

PCI Express - Architecture, Implementation and Measurement.

Jan Marjanovic (MTCA Tech Lab/DESY)
@janmarjanovic

2019-10-22

Design and Validation of Multi-Gigabit Systems
DESY, Hamburg

microTCA
TECHNOLOGY LAB

HELMHOLTZ
RESEARCH FOR GRAND CHALLENGES



- ▶ MicroTCA and PCI Express
- ▶ PCI Express protocol
- ▶ PCI Express protocol - physical layer
- ▶ PCI Express protocol - configuration
- ▶ FPGA implementation
- ▶ Examples
- ▶ Further resources and summary

This is a hardware event - we focus on the hardware side

Other resources available on system level:

▶ **MTCA training courses: Basic and Advanced**

https://techlab.desy.de/services/training/index_eng.html

▶ **tutorial talks:**

MTCAWS 2018:

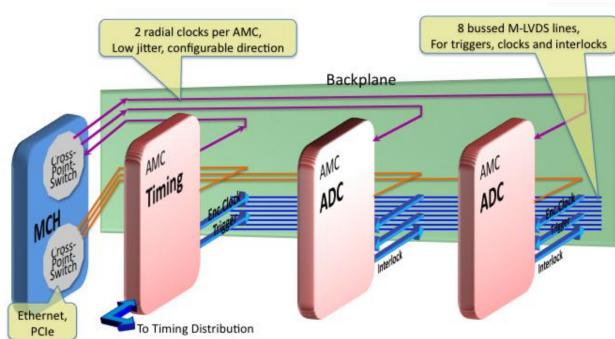
https://github.com/MicroTCA-Tech-Lab/damc-tck7-fpga-bsp/raw/demo-mtcaws2018/docs/demo/mtcaws2018_marjanovic_dma.pdf

MTCAWS China 2019:

<https://indico.ihep.ac.cn/event/9651/contribution/24/material/slides/0.pdf>

MicroTCA and PCI Express

A typical ADC system in MicroTCA is shown here:



**Figure 6-6: A 'typical' analog front-end.
CPU and further AMC's are not shown**

from PICMG® MicroTCA.4 Enhancements for Rear I/O and Timing R1.0

A typical board in such a system is shown here:

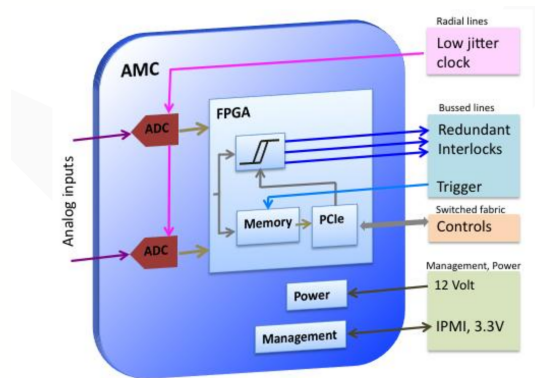
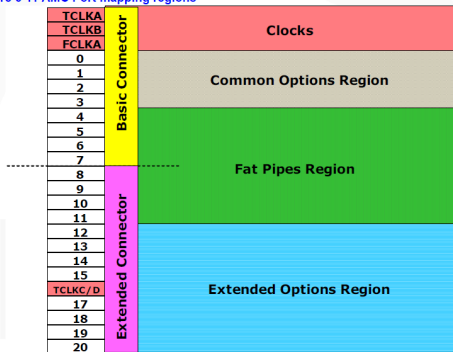


Figure 6-7: An example of a fast front-end (ADC AMC0)

from PICMG® MicroTCA.4 Enhancements for Rear I/O and Timing R1.0

AMC standard provides some guidance on how the ports should be assigned:

Figure 6-11 AMC Port mapping regions



from PICMG® Advanced Mezzanine Card AMC.O Specification R2.0

MicroTCA backplane connections as specified by MicroTCA.4 standard:

- REQ 6-6** Port 0 **should** be used for base Ethernet interface. Port 1 is connected to the optional (redundant) MCH.
- REQ 6-7** Port 4 to 7 **should** be used for PCIe.
- REQ 6-8** Port 12 to 15 **may** be used for application specific wire-ring.
- REQ 6-9** Ports 17- 20 **should** be wired as a bus according Section 6.4.
- REQ 6-10** FCLKA (Clk3) **should** be used for PCIe clock distribution as defined in MTCA.O R1.0

...

from PICMG® Specification MTCA.4 Revision 1.0: MicroTCA Enhancements for Rear I/O and Precision Timing

MicroTCA backplane connections as specified by MicroTCA.4 standard:

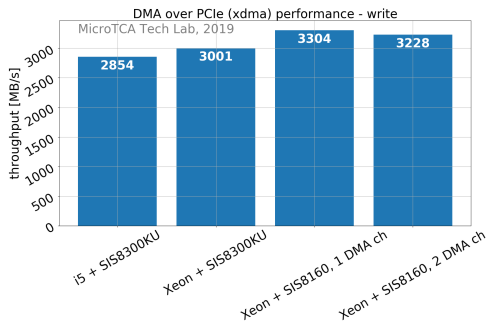
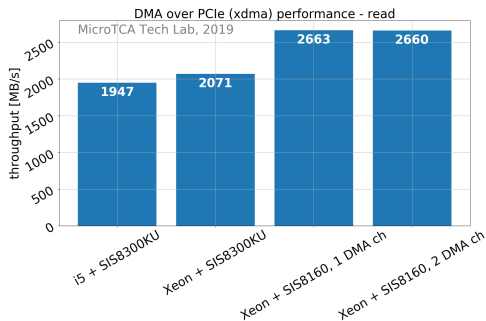
- REQ 6-6** Port 0 **should** be used for base Ethernet interface. Port 1 is connected to the optional (redundant) MCH.
- REQ 6-7** Port 4 to 7 **should** be used for PCIe.
- REQ 6-8** Port 12 to 15 **may** be used for application specific wire-ring.
- REQ 6-9** Ports 17- 20 **should** be wired as a bus according Section 6.4.
- REQ 6-10** FCLKA (Clk3) **should** be used for PCIe clock distribution as defined in MTCA.O R1.0

...

from PICMG® Specification MTCA.4 Revision 1.0: MicroTCA Enhancements for Rear I/O and Precision Timing

- ▶ High-throughput (up to 64 Gbps (gen 3, x8 link))
- ▶ Low-latency
- ▶ Software model (memory-based transaction, method to detect devices on the bus)

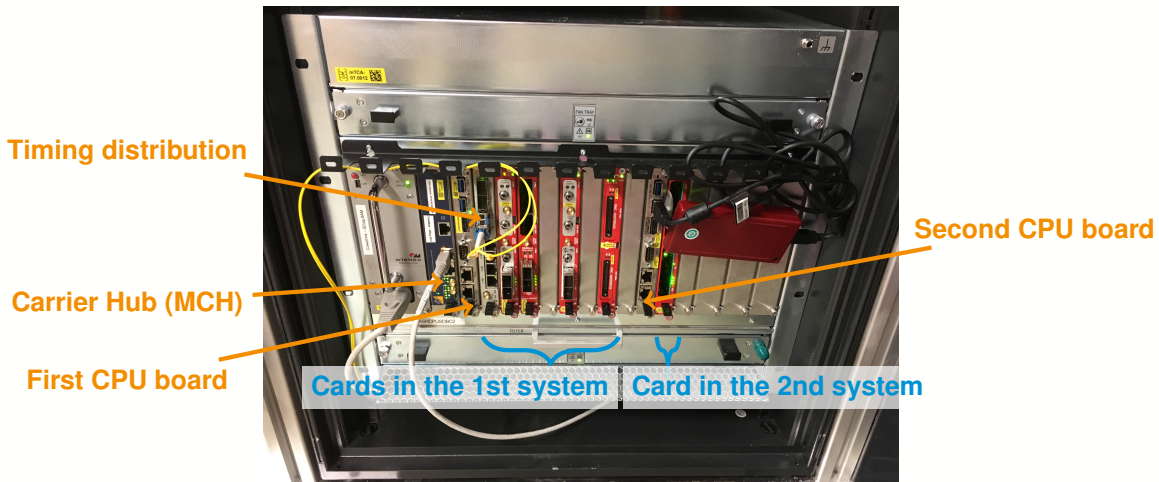
Example of the throughput achieved with two of the boards in MicroTCA:



A lot of boards available on the market use PCI Express as main/only communication protocol.
FPGAs also allow implementation of other protocols (e.g. 10 and 40 Gb Ethernet and Serial Rapid IO).

