



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

CXL 2.0 + Linux + QEMU = Yes



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

CXL 2.0 + Linux + QEMU = ~~Yes~~



**LINUX** September 20-24, 2021  
**PLUMBERS**  
**CONFERENCE**

# Introduction

Last slide first!!!



JAKE-CLARK.TUMBLR



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

# Agenda

- Introduction
- CXL 2.0 Background
  - Linux Driver Details
- QEMU
- Future



# Intro/Links

- Communications
  - #cxl on OFTC
  - [linux-cxl@vger.kernel.org](mailto:linux-cxl@vger.kernel.org)
- Drivers
  - <https://git.kernel.org/pub/scm/linux/kernel/git/cxl/cxl.git/>
- QEMU
  - <https://gitlab.com/bwidawsk/qemu>
  - [https://github.com/pmempv/run\\_qemu](https://github.com/pmempv/run_qemu)
- Userspace
  - <https://github.com/pmempv/ndctl/tree/cxl-2.0v3>
  - [https://gitlab.com/bwidawsk-cxl/cxl\\_rs](https://gitlab.com/bwidawsk-cxl/cxl_rs)



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

# CXL 2.0 Background



# CXL Features

## Challenges

Industry trends driving demand for faster data processing and next-gen data center performance

Increasing demand for heterogeneous computing and server disaggregation

Need for increased memory capacity and bandwidth

Lack of open industry standard to address next-gen interconnect challenges

## CXL

An open  
industry-supported  
cache-coherent  
interconnect for  
processors, memory  
expansion and  
accelerators

## Coherent Interface

Leverages PCIe® with 3  
mix-and-match protocols

## Low Latency

.Cache and .Memory targeted at  
near CPU cache coherent latency

## Asymmetric Complexity

Eases burdens of cache coherent  
interface designs



# CXL Features

## Challenges

Industry trends driving demand for faster data processing and next-gen data center performance

Increasing demand for heterogeneous computing and server disaggregation

Need for increased memory capacity and bandwidth

Lack of open industry standard to address next-gen interconnect challenges

**CXL**  
An open  
industry-supported  
cache-coherent  
interconnect for  
processors, memory  
expansion and  
accelerators

## Coherent Interface

Leverages PCIe® with 3 mix-and-match protocols

## Low Latency

.Cache and .Memory targeted at near CPU cache coherent latency

## Asymmetric Complexity

Eases burdens of cache coherent interface designs





**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**



arm



facebook



Google



intel®

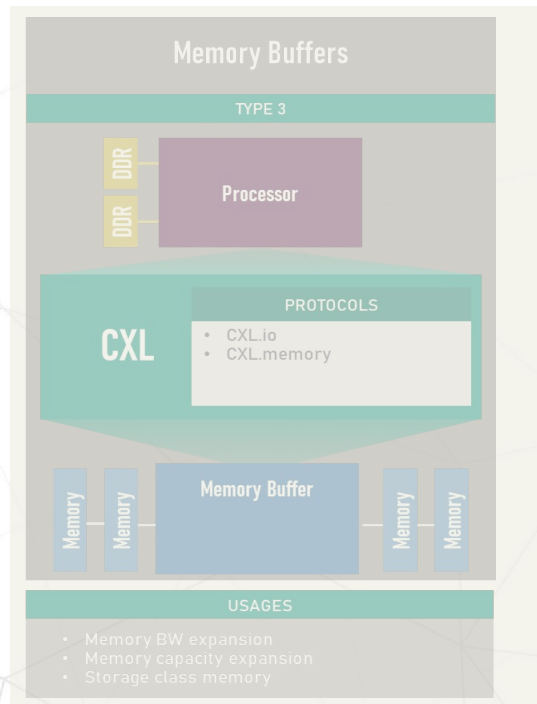
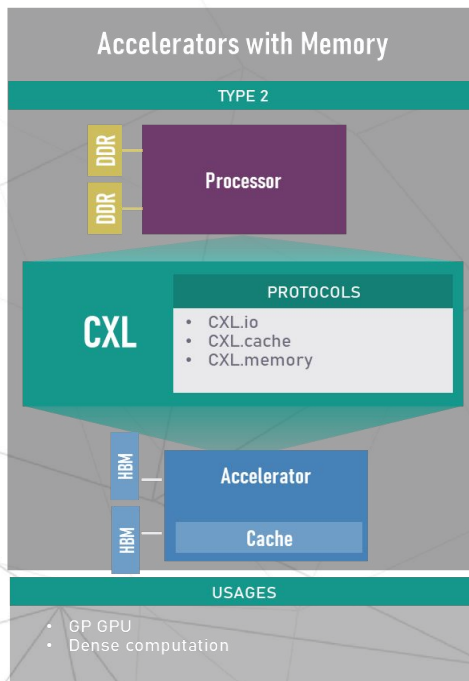
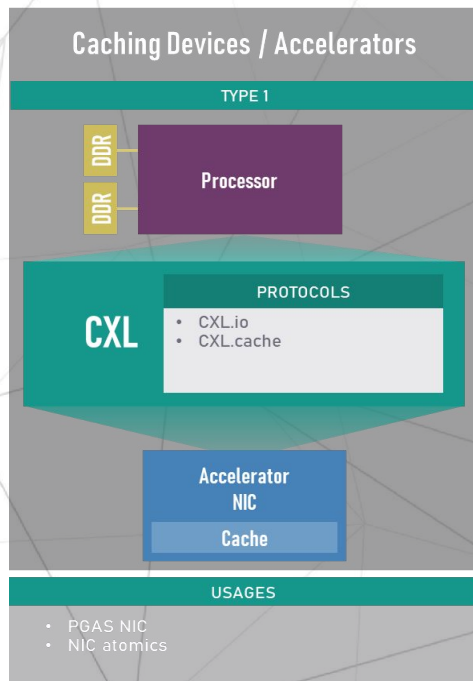


