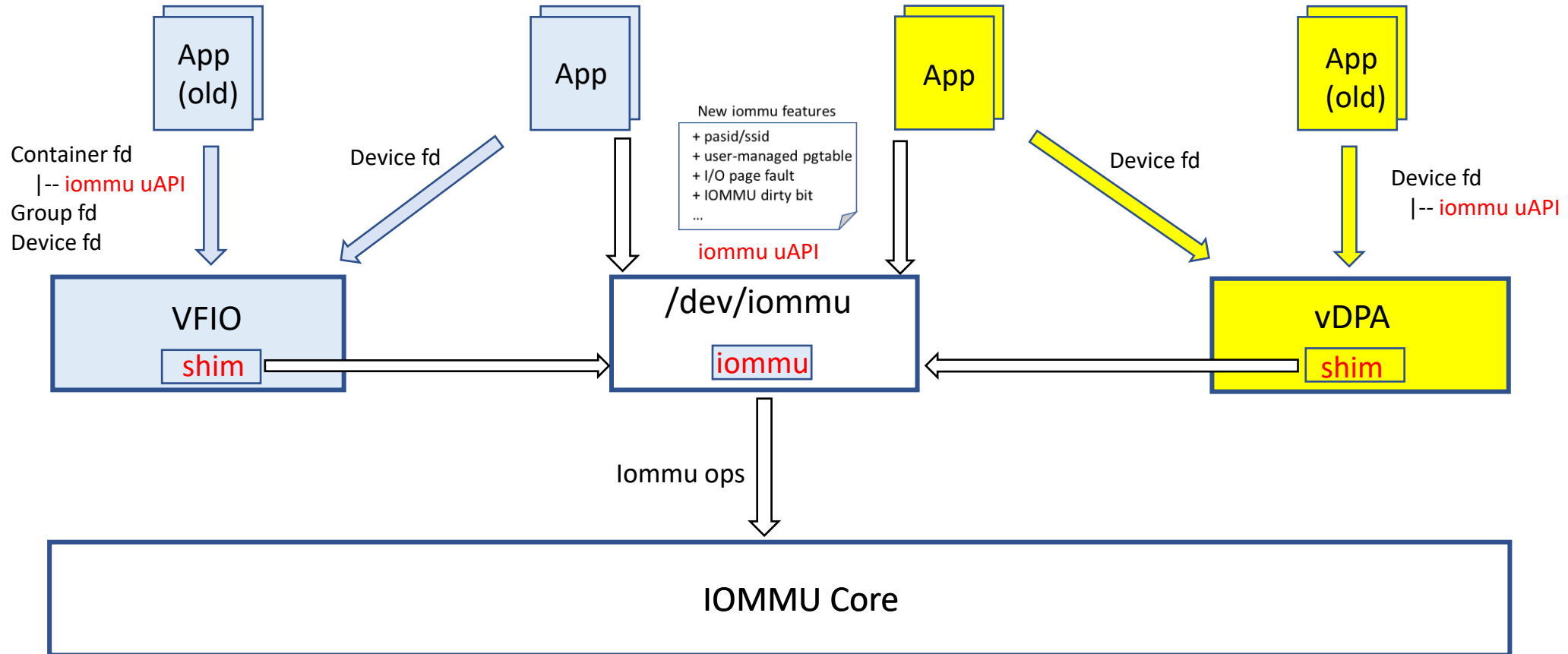# Unified I/O Page Table Management for Passthrough Devices