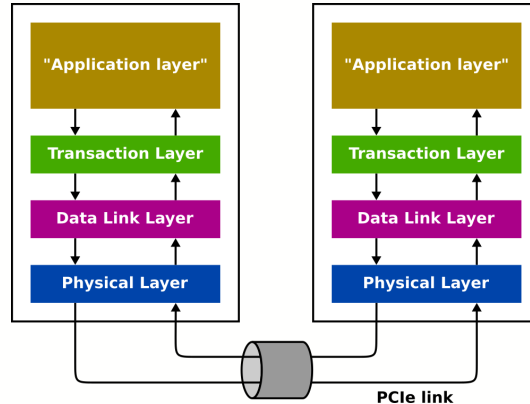


Example from FLASH: PCI Express switch enables splitting crate in two separate systems, but with shared timing and interlocks.



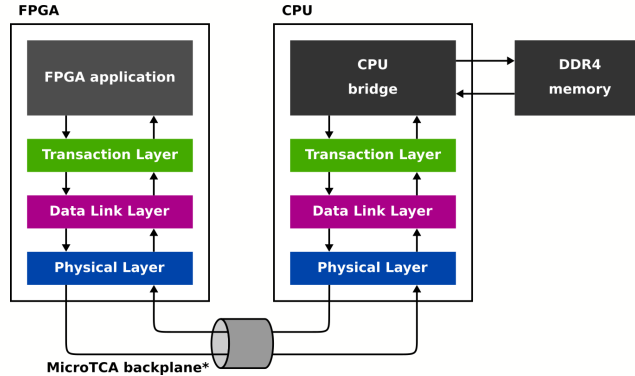
PCI Express protocol

PCIe express is high-speed protocol for PC extension cards. Different form-factors are also available. Link out of 1, 2, 4, 8 or 16 lanes, lane rate 2.5 GT/s, 5.0 GT/s, 8.0 GT/s and 16.0 GT/s.



inspired by Figure 2-12 from PCI Express Technology 3.0

Concrete example:



* real PCIe link in MicroTCA system includes a PCIe switch in MicroTCA Carrier Hub (MCH)

Physical Layer

- ▶ two differential pairs (TX and RX) per lane
- ▶ system clock (typical 100 MHz)
- ▶ 8b/10b (2.5 and 5.0 GT/s) or 128b/130b encoding (8.0 GT/s)
- ▶ Link Training and Status State Machine (LTSSM)
- ▶ in FPGAs implemented with transceivers + some logic

Data Link Layer

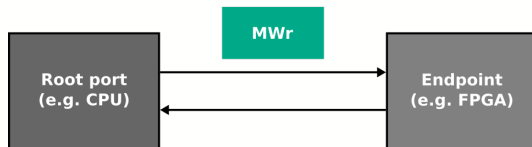
- ▶ Exchange of flow-control credits (= free buffers in receiver)
- ▶ CRC check, ACK/NAK, re-transmission

Transaction Layer

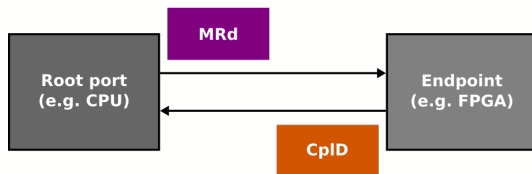
- ▶ "Real application" packets - memory read and write, ...

Packets can be one of several types: Memory Read, Memory Write, Completion, Completion with Data, Messages, Configuration Accesses and some legacy types.

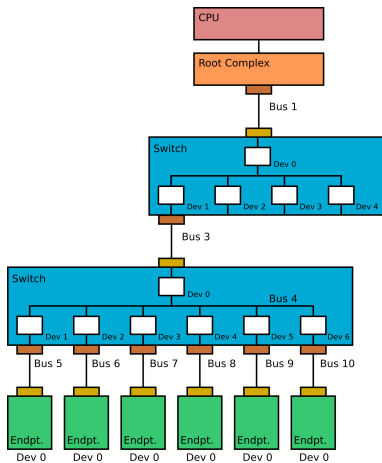
CPU writes to a device registers



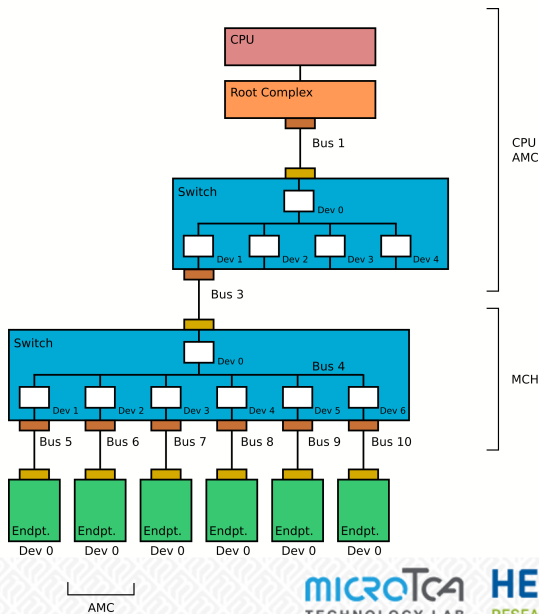
CPU reads from a device registers



Since PCIe is a point-to-point link, **switches** are used to connect more than 1 device to the CPU. Each PCIe device has its own **address** (bus:device.function), assigned during the **enumeration** process.



In the switch on the MCH one or more slots are configured as Upstream AMC.



Configuration Register Space defined for each Function. The first 256 bytes are the same as for PCI, a total of 4 kB for PCIe Extended Configuration Area.

The configuration space contains information such as: Vendor and Device ID, Link Status, Base Address Registers, ...

No driver ¹ is needed to read this space. Drivers in Linux use Vendor ID and Device ID (and Subsystem Vendor ID and Subsystem Device ID) to claim devices.