CXL Board of Directors

**CXL** / Compute  
Express  
Link™

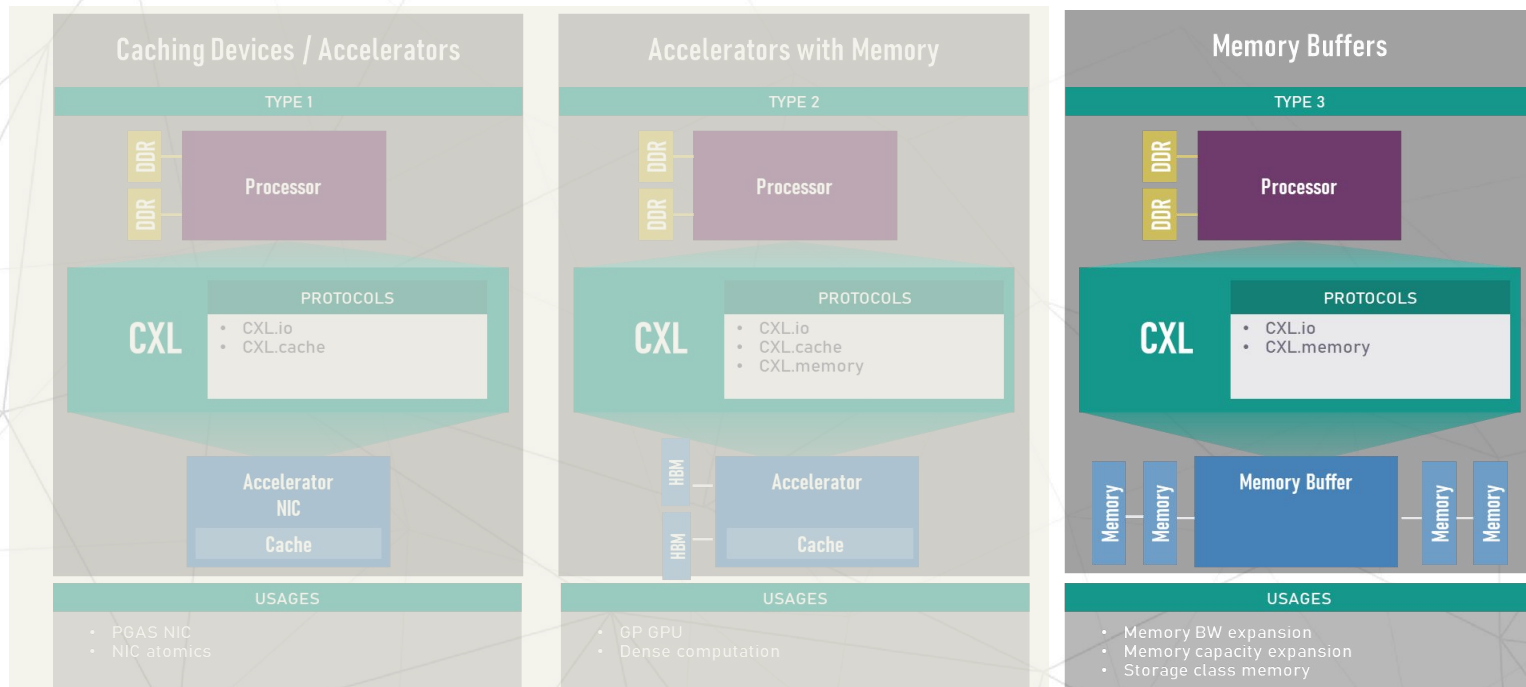


# Usage Models





# Usage Models



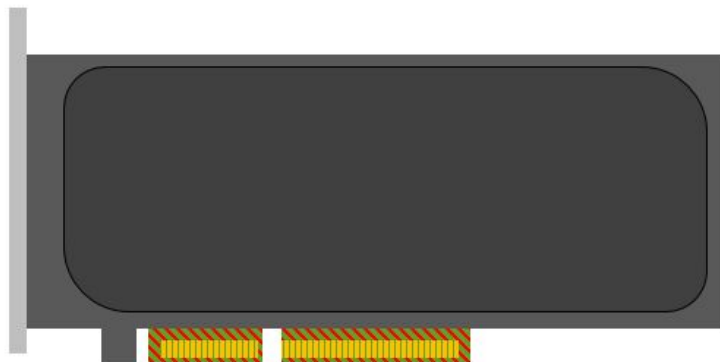
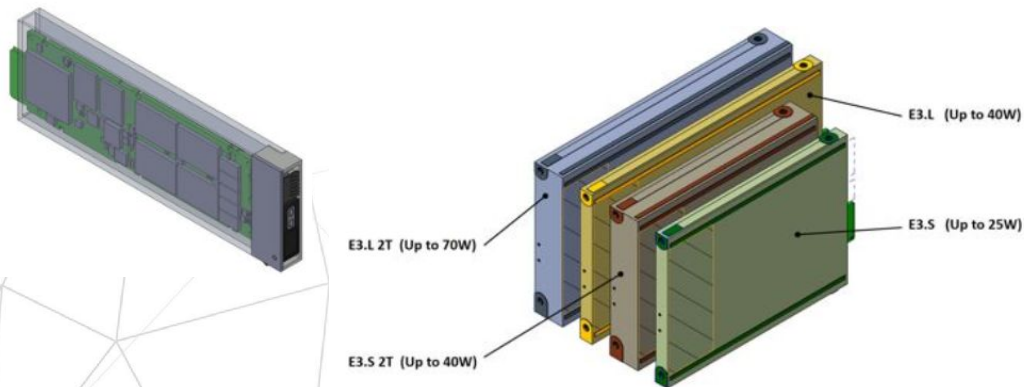


# Form Factors

Ex #1: EDSFF E1.S

Ex #2: EDSFF E3.S / E3.L

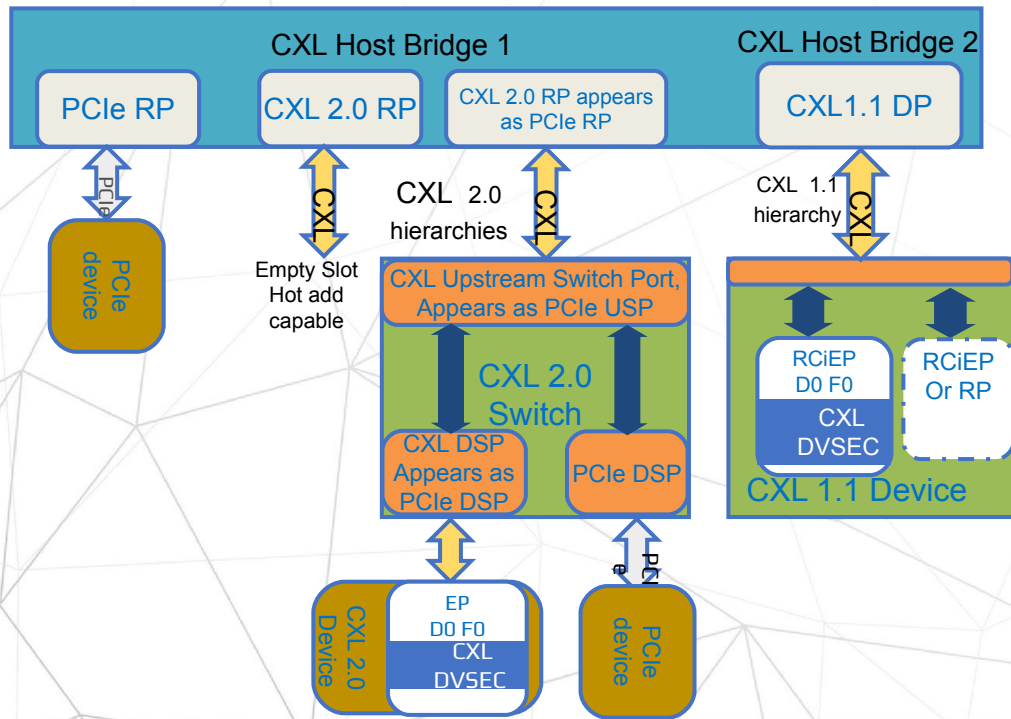
Ex #3: Add-in Card (AIC)



Factors	EDSFF E1.S	EDSFF E3.S / E3.L	AIC
Area vs. DDRx Server DIMM	<ul style="list-style-type: none"><li>• Smaller</li></ul>	<ul style="list-style-type: none"><li>• Larger</li></ul>	<ul style="list-style-type: none"><li>• Larger (larger than E3.S/L)</li></ul>
Expected Max. Power Range	<ul style="list-style-type: none"><li>• 12 ~ 25W</li></ul>	<ul style="list-style-type: none"><li>• 25W ~ 40W (1T), 40W ~ 70W (2T)</li></ul>	<ul style="list-style-type: none"><li>• Similar range compared to E3.S/L</li></ul>