Kevin Tian, Baolu Lu

# Agenda

- A brief background

- Development plan

- Manage security context for user-initiated DMAs
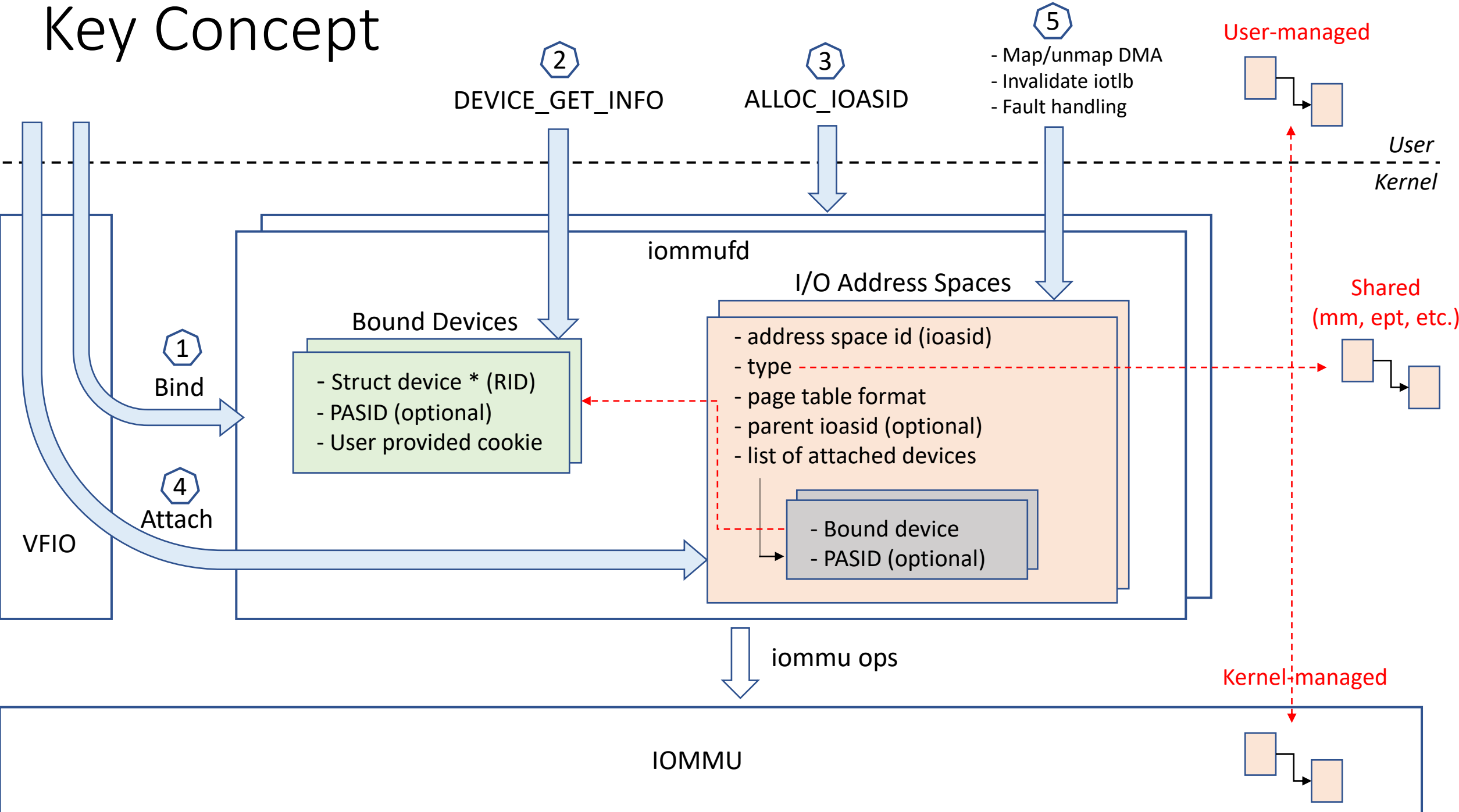
- Miscellaneous opens if time allows
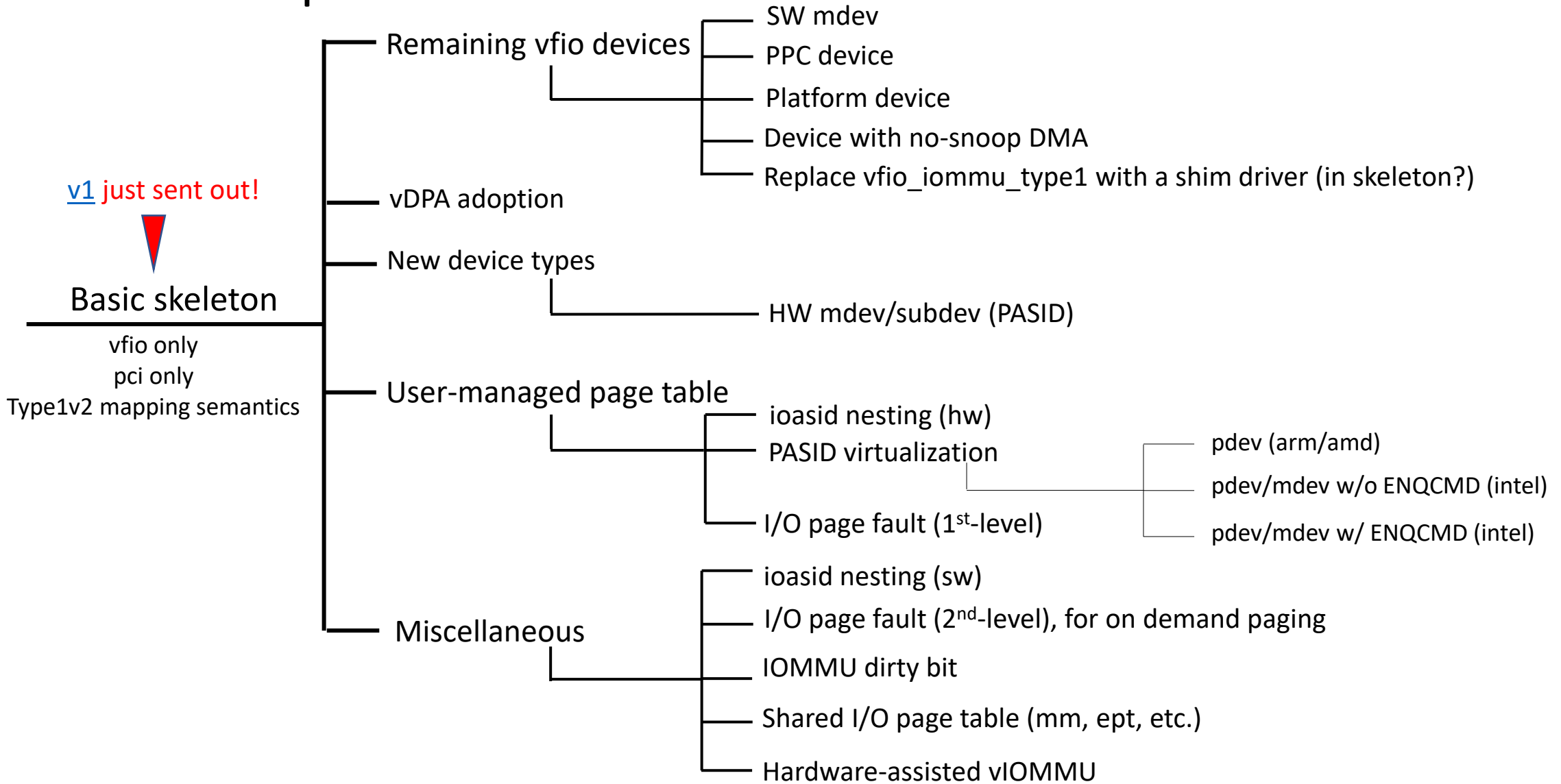
# Proposal for a Unified Framework