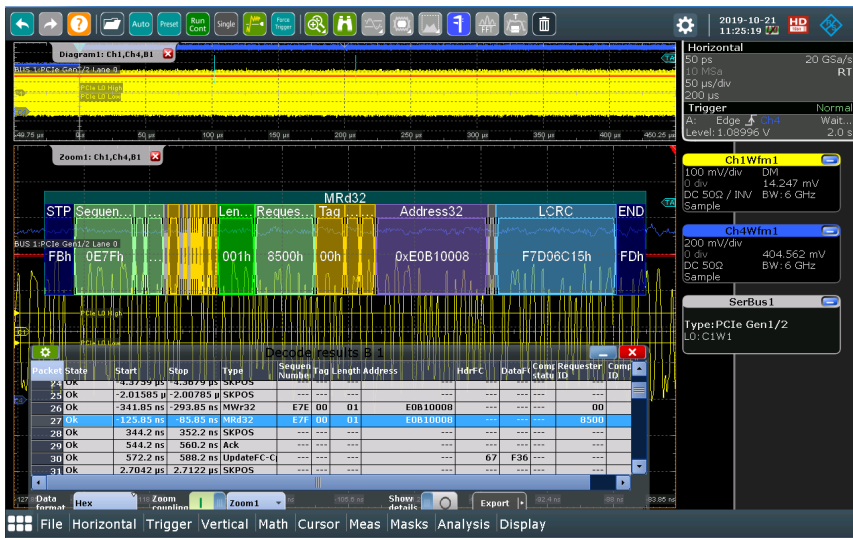
31		16 15		0		
Device ID		Vendor ID		00h		
Status		Command		04h		
Class Code			Revision ID			08h
BIST	Header Type	Lat. Timer	Cache Line S.			0Ch
Base Address Registers						10h 14h 18h 1Ch 20h 24h
Cardbus CIS Pointer						28h
Subsystem ID			Subsystem Vendor ID			2Ch
Expansion ROM Base Address						30h
Reserved				Cap. Pointer		34h
Reserved						38h
Max Lat.	Min Gnt.	Interrupt Pin	Interrupt Line			3Ch

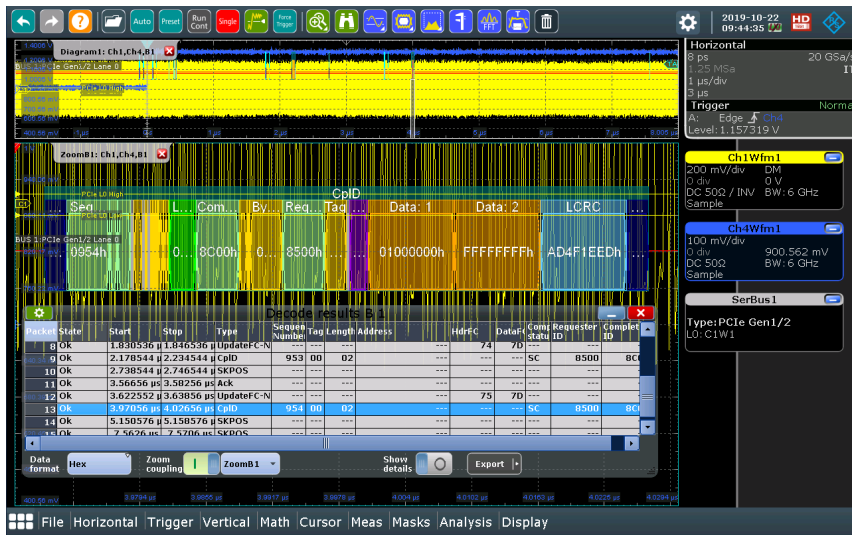
By Vijay Kumar Vijaykumar - Own work, Public Domain,
<https://commons.wikimedia.org/w/index.php?curid=3181779>

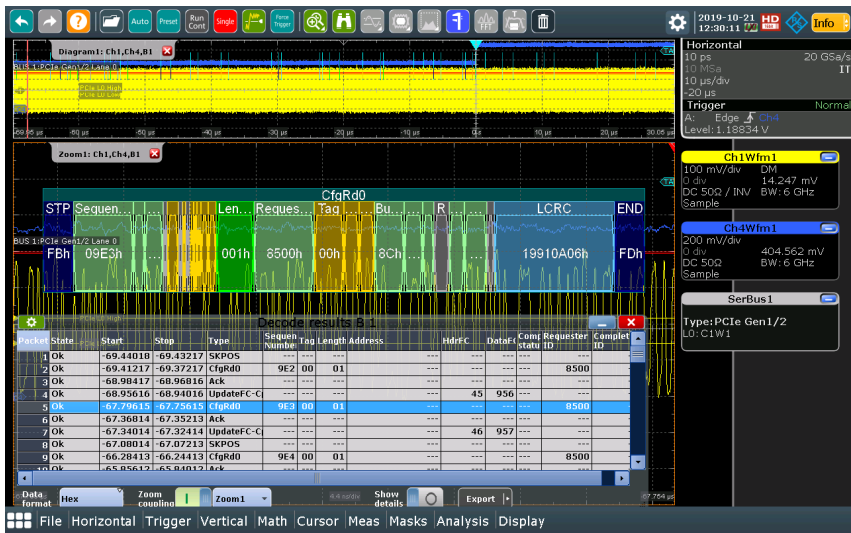
¹there is a driver embedded in OS, usually, i.e. in GNU/Linux