# CXL Topology

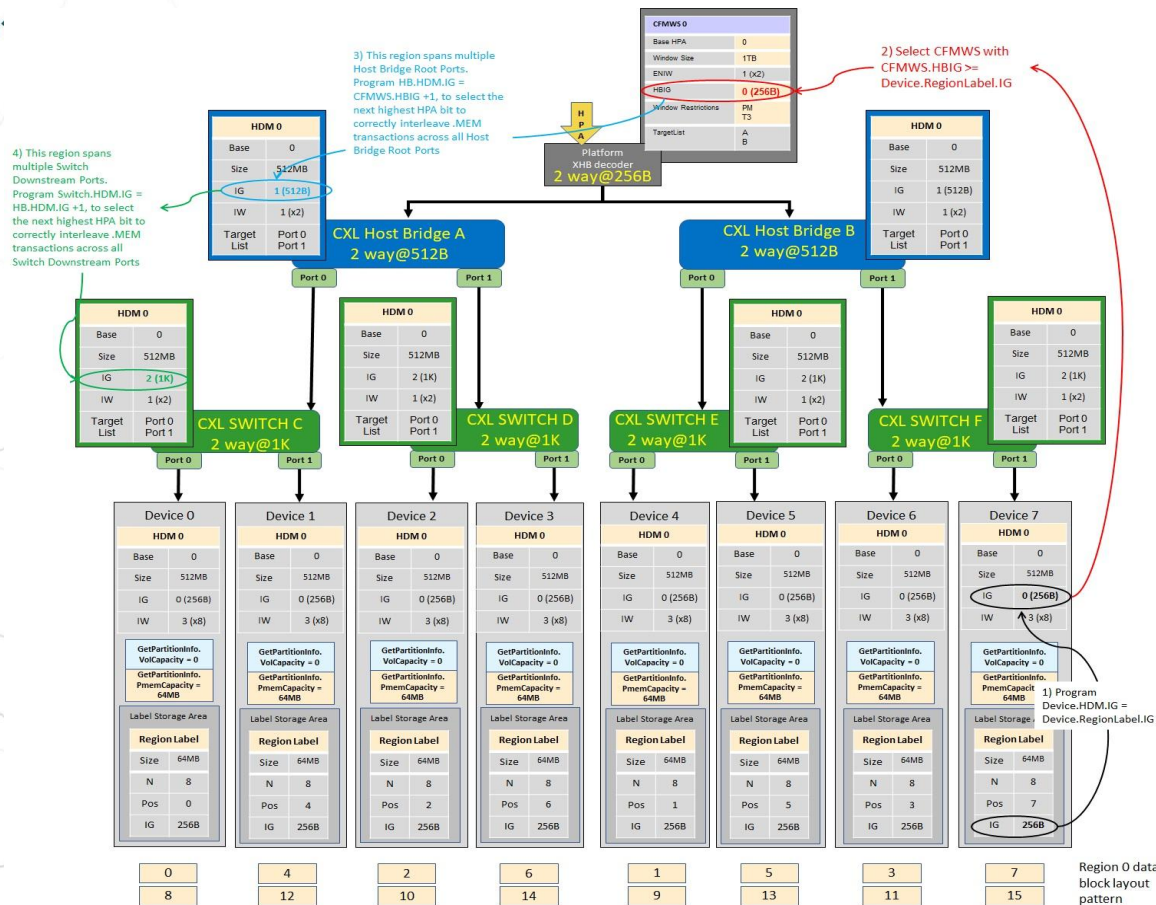


- CXL 2.0 hierarchy appears like PCIe hierarchy
  - Legacy PCI SW and CXL SW sees a RP or DSP with Endpoints below
  - CXL link/interface errors are signaled to RP, not RCEC
  - Port Control Override registers prevent legacy PCIe software from unintentionally resetting the device and the link
- Interleaving
  - Cross host bridge
  - Switch
  - Device



# LINUX September 20-24, 2021 PLUMBERS CONFERENCE

## Interleave

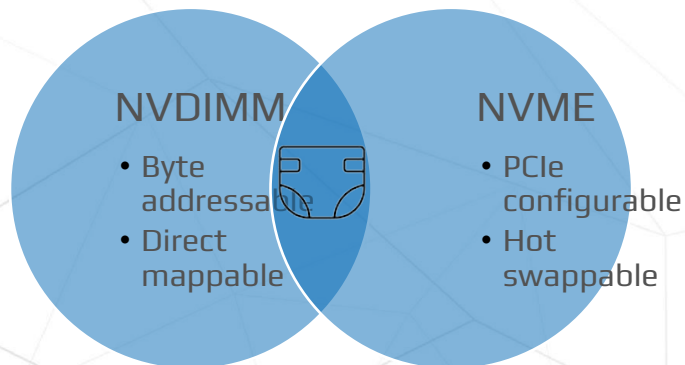






# CXL Persistent Memory Devices

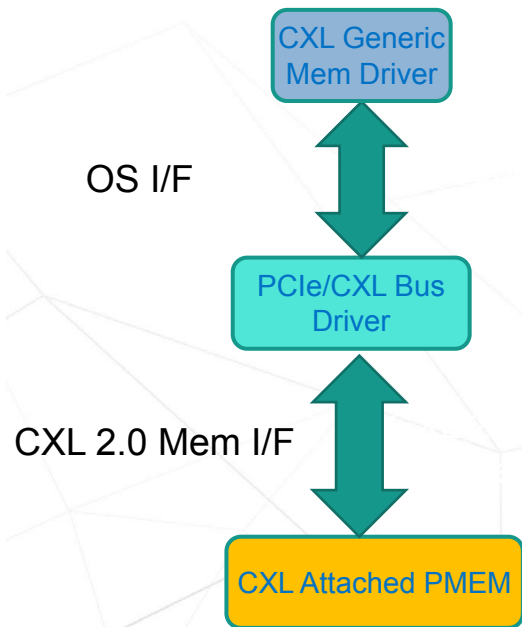
PCI and NVDIMM had a coherent byte addressable baby... with atomics.





- Persistent memory devices rely on System Software for provisioning and management
- CXL 2.0 introduces a standard register interface
- A generic memory device driver simplifies software enabling
- Architecture Elements
  - Defined as number of discoverable Capabilities
  - Capabilities includes Device Status and standard mailboxes, accessed via MMIO registers
  - Standardized mailbox commands that cover errors/health, alerts, partitioning, passphrases etc.
  - Allow Vendor specific extensions