* Thank Jason Gunthorpe for initiating this idea!

* Please refer to [link] for a full design proposal

# Key Concept



② DEVICE_GET_INFO

③ ALLOC_IOASID

⑤
- Map/unmap DMA
- Invalidate iotlb
- Fault handling

User-managed

*User*
*Kernel*

iommufd

I/O Address Spaces

Bound Devices

① Bind

VFIO

- Struct device * (RID)
- PASID (optional)
- User provided cookie

④ Attach

- address space id (ioasid)
- type
- page table format
- parent ioasid (optional)
- list of attached devices

  - Bound device
  - PASID (optional)

Shared
(mm, ept, etc.)

iommu ops

Kernel-managed

IOMMU

# Development Plan

Basic skeleton
- vfio only
- pci only
- Type1v2 mapping semantics

v1 just sent out!

- Remaining vfio devices
  - SW mdev
  - PPC device
  - Platform device
  - Device with no-snoop DMA
  - Replace vfio_iommu_type1 with a shim driver (in skeleton?)
- vDPA adoption
- New device types
  - HW mdev/subdev (PASID)
- User-managed page table
  - ioasid nesting (hw)
  - PASID virtualization
    - pdev (arm/amd)
    - pdev/mdev w/o ENQCMD (intel)
    - pdev/mdev w/ ENQCMD (intel)
  - I/O page fault (1$^{st}$-level)
- Miscellaneous
  - ioasid nesting (sw)
  - I/O page fault (2$^{nd}$-level), for on demand paging
  - IOMMU dirty bit
  - Shared I/O page table (mm, ept, etc.)
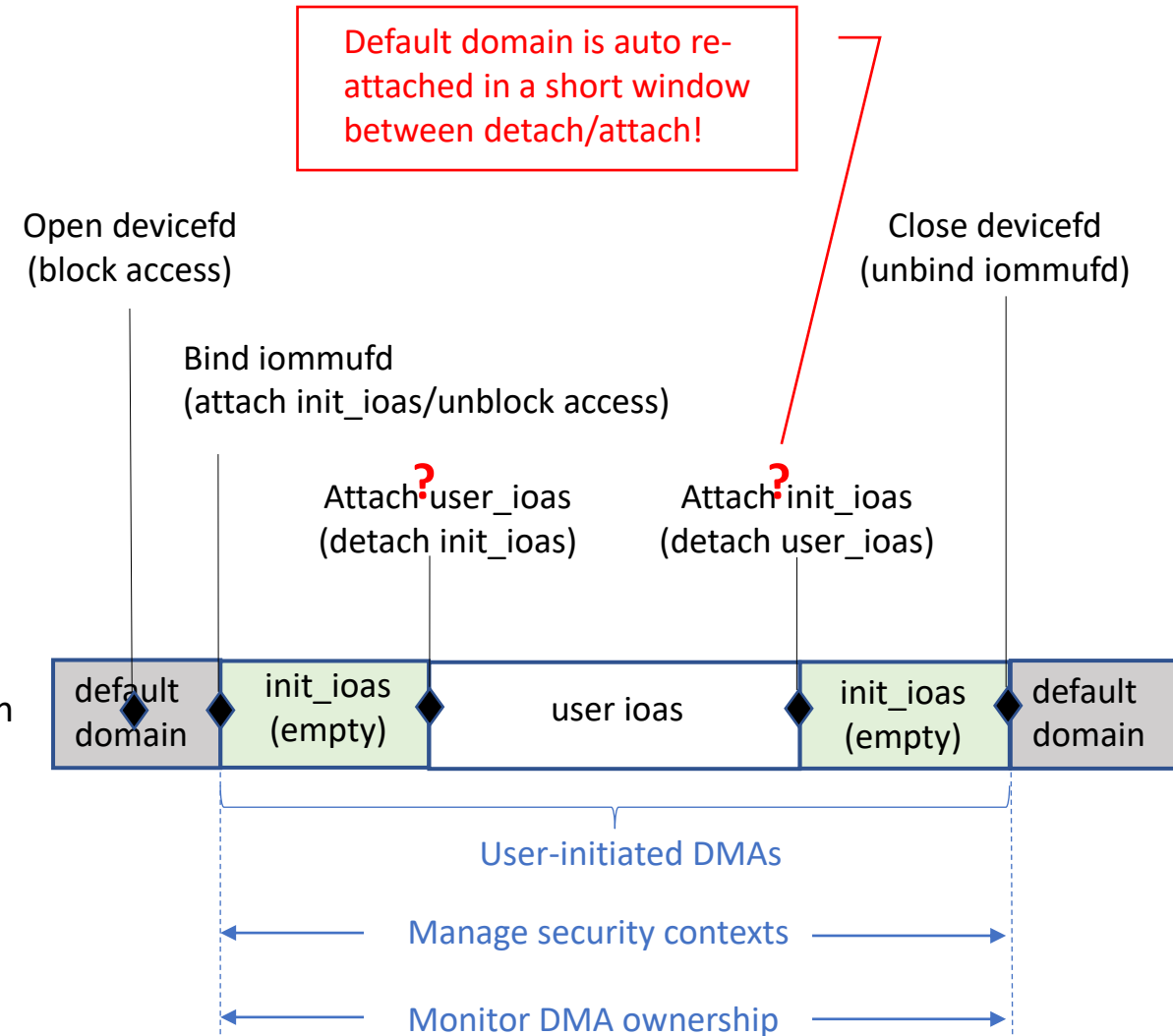  - Hardware-assisted vIOMMU