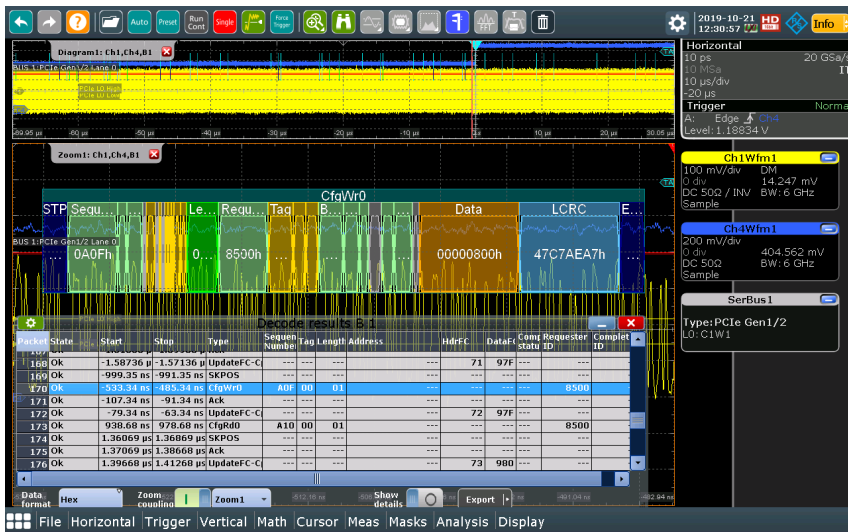
TLP examples



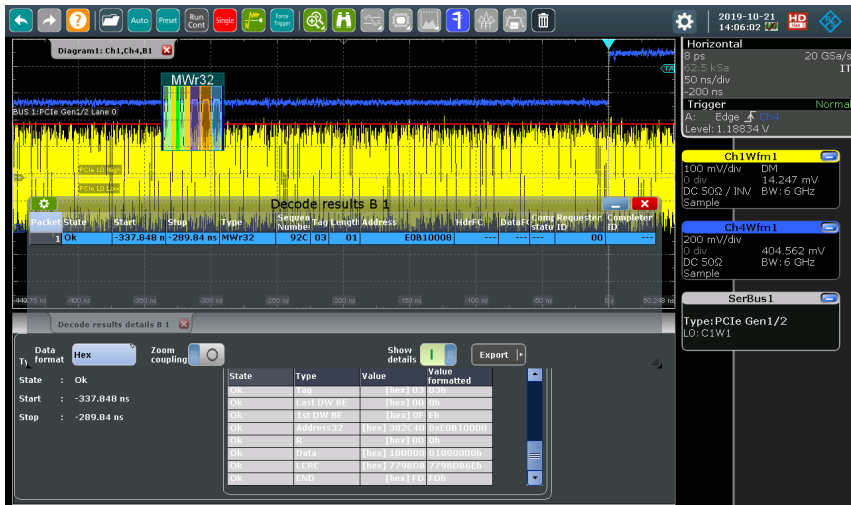








Turn on an LED with Memory Write



Software "walking" Configuration Area

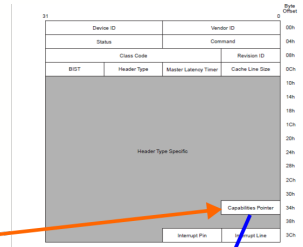
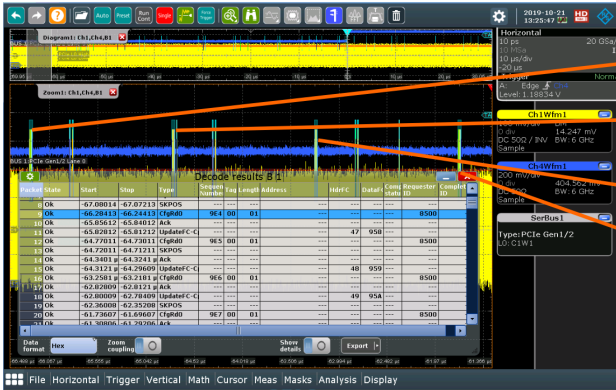


Figure 7-4: Common Configuration Space Header

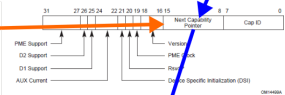
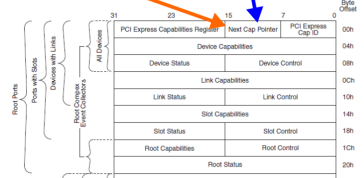


Figure 7-7: Power Management Capabilities Register



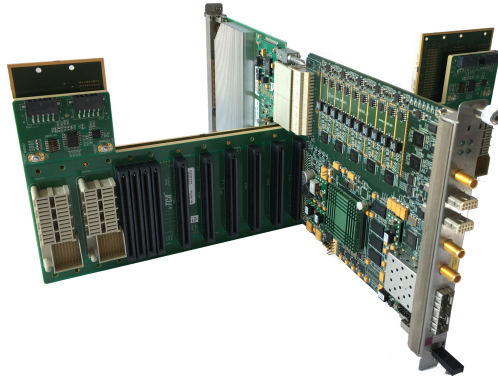
Figure 6-11: MSI Capability Structure



PCI Express protocol Physical Layer

In a typical MicroTCA systems a signal integrity is an important topic.

- ▶ Each high speed signal passes 2 card-edge connectors.
- ▶ Components from different vendors are plugged together



Link not running at the maximum possible rate:

PCIe Link Status Menu

	AMC1	AMC2	AMC3	AMC4	AMC5	AMC6	AMC7	AMC8	AMC9	AMC10	AMC11	AMC12	OPT1	RTM
	4..7	4..7	4..7	4..7	4..7	4..7	4..7	4..7	4..7	4..7	4..7	4..7		
	-	-	-	-	-	-	-	-	x4	x4	-	-	x8	-
Link Speed	-	-	-	-	-	-	-	-	5 GT/s	2.5 GT/s	-	-	8 GT/s	-

High rate of DLLP errors (rule of thumb: 1 correctable (re-trans) error per minute is OK):

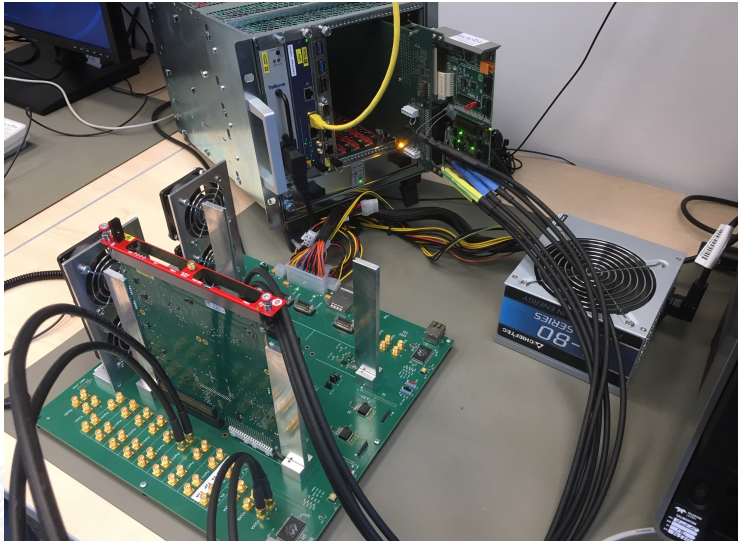
PCIe Error Counters Menu

	AMC1	AMC2	AMC3	AMC4	AMC5	AMC6
	4..7	4..7	4..7	4..7	4..7	4..7
	x4	x4	x4	x4	x4	x4
RCV_CNT	0	255	0	0	0	0
Bad TLP	0	0	0	0	0	0
Bad DLLP	0	65	0	0	0	0

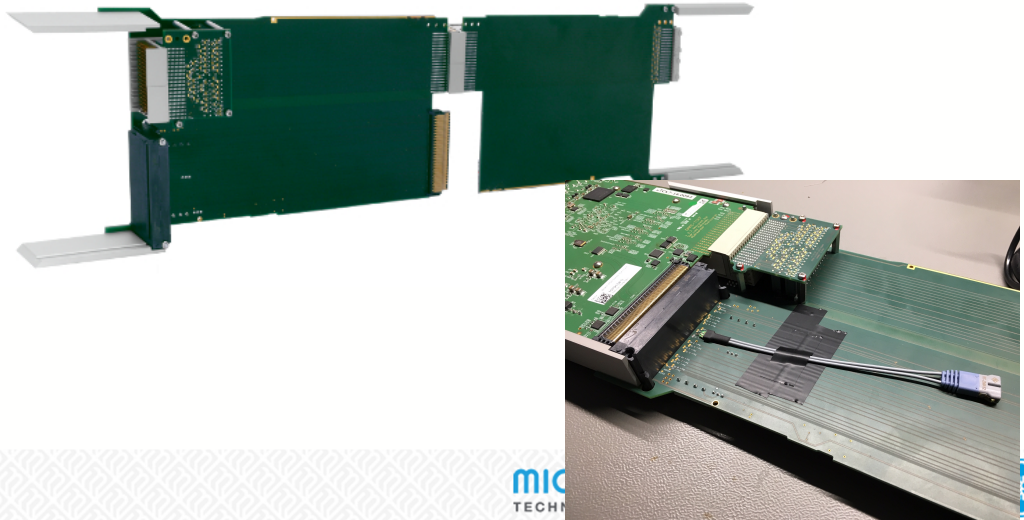
Reset Refresh

When developing or debugging boards in MicroTCA (or any other form factor) having test harnesses is extremely useful.