# PMEM





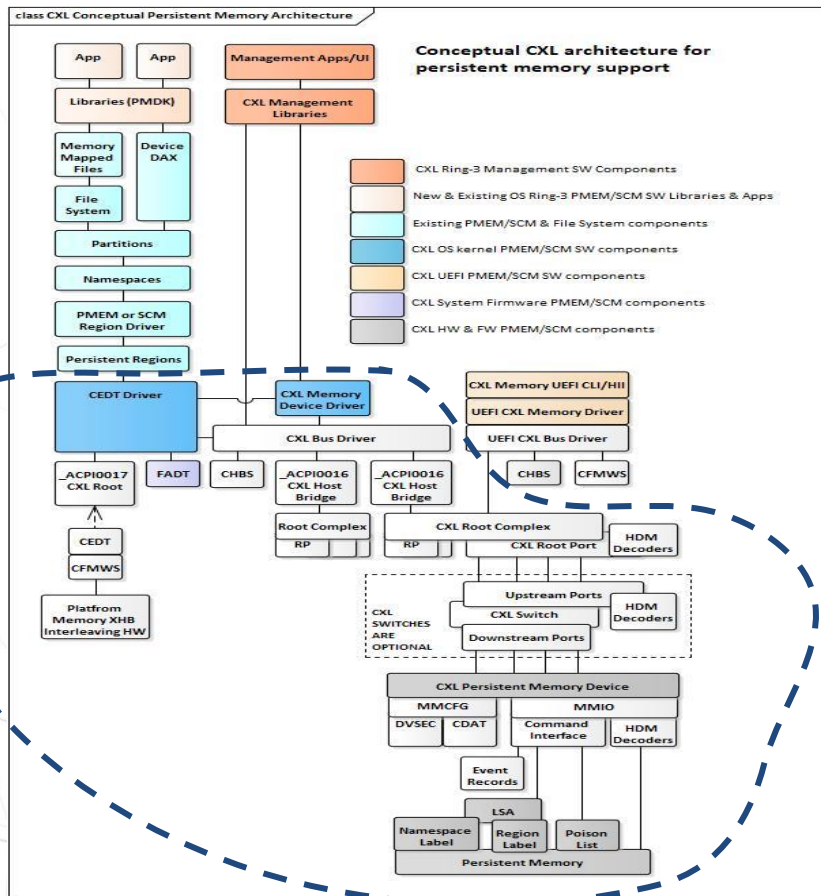


**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

# Linux Driver Details

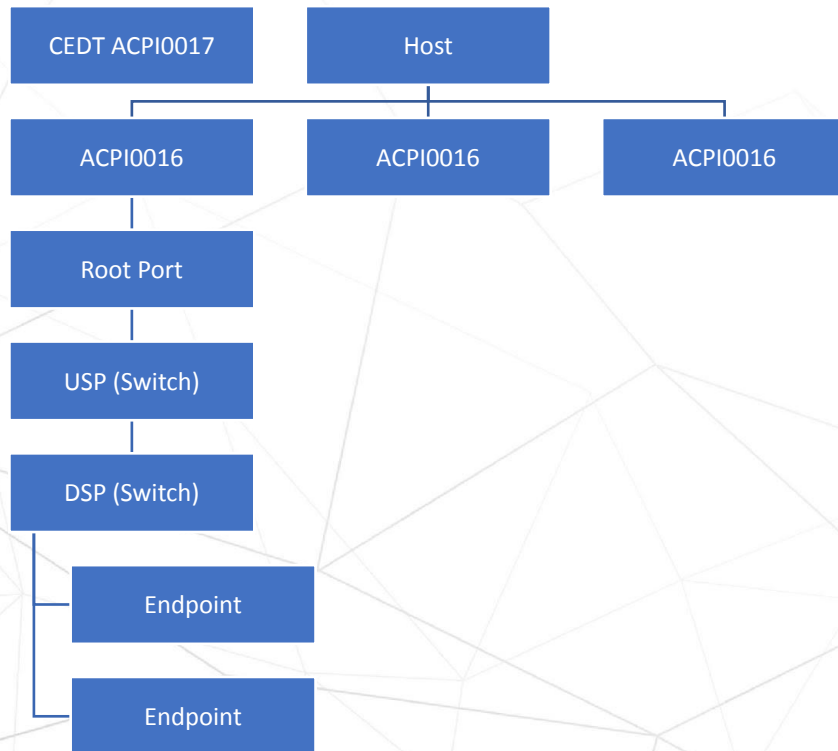


# Software Responsibilities



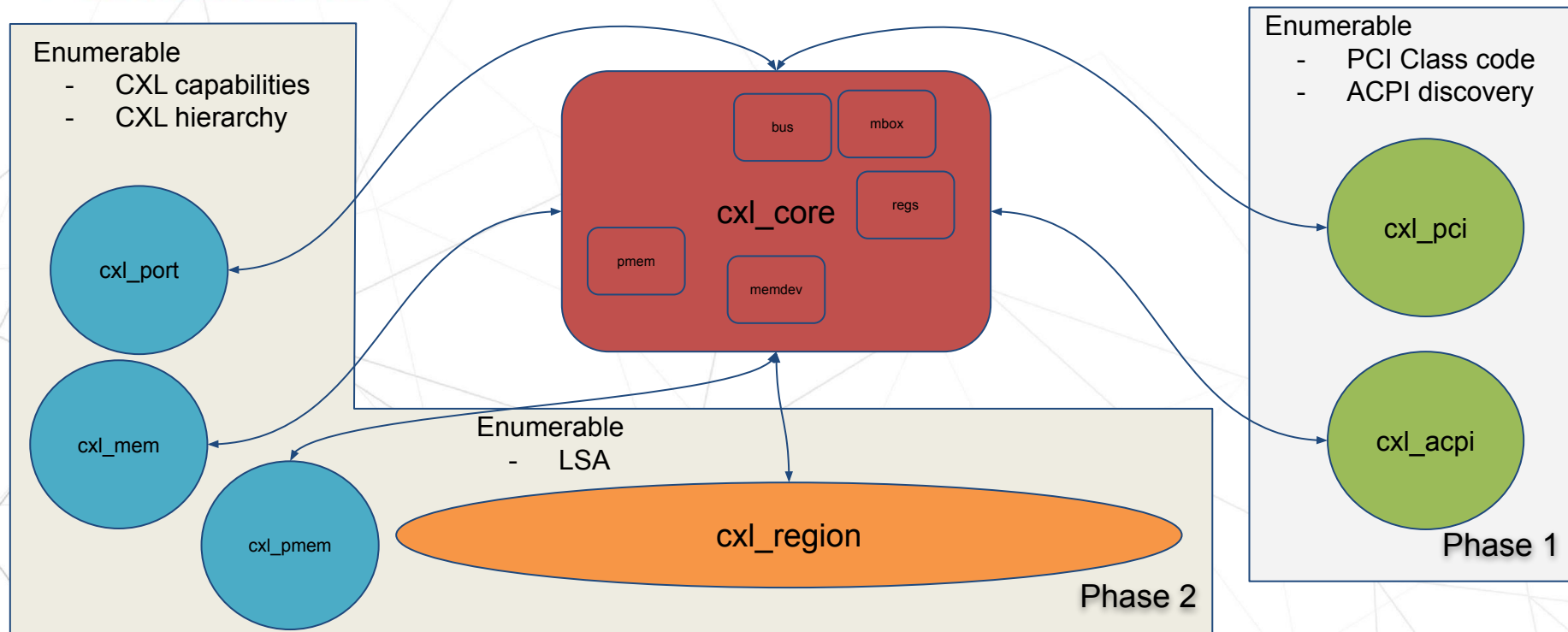


# SW Enumerable Components





# Linux Drivers





# cxl\_core

- Maintains cxl\_driver infra
- Interfaces with LIBNVDIMM
- Manages device (sysfs)
  - Services to add devices, ie. cxl\_decoder\_add()
- IOCTL interface
- Common functionality
  - Mailbox controls (session layer)
  - Register mapping

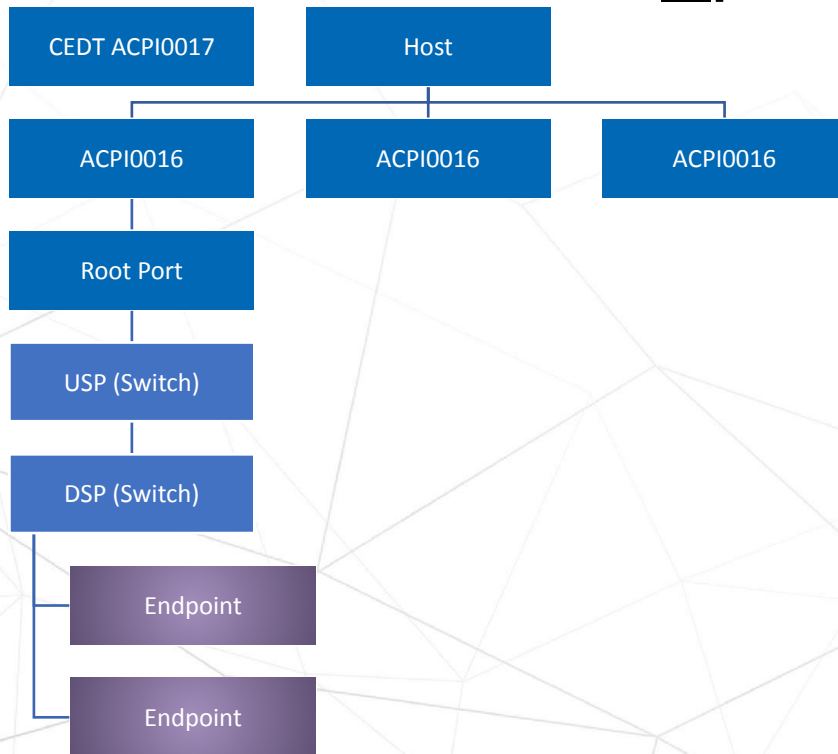


## cxl\_pci

- Probed like a typical PCI device
  - { PCI\_DEVICE\_CLASS((PCI\_CLASS\_MEMORY\_CXL << 8 | CXL\_MEMORY\_PROGIF), ~0)}
- CXL device manageability
  - Implements mailbox transport (CXL) protocol
- Enumerates CXL device for subsequent driver
  - cxl\_mem can't run until cxl\_pci is done
- Attestation/Security/Whatever



# cxl\_pci





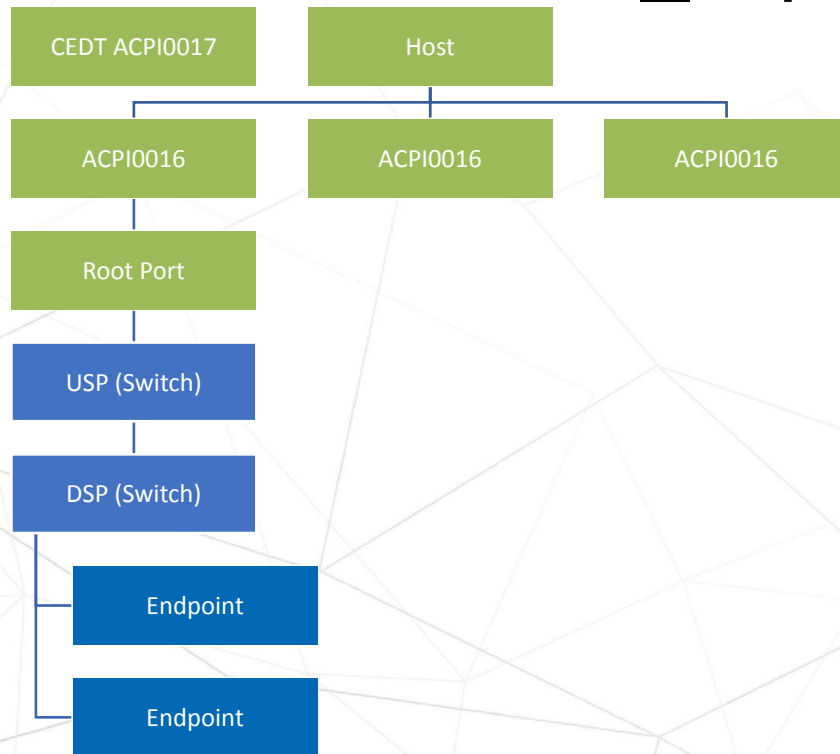
# cxl\_acpi

- Probed like a typical ACPI device
  - { "ACPI0017", (unsigned long) &native\_acpi0017 },
- Platform specific CXL enumeration
  - Mostly specified in UEFI, and CXL
- ACPI0017 starts enumeration of CXL ports
  - CEDT
  - “Root level” ports (platform)
  - Hostbridges and root ports