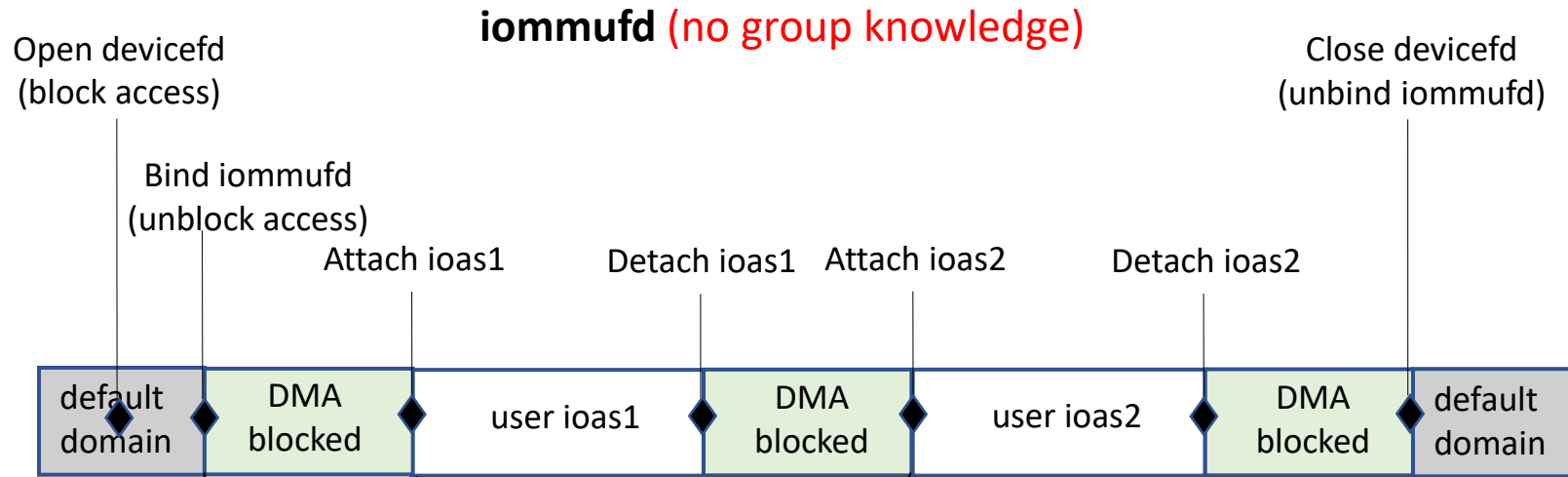Need community collaboration. Many tasks can be done in parallel!

# Manage User-initiated DMAs

- Establish a secure context (iommu unmanaged domain) which restricts user-initiated DMAs to
  - process memory
  - sibling devices in the same group

- Guarantee exclusive DMA ownership on the group, i.e. devices in the group must be
  - Bound to the owner driver (e.g. vfio), or
  - Bound to a driver known DMA-safe (e.g. pci-stub), or
  - Driverless

- iommufd can copy what vfio does today, with one exception
  - Need manage multiple security contexts due to decoupled bind/attach

- Current IOMMU API has problem on the transition between unmanaged domains
  - Default domain is automatically re-attached after detaching from previous context

Need cooperation from IOMMU core!

# Manage User-initiated DMAs (Cont.)

**iommufd** (no group knowledge)

Open devicefd
(block access)

Bind iommufd
(unblock access)

Close devicefd
(unbind iommufd)

Attach ioas1    Detach ioas1    Attach ioas2    Detach ioas2

| default domain | DMA blocked | user ioas1 | DMA blocked | user ioas2 | DMA blocked | default domain |

Device-centric IOMMU API

**IOMMU Core** (user-dma awareness)

iommu_device_init_user_dma():
- If the first device in the group
    * Validate and start monitoring group DMA ownership
    * Mark the group for user-dma
    * Block DMA for the entire group
- Else
    * refcount_inc(user_dma_cnt)

iommu_attach_device():
- If the first device in the group
    * Attach the group to ioas
- Else
    * Refcount_inc(attach_cnt)