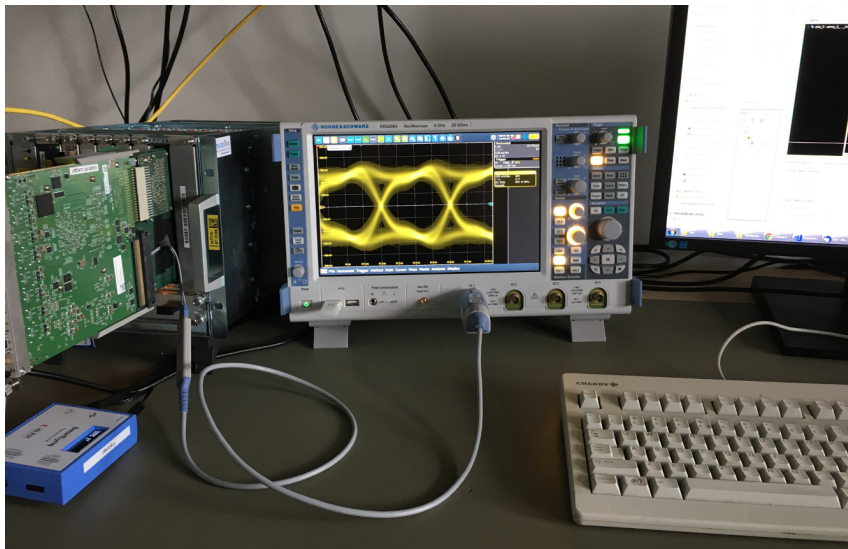
One example is an AMC test harness (PCIe connected with coax cables - 6 in total):



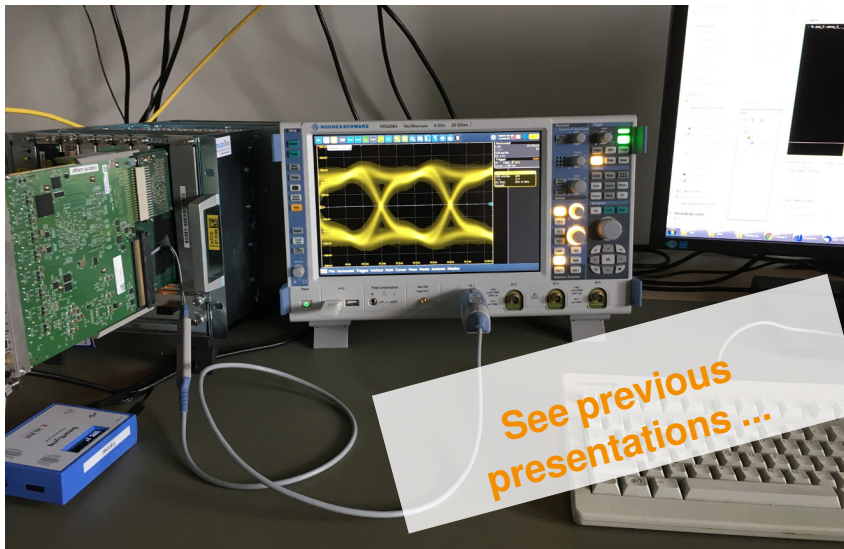
Extender boards (e.g. from N.A.T. - <https://www.nateurope.com/products/NAMC-EXT-RTM.html>) is a good equipment when we want to observe the data on the backplane:



With oscilloscope, PCIe 2.5 GT/s:



With oscilloscope, PCIe 2.5 GT/s:



In System IBERT allows measuring link quality **during operation** (everything in grey is not needed for normal operation of the link and is there only to be used for eye scan):

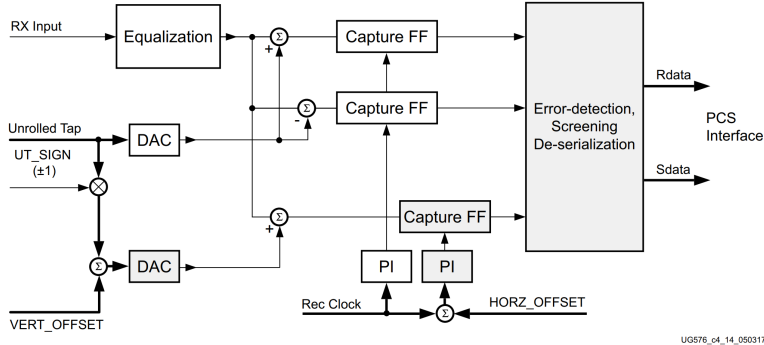
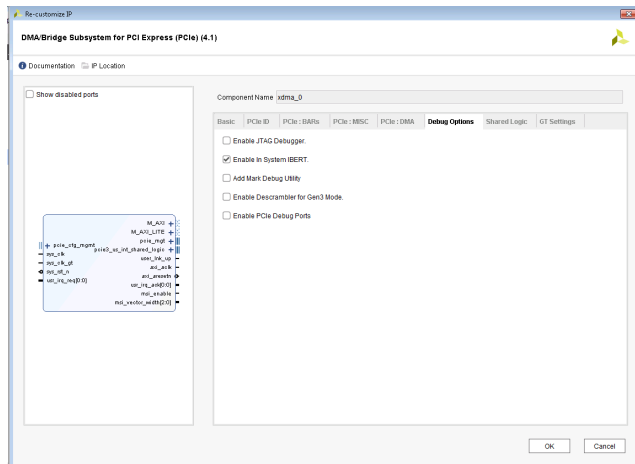


Figure 4-19: PMA Architecture to Support Eye Scan

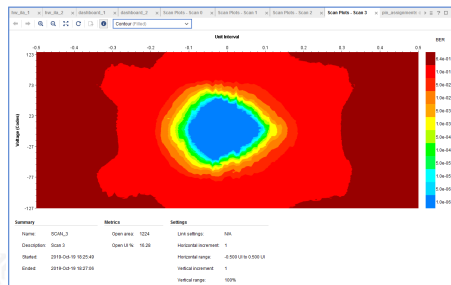
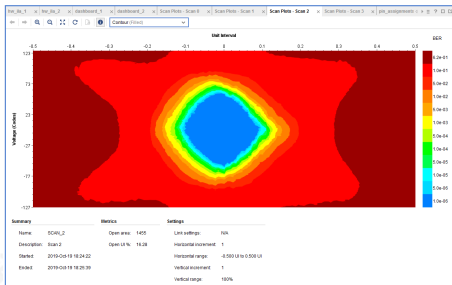
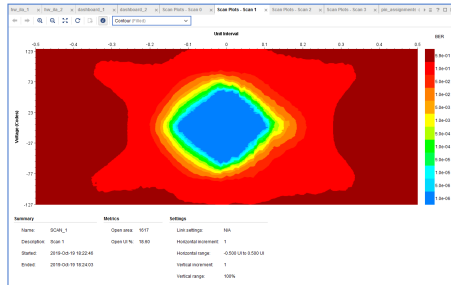
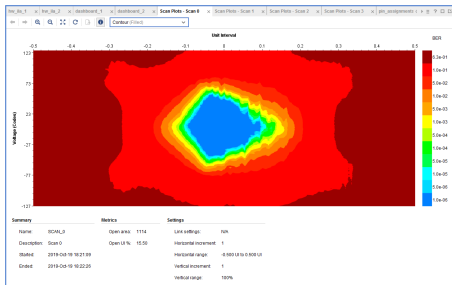
from

https://www.xilinx.com/support/documentation/user_guides/ug576-ultrascale-gth-transceivers.pdf