# cxl\_acpi



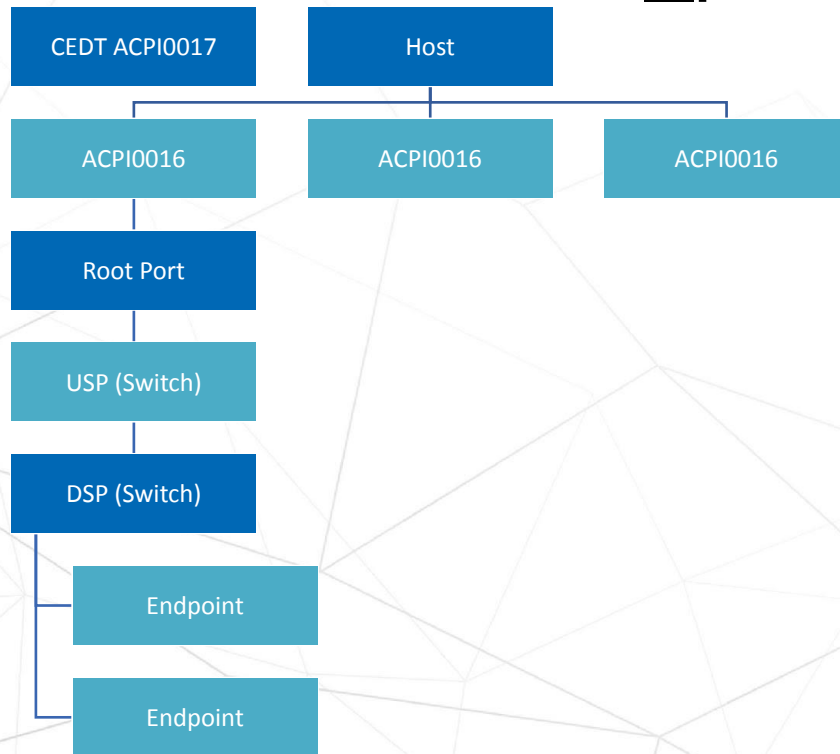


# cxl\_port

- Ports are created for all components with an “upstream port”
  - Hostbridge
  - Switch
  - Endpoint
- Enumeration and control and control of decoder resources
  - Provides as a service for other drivers



# cxl\_port



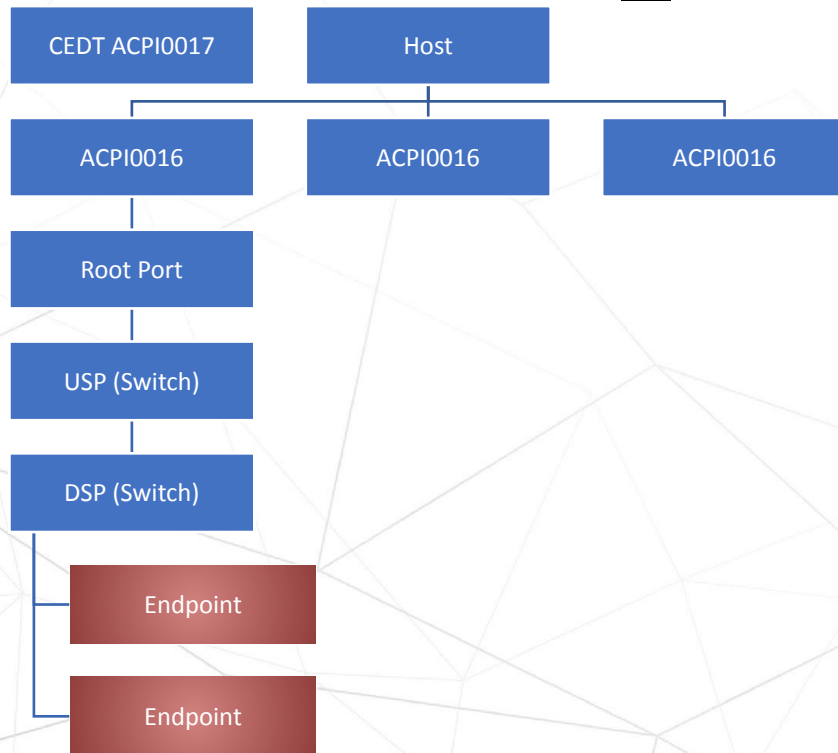


# cxl\_mem

- connects a device enumerated with cxl\_pci to CXL.mem.
  - Implements device functionality not handled by cxl\_pci
- “exports” if device is CXL.mem routed and enabled