iommu_detach_device():
- If the last device in the group
    * Detach the group from ioas
    * Block DMA instead of re-attaching to default domain
- Else
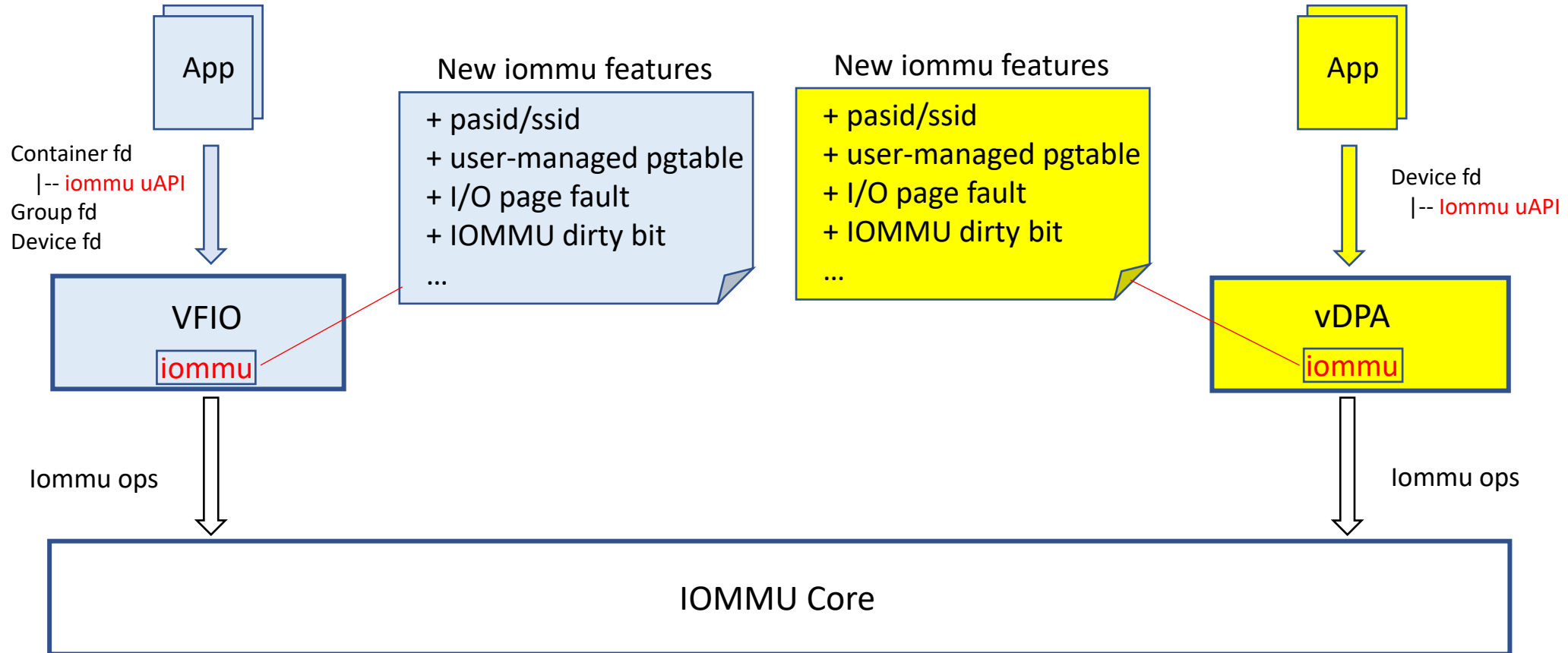    * Refcount_dec(attach_cnt)

iommu_device_exit_user_dma():
- If the last device in the group
    * Clear user-dma flag
    * Re-attach the group to default domain
    * Stop monitoring group DMA ownership
- Else
    * Refcount_dec(user_dma_cnt)

# If time allows

- ioasid naming conflict
  - fd-local software handle vs. hw asid (pasid/ssid, /drivers/iommu/ioasid.c)
- Module name and devnode
  - iommufd (/dev/iommu) vs. uiommu (/dev/uiommu)
  - Other options?
- /dev/vfio/devices/ hierarchy
  - A plain layout mixing all types together
    - /dev/vfio/devices/0000:00:14.2 (pci)
    - /dev/vfio/devices/PNP0103:00 (platform)
    - /dev/vfio/devices/83b8f4f2-509f-382f-3c1e-e6bfe0fa1001 (mdev)
  - Subdirectories based on device types
    - pci, platform, ccw, etc.
    - Pdev vs. mdev
- Do we need to build iommufd as a separated module?
- Convert vfio_iommu_type1 to a shim driver

# Backup

# Current Situation



App

Container fd
|-- iommu uAPI
Group fd
Device fd

New iommu features

+ pasid/ssid
+ user-managed pgtable
+ I/O page fault
+ IOMMU dirty bit
...

New iommu features

+ pasid/ssid
+ user-managed pgtable
+ I/O page fault
+ IOMMU dirty bit
...

App

Device fd
|-- Iommu uAPI

VFIO

iommu

vDPA

iommu

Iommu ops

Iommu ops

IOMMU Core

Not a scalable architecture moving forward!

# Manage User-initiated DMAs

**Current VFIO**