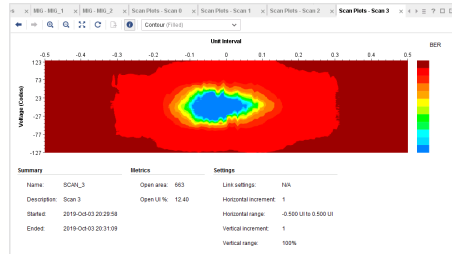
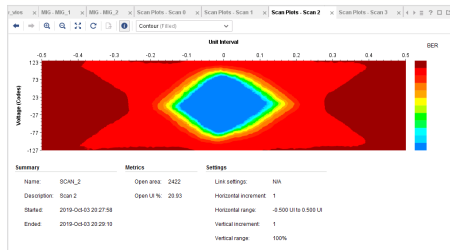
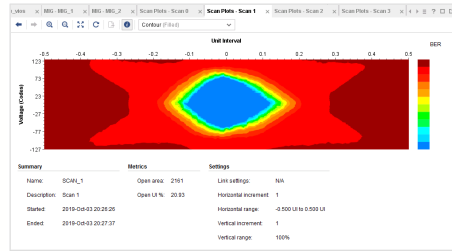
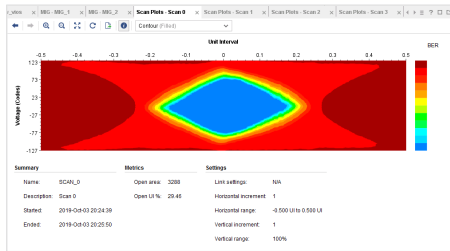
In System IBERT is enabled in IP configuration:



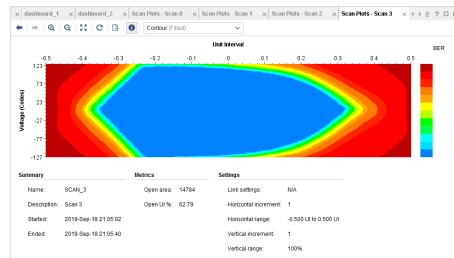
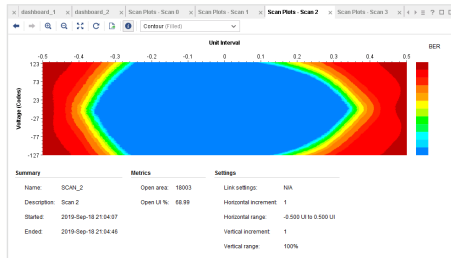
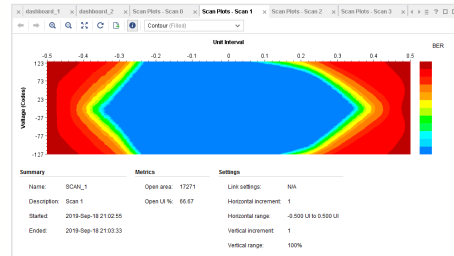
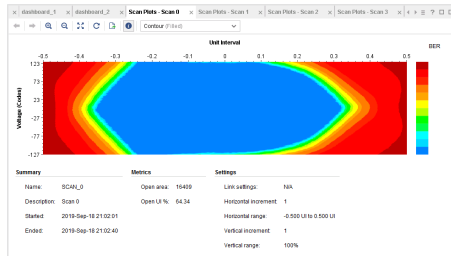
ELMA crate (12 slot) with SIS8300-KU at 8 GT/s



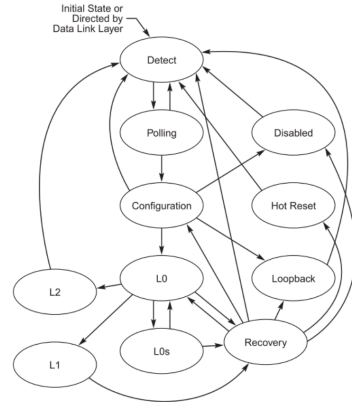
Schroff crate (12 slot) with SIS8160 at 8 GT/s



Schroff crate (7 slot) with SIS8300-KU at 5 GT/s (limited to 5 GT/s from slot configuration)



Link Training and State State Machine is a state machine running on both devices on the link. At the start up it negotiates link speed and link width, and removes the skew between lanes.



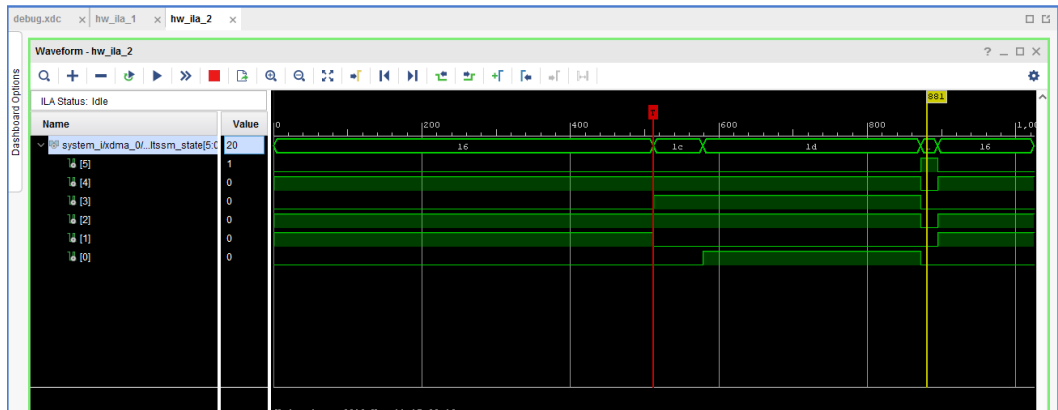
Very useful PCIe debug guide from Xilinx (recomends observing LTTSM state with Integrated Logic Analyzer - ILA):

https://www.xilinx.com/Attachment/Xilinx_Answer_56616_7_Series_PCIe_Link_Training_Debug_Guide.pdf

Here is an example of LTSSM state observed with ILA.

This example shows LTSSM leaving state L0 (the normal working state, encoded as 0x16) and cycling through states:

0x1C - Recovery Rcvrlock, 0x1D - Recovery Rcvrcfg, 0x20 - Recovery Idle.



PCI Express protocol Configuration

Two utilities from pciutils package (<https://github.com/pciutils/pciutils>) are useful to inspect and configure PCI Express system (on GNU/Linux):

- ▶ `lspci` - list PCI and PCIe devices, their organization and current status
- ▶ `setpci` - read and write individual registers in Configuration Space

+

Tools to access the interface from the other (back) side:

N.A.T. MCH (PCIe switch) interface

```

telnet MCH100191.tech.lab x
File Edit View Search Terminal Help
[11] : show switch port state
[12] : set switch port state
[13] : perform PRBS test
[14] : read FPGA register
[15] : write FPGA register
[16] : read Si5338 register
[17] : write Si5338 register
[18] : reset PCIe PCB
[?] : ?: help
[h] : h: help
[q] : q: quit submenu
PCIE (RET=0/0x0): 8
select hub module (0=MCH1, 1=MCH2) (RET=0/0x0):
Enter port (RET=2/0x2):
Enter address (RET=62/0x3e):
Enter access mode (0=TP, 1=NT-L or 2=NT-V) (RET=0/0x0):
HUB REG 0x0000003e = 0x0013010b
  
```

Xilinx CFG interface

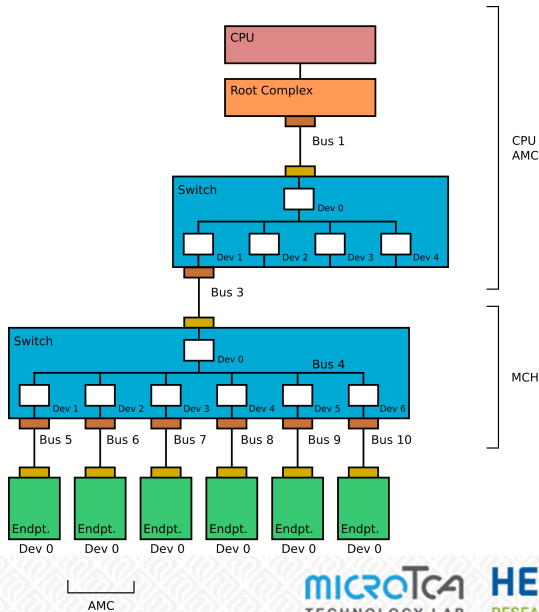
```

||- pcie_cfg_mgmt
  ||- ▶ cfg_mgmt_addr[18:0]
  ||- ▶ cfg_mgmt_byte_enable[3:0]
  ||- ◀ cfg_mgmt_read_data[31:0]
  ||- ▶ cfg_mgmt_read
  ||- ◀ cfg_mgmt_read_write_done
  ||- ▶ cfg_mgmt_type1_cfg_reg_access
  ||- ▶ cfg_mgmt_write_data[31:0]
  ||- ▶ cfg_mgmt_write
  ||- sys_clk
  ||- sys_rst_n
  ||- usr_irq_req[0:0]
  M_A\
  pc
  us
  ax
  usr_irq
  m
  msi_vector_
  