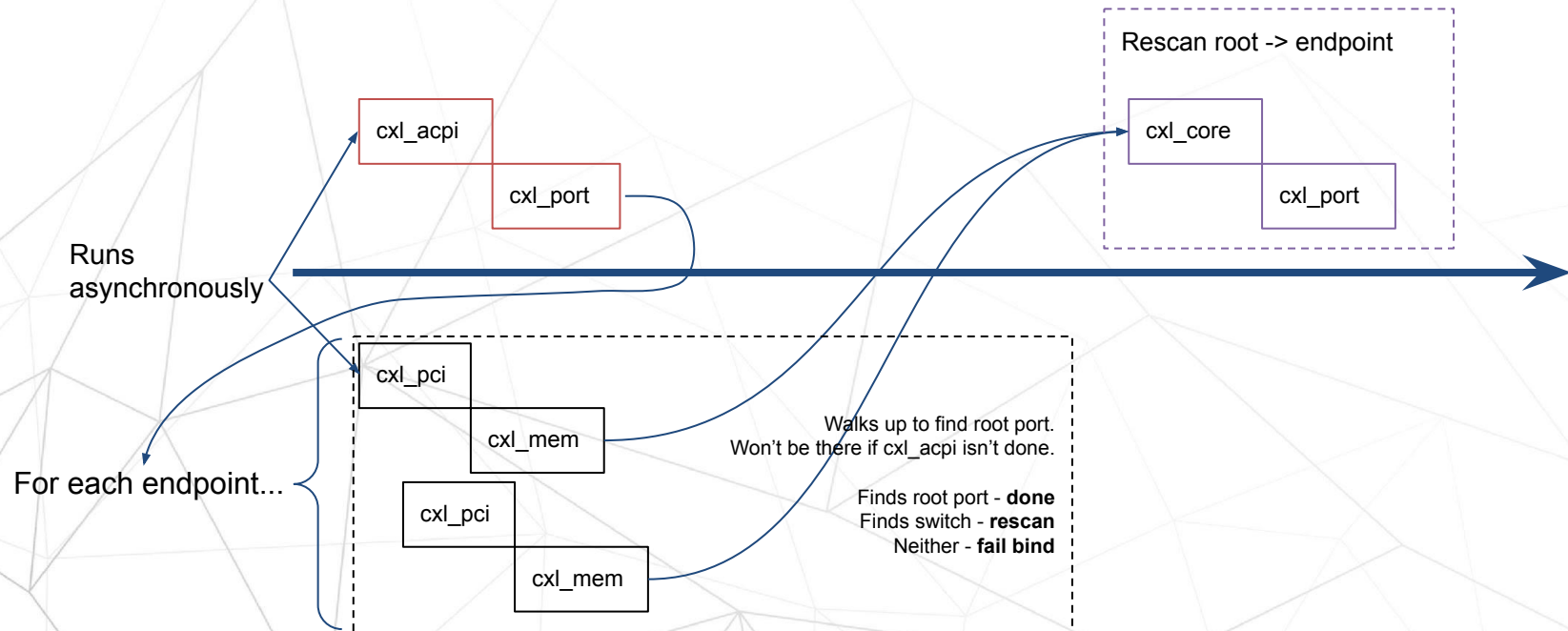


# cxl\_mem





# Enumeration Timeline





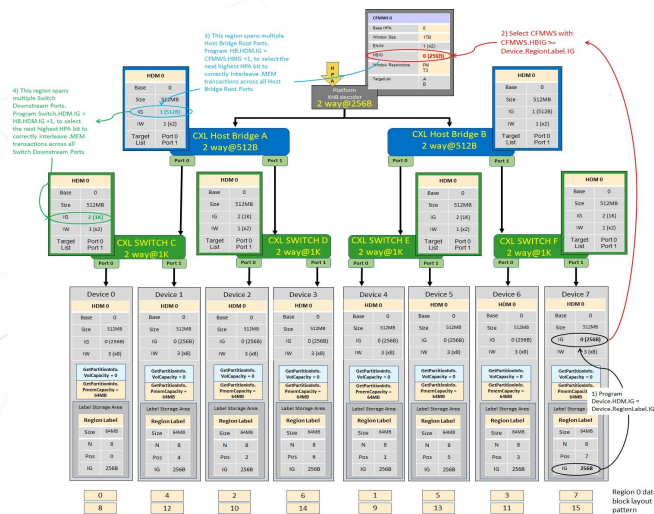
# Volatile vs. Persistent

		Persistent	Volatile	
1	Configured by BIOS	✗	✓	<ul style="list-style-type: none"><li>• BIOS configures all volatile capacities</li><li>• BIOS may check PMEM devices, but will not configure regions.<ul style="list-style-type: none"><li>◦ BIOS may configure bootable PMEM</li></ul></li></ul>
2	OS managed	✓	😊	<ul style="list-style-type: none"><li>• OS initializes persistent regions</li><li>• OS may create new persistent regions</li><li>• Manages hotplug, error, and reset sequences for both</li></ul>
3	Requires CXL.mem capability	✓	✓	



# Regions

- Region
  - Interleave set of devices
  - Parameters (IG, HPA, etc)
  - Stored in the Label Storage Area
- Creation
  - Via sysfs ABI
  - Provisioned offline
    - Manufacture time
- Responsibilities
  - Validating region configuration
  - Programming HDM decoders





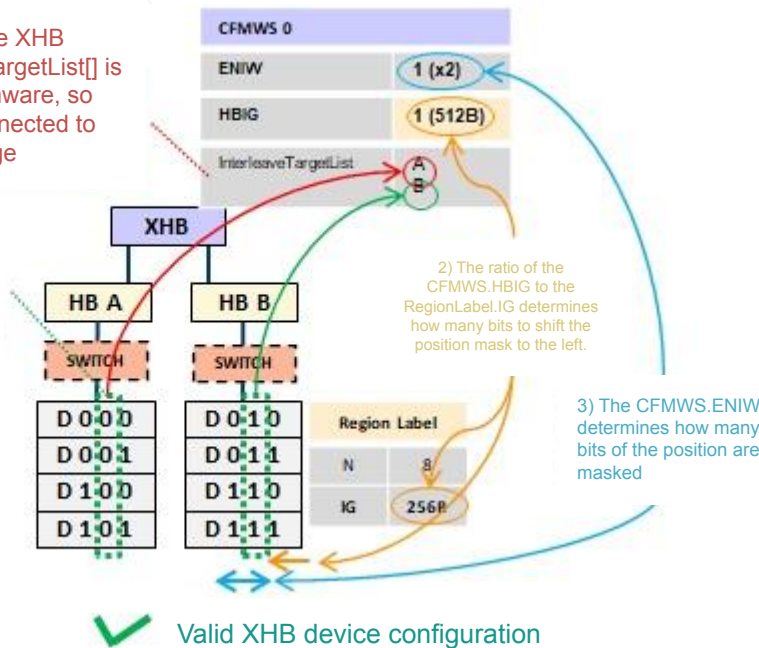


# Region Validation

1) The ordering of the XHB CFMWS.InterleaveTargetList[] is fixed by System Firmware, so devices must be connected to the correct host bridge

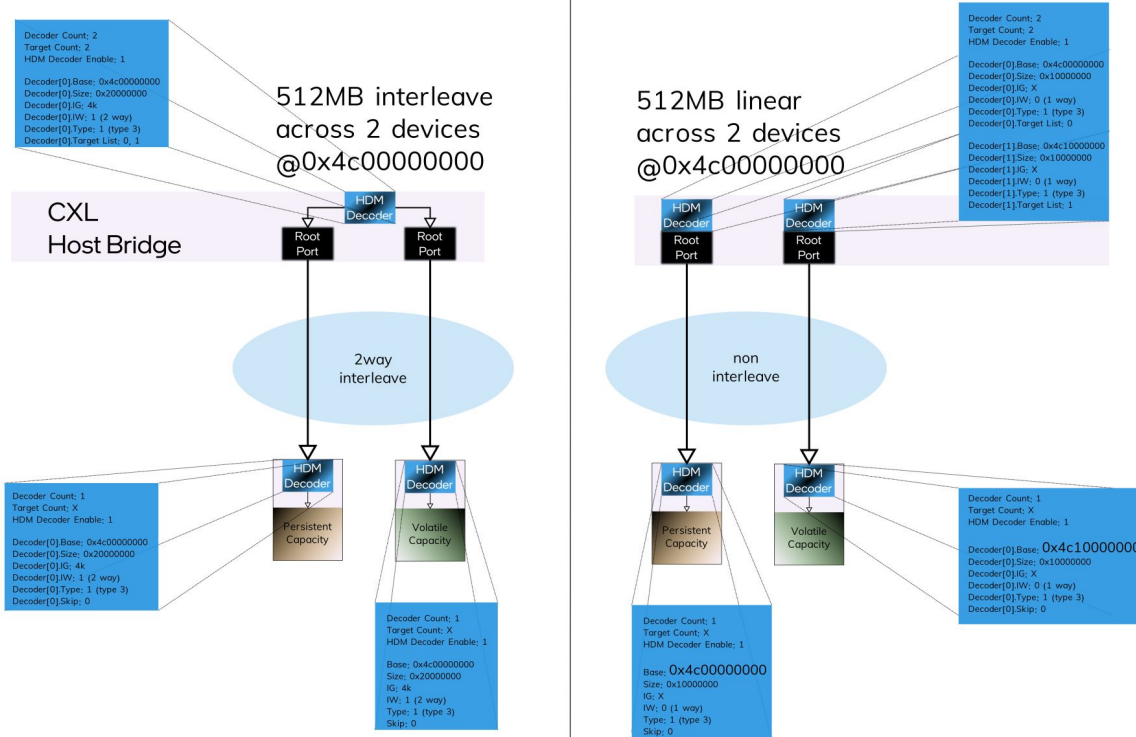
Verify each device has the same masked position value

$\text{Device}[x].\text{RegionLabel.Position} \gg (\text{CFMWS.HBIG} - \text{Device}[x].\text{RegionLabel.IG}) \& ((2^{\text{CFMWS.ENIW}} - 1))$





# Decoder Programming





# Linux Interfaces

- **Sysfs**
  - `/sys/bus/cxl/`
- **IOCTL**
  - QUERY
  - SEND
  - Managed command set
  - RAW escape command
- **Future**
  - Region Creation ABI



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

QEMU



# Review Goals

- 1) Upstream Linux Driver
  - a) 0 days of spec release (v5.12)
    - i) Ease customer adoption
  - b) Backportable
  - c) Platform aiding hw implementation and validation
    - i) Validate the spec for driver usage
- 2) Reusable past driver bringup
  - a) infra for regression testing
  - b) Virtualization
- 3) Scalable
  - a) Community Contributions & Fixes



# Pre-silicon State of the Art

- Hardware
  - No 2.0 FPGAs available
  - 1.1 is limited use and not readily available.
- Internal Simulation
  - Delays/
  - Can't work with community
- Prior art
  - nfit\_test
  - QEMU CCIX patches



# Pre-silicon State of the Art

- Hardware
  - No 2.0 FPGAs available
  - 1.1 is limited use and not readily available.
- Internal Simulation
  - Delays/
  - Can't work with community
- Prior art
  - nfit\_test
  - QEMU CCIX patches

