of time spend on TDP pre-population
  - ~2s with huge page
  - ~32s when huge page is disabled
- In concept can reach equal boot time performance as before sharing by reducing TDP root update count.

# TODOs

- Snoop bit handling

- Unified dirty page tracking

- Nested VM (vIOMMU, virtual EPT/NPT)

- Performance optimization
  - Page table root update reduction,
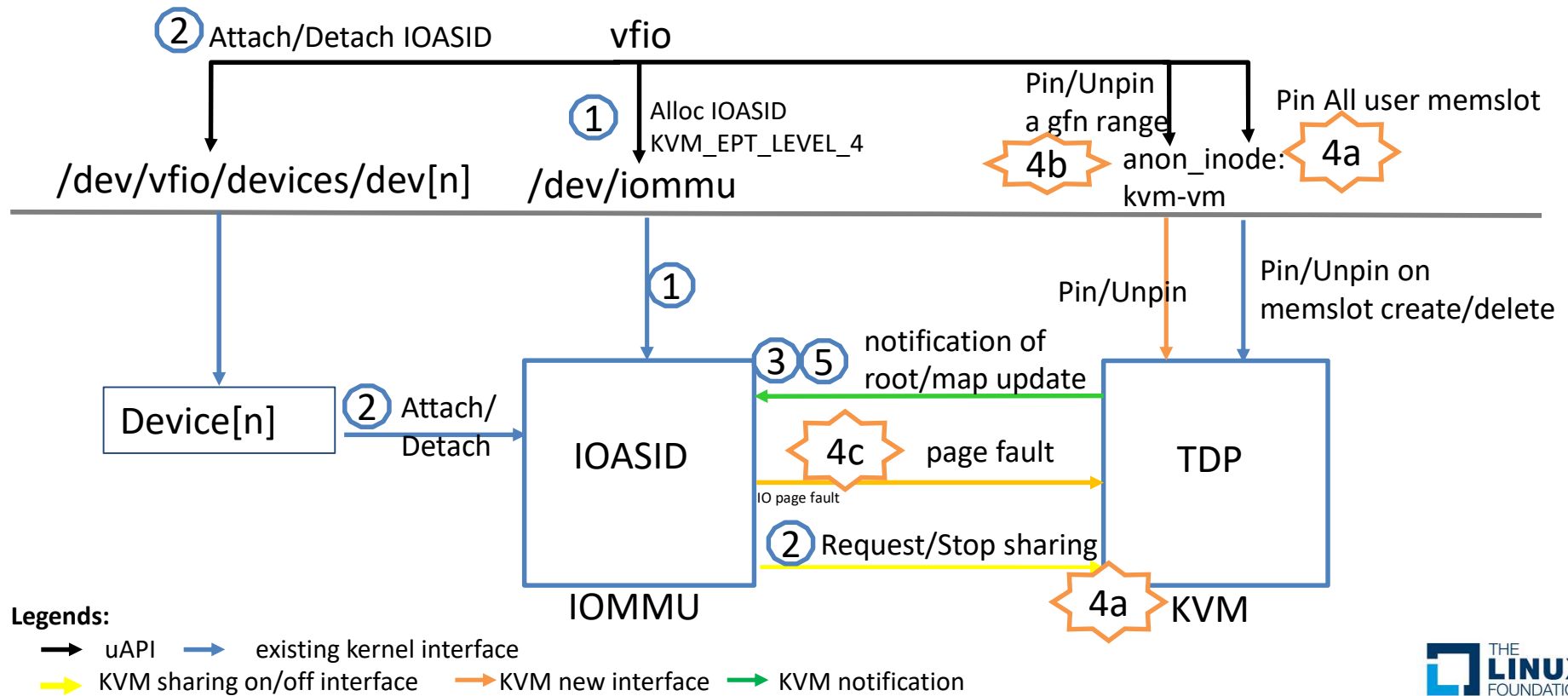  - Huge page support for P2P, etc.

# Why it is KVM manages the shared table

- CPU side has more restrictions in page size
  - Check guest MTRR
  - NX huge page workaround

- CPU side has extra GFN ranges to access
  - Private memslots in kernel space

- IOMMU page tables are not always present.

# Overall Design

# Overall Design (alternative)