```

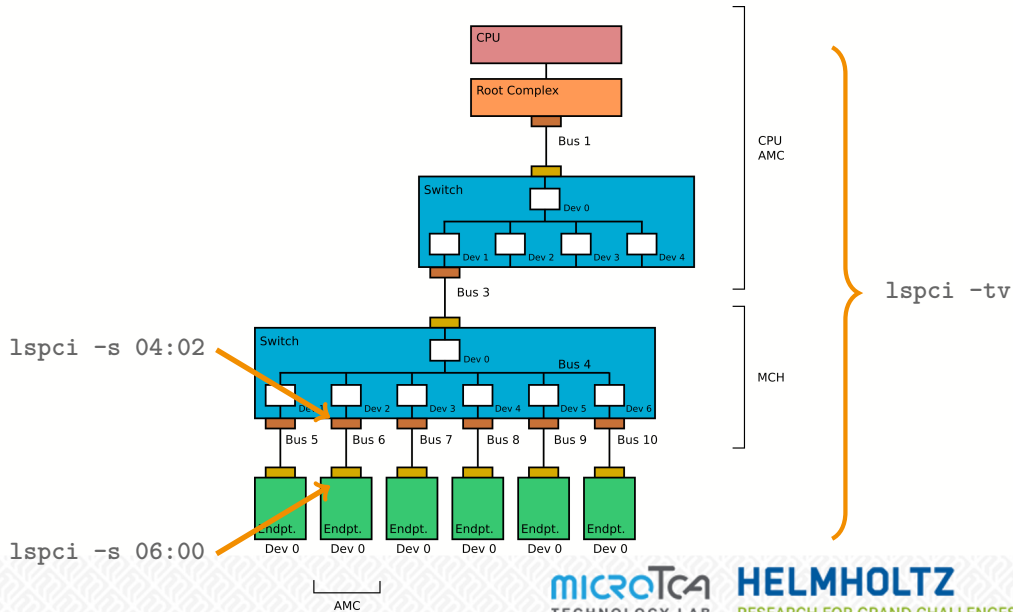
Example on Concurrent Tech AM G6x board:

```
$ lspci
00:00.0 Host bridge: Intel Corporation Xeon E3-1200 v6/7th Gen Core Processor Host Bridge /
00:01.0 PCI bridge: Intel Corporation Xeon E3-1200 v5/E3-1500 v5/6th Gen Core Processor
...
00:1f.3 Audio device: Intel Corporation CM238 HD Audio Controller (rev 31)
00:1f.4 SMBus: Intel Corporation 100 Series/C230 Series Chipset Family SMBus (rev 31)
01:00.0 PCI bridge: PLX Technology, Inc. Device 8725 (rev ca)
01:00.1 System peripheral: PLX Technology, Inc. Device 87d0 (rev ca)
...
02:01.0 PCI bridge: PLX Technology, Inc. Device 8725 (rev ca)
...
03:00.0 PCI bridge: PLX Technology, Inc. PEX 8748 48-Lane, 12-Port PCI Express Gen 3 (8 / )
04:01.0 PCI bridge: PLX Technology, Inc. PEX 8748 48-Lane, 12-Port PCI Express Gen 3 (8 / )
04:03.0 PCI bridge: PLX Technology, Inc. PEX 8748 48-Lane, 12-Port PCI Express Gen 3 (8 / )
...
05:00.0 Serial controller: Xilinx Corporation Device 7024
0d:00.0 Ethernet controller: Intel Corporation Ethernet Controller X710 for 10GbE SFP+ (
...
```

Several arguments can be passed to lspci to obtain different information from the system.



Several arguments can be passed to lspci to obtain different information from the system.




```
$ sudo lspci -s 04:03.0 -vv
04:03.0 PCI bridge: PLX Technology, Inc. PEX 8748 48-Lane, 12-Port PCI Express Gen 3 (8 / )
Control: I/O+ Mem+ BusMaster+ SpecCycle- MemWINV- VGASnoop- ParErr- Stepping- S -
[...]
Bus: primary=04, secondary=06, subordinate=06, sec-latency=0
[...]
Capabilities: [68] Express (v2) Downstream Port (Slot+), MSI 00
    DevCap: MaxPayload 512 bytes, PhantFunc 0
            ExtTag- RBE+
    DevCtl: Report errors: Correctable+ Non-Fatal+ Fatal+ Unsupported+
            RlxdOrd+ ExtTag- PhantFunc- AuxPwr- NoSnoop+
            MaxPayload 256 bytes, MaxReadReq 128 bytes
    DevSta: CorrErr+ UncorrErr- FatalErr- UnsuppReq+ AuxPwr- TransPend-
    LnkCap: Port #3, Speed 8GT/s, Width x4, ASPM L1, Exit Latency L0s <2us, <4
            ClockPM- Surprise+ LLActRep+ BwNot+ ASPMOptComp+
    LnkCtl: ASPM Disabled; Disabled- CommClk-
            ExtSynch- ClockPM- AutWidDis- BWInt- AutBWInt-
    LnkSta: Speed 5GT/s, Width x4, TrErr- Train- SlotClk- DLActive+ BWMgmt+
    SltCap: AttnBtn+ PwrCtrl+ MRL+ AttnInd- PwrInd- HotPlug+ Surprise-
            Slot #3, PowerLimit 25.000W; Interlock- NoCompl-
    SltCtl: Enable: AttnBtn+ PwrFlt- MRL- PresDet- CmdCplt+ HPirq+ LinkChg+
            Control: AttnInd Unknown, PwrInd Unknown, Power- Interlock-
    SltSta: Status: AttnBtn- PowerFlt- MRL- CmdCplt- PresDet+ Interlock-
            Changed: MRL+ PresDet- LinkState-
    [...]
[...]
Kernel driver in use: pcieport
Kernel modules: shpchp
```



```
$ sudo lspci -s 06:00 -vv
06:00.0 Serial controller: Xilinx Corporation Device 7024 (prog-if 01 [16450])
  Subsystem: Xilinx Corporation Device 0007
  Physical Slot: 3
  Control: I/O- Mem- BusMaster- SpecCycle- MemWINV- VGASnoop- ParErr- Stepping- SERR-
  Status: Cap+ 66MHz- UDF- FastB2B- ParErr- DEVSEL=fast >TAbort- <TAbort- <MAbort- >
  Interrupt: pin A routed to IRQ 16
  Region 0: Memory at 92800000 (32-bit, non-prefetchable) [size=8M]
  Region 1: Memory at 93000000 (32-bit, non-prefetchable) [size=64K]
  [...]
  Capabilities: [60] Express (v2) Endpoint, MSI 00
    DevCap: MaxPayload 512 bytes, PhantFunc 0, Latency LOs <64ns, L1 unlimi
      ExtTag+ AttnBtn- AttnInd- PwrInd- RBE+ FLReset-
    DevCtl: Report errors: Correctable- Non-Fatal- Fatal- Unsupported-
      RlxdOrd+ ExtTag+ PhantFunc- AuxPwr- NoSnoop+
      MaxPayload 256 bytes, MaxReadReq 512 bytes
    DevSta: CorrErr- UncorrErr- FatalErr- UnsuppReq- AuxPwr- TransPend-
    LnkCap: Port #0, Speed 5GT/s, Width x4, ASPM L0s, Exit Latency L0s unli
      ClockPM- Surprise- LLActRep- BwNot- ASPMOptComp-
    LnkCtl: ASPM Disabled; RCB 64 bytes Disabled- CommClk-
      ExtSynch- ClockPM- AutWidDis- BWInt- AutBWInt-
    LnkSta: Speed 5GT/s, Width x4, TrErr- Train- SlotClk+ DLActive- BWMgmt-
    DevCap2: Completion Timeout: Range B, TimeoutDis-, LTR-, OBFF Not Suppo
    DevCtl2: Completion Timeout: 50us to 50ms, TimeoutDis-, LTR-, OBFF Disa
    LnkCtl2: Target Link Speed: 5GT/s, EnterCompliance- SpeedDis-
      [...]
  Kernel driver in use: xdma
  Kernel modules: xdma
```