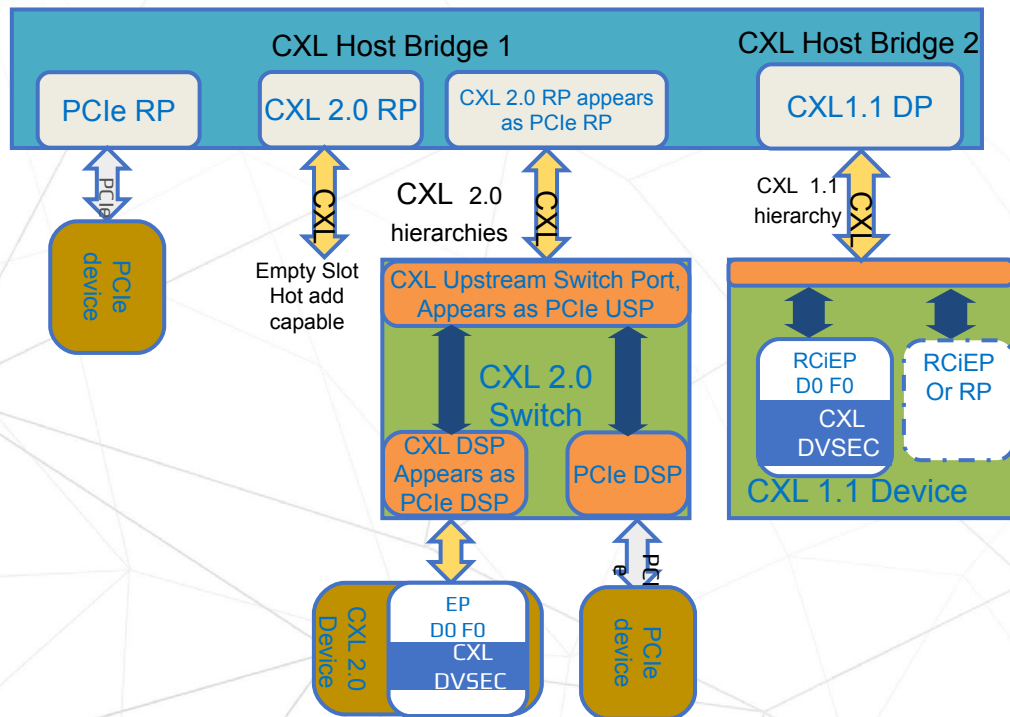


[This Photo](#) by Unknown Author is licensed under [CC BY-SA](#)





# CXL Arch Review



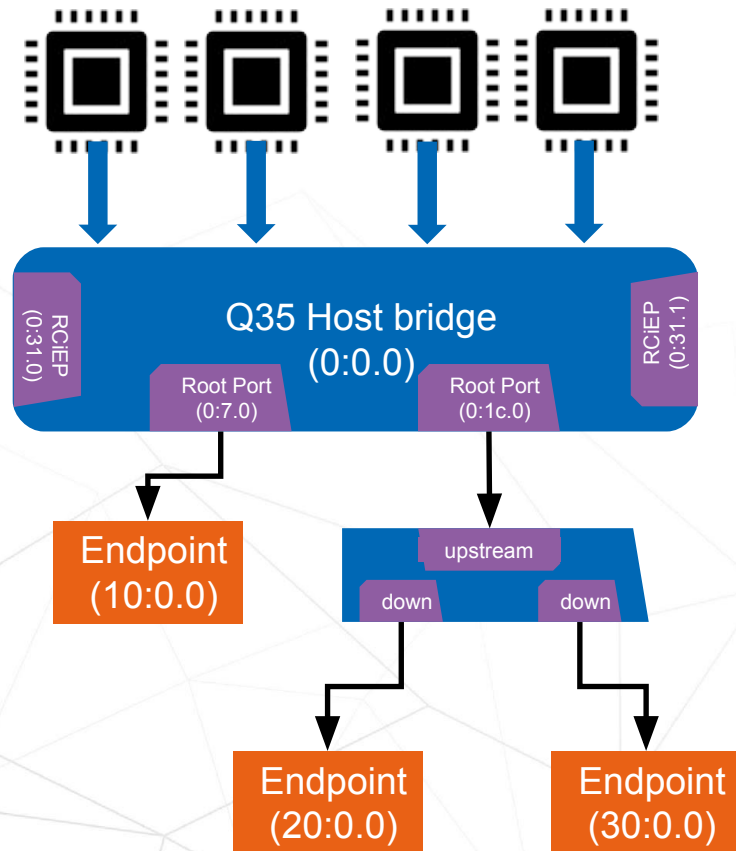




# PCIe in QEMU

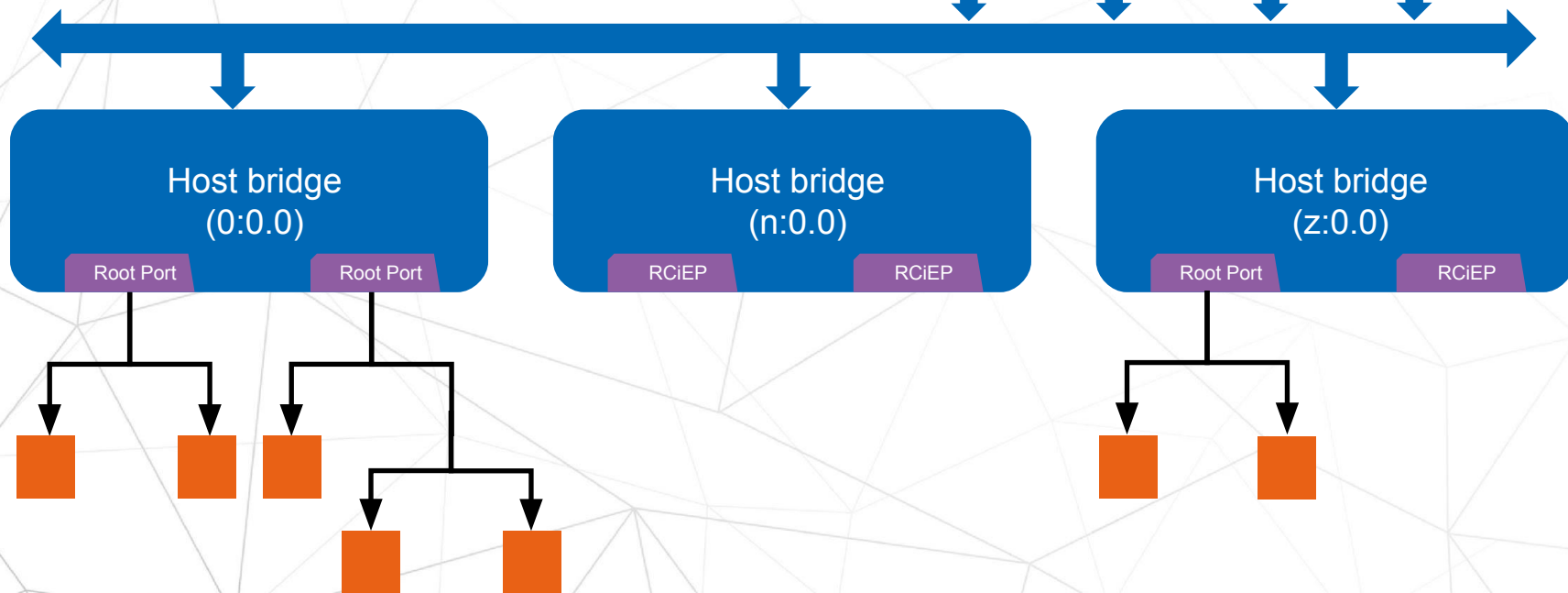
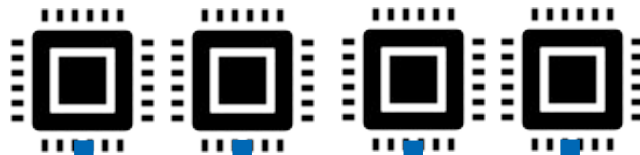
## What we all know and love

- Single root complex
  - Endpoints
  - Root ports
  - Switches
- All traffic is funneled to the single host bridge
  - QPI/UPI (not modeled)