Using setpci to enter the compliance mode

Another system: PCIe up-link (Fujitsu server) and Struck SIS8300-KU AMC

```
$ lspci -tv
[...]
```

+-[0000:85]	-+-00.0-[86-8f]	---+00.0-[87-8f]	---+08.0-[88]	---	\-09.0-[89-8f]	---00.0-[8a-8f]	---+00.0-[8b]	---	+02.0-[8c]	----00.0	Xilinx 8021
									+04.0-[8d]	--	
									+06.0-[8e]	--	
									\-08.0-[8f]	--	
									+00.1	PLX Technology, Inc. Device 87d0	
									+00.2	PLX Technology, Inc. Device 87d0	
									+00.3	PLX Technology, Inc. Device 87d0	
									+00.4	PLX Technology, Inc. Device 87d0	
									\-00.4	PLX Technology, Inc. Device 87d0	

```
[...]
```

To enter the compliance mode, we must set **Enter Compliance** bit and then perform the **Hot Reset**.

PCI EXPRESS BASE SPECIFICATION, REV. 3.0

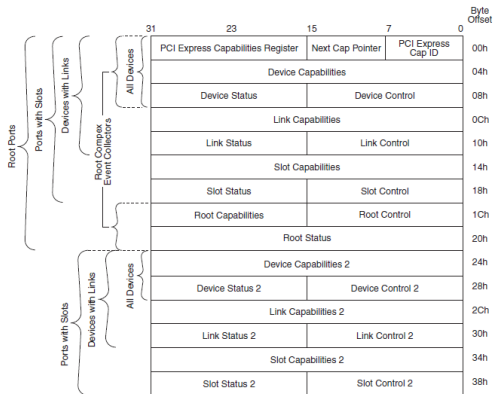
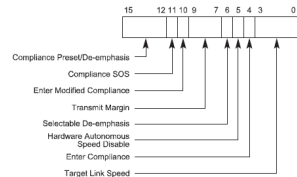


Figure 7-10: PCI Express Capability Structure

OM14318B

7.8.19. Link Control 2 Register (Offset 30h)



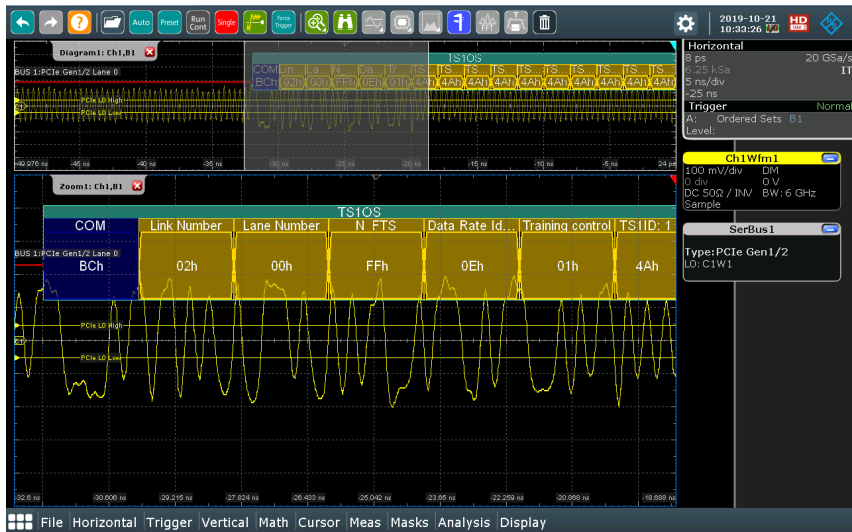
A-0573A

Figure 7-28: Link Control 2 Register

Commands to enter the compliance mode:

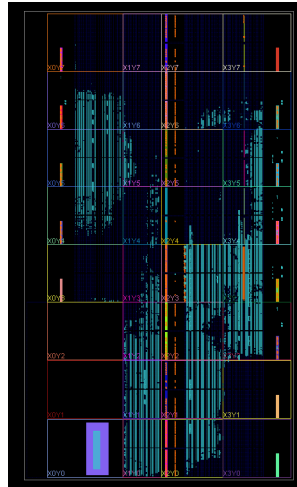
```
$ sudo setpci -s 8a:02 CAP_EXP+0x30.W  
0003  
$ sudo setpci -s 8c:00 CAP_EXP+0x30.W  
0002  
$ sudo setpci -s 8a:02 BRIDGE_CONTROL.W  
0013  
  
$ sudo setpci -s 8a:02 CAP_EXP+0x30.W=0x13  
$ sudo setpci -s 8c:00 CAP_EXP+0x30.W=0x0012  
$ sudo setpci -s 8a:02 BRIDGE_CONTROL.W=0x53
```

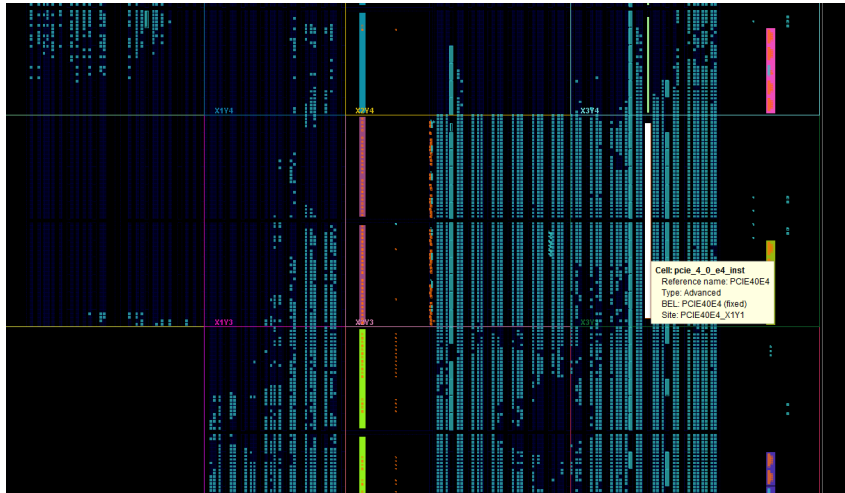
We are now in compliance mode (both link partners are transmitting TS1 Ordered Sets):