# PCIe ~2014



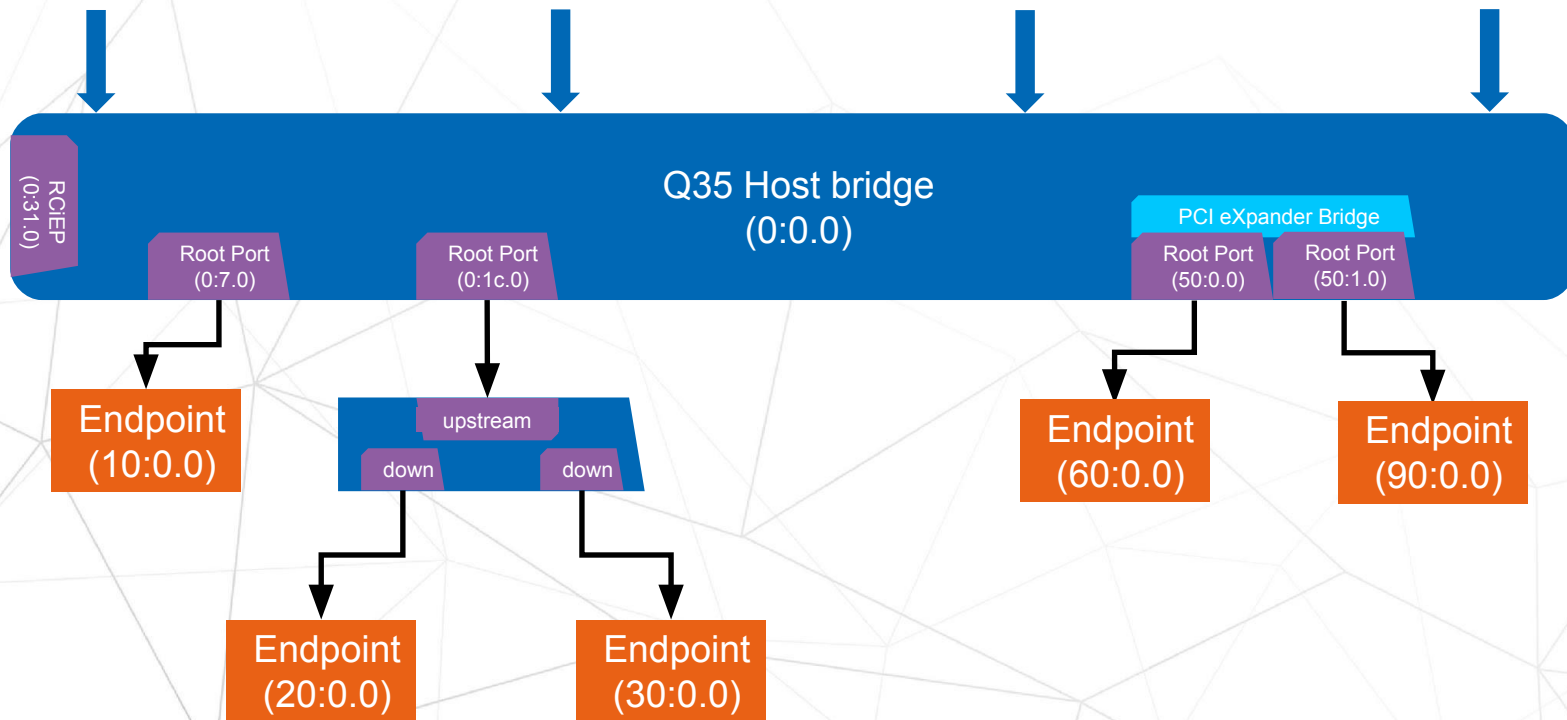


# Options

- Hacks to make Q35 - CXL 2.0
  - Limited potential for interleave scenarios
  - Touching Q35 is risky.
  - Mistakes make everything work incorrectly.
- Replace Q35 with something newer
  - Still Risky.
  - A lot of work for not much gain
  - What good does modeling UPI do for QEMU?
  - Doubtful community wants it (support burden).



# PCI eXpander Bridge (PXB)

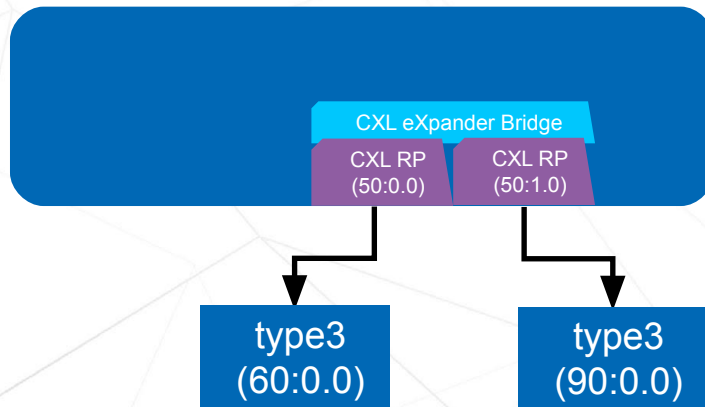




# LINUX September 20-24, 2021

## PLUMBERS CONFERENCE

- CXL Type 3 device
  - `hw/mem/cxl_type3.c`
- CXL Root Port
  - `hw/pci-bridge/cxl_root_port.c`
- CXL PXB
  - `hw/pci-bridge/pci_expander_bridge.c`





# LINUX PLUMBERS CONFERENCE

September 20-24, 2021

- NVDIMM & PCI had a baby...
- Inherits from both interfaces
- Mailbox handling

```
static const TypeInfo ct3d_info = {  
    .name = TYPE_CXL_TYPE3_DEV,  
    .parent = TYPE_PCI_DEVICE,  
    .class_init = ct3_class_init,  
    .instance_size = sizeof(CXLType3Dev),  
    .instance_init = ct3_instance_init,  
    .instance_finalize = ct3_finalize,  
    .interfaces = (InterfaceInfo[]) {  
        { TYPE_MEMORY_DEVICE },  
        { INTERFACE_CXL_DEVICE },  
        { INTERFACE_PCIE_DEVICE },  
        {}  
    },  
};
```



# LINUX PLUMBERS CONFERENCE

September 20-24, 2021

```
archlinux ~ # cxl list  
[  
  {  
    "memdev":"mem0",  
    "pmem_size":268435456,  
    "ram_size":0,  
    "fw_revision":"BFWF VERSION 00",  
    "partition_align":0,  
    "lsa_size":0  
  }  
]
```

ndctl/cxl

libcxl  
• IOCTL

Linux

cxl\_pci.ko  
• Mailbox  
MMIO

QEMU

cxl-mailbox-  
utils  
• Mailbox  
MMIO



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

# November 10th 2020

Spec  
Released

Linux  
Patches  
submitted

QEMU  
patches  
submitted

## PATCHBOMB ALL THE THINGS







# LINUX September 20-24, 2021 PLUMBERS CONFERENCE

- QEMU v3 patches sent
  - v4 is ready for submission
  - Community contributions for DOE, CDAT, and SPDM
  - High bar for adding more
    - Nothing exists quite like a CXL memory device
      - Volatile + Persistent capacities
      - Interleaving at multiple levels
- Linux phase 1 driver merged in 5.12
  - Phase 2 actively being developed.
  - Community contributions for DOE and CDAT
- Spec issues found and the fixes prototyped in QEMU



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

Future



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**

# CXL RUST FTW!!!!!!!!!!!!!!

**CXL**

Compute  
Express  
Link™





# What to do

- Community didn't adopt
  - Minimal feedback
  - Major reworks for interleave
  - cxl\_test came along
- External contributions
  - DOE
  - CDAT
  - SPDM
- Commercial Adoption...



# LINUX PLUMBERS CONFERENCE

September 20-24, 2021



JAKE-CLARK.TUMBLR



**LINUX** September 20-24, 2021  
**PLUMBERS**  
**CONFERENCE**

# We're Hiring...

<https://jobs.intel.com/ShowJob/Id/3089754/Linux-Kernel-Development-Engineer>



**LINUX** September 20-24, 2021  
**PLUMBERS  
CONFERENCE**





# LINUX PLUMBERS CONFERENCE

September 20-24, 2021

- Linux

- DPA mapping (WIP)
  - Libnvdimm integration
- Interleave
  - Provisioning
  - Recognition
- Hotplug
  - Hot add
  - Managed remove
- Asynchronous mailbox

- Userspace

- Testing

- QEMU

- Better tests
- Upstream/Downstream Ports
- Interleave Support
  - Host bridge
  - switch
- More firmware commands
- Hotplug support
- Error testing
- Interrupt support
- **Memory class device overhaul**
- -----
- Make Q35 CXL capable
- CXL type 1 and 2 devices
- CXL 1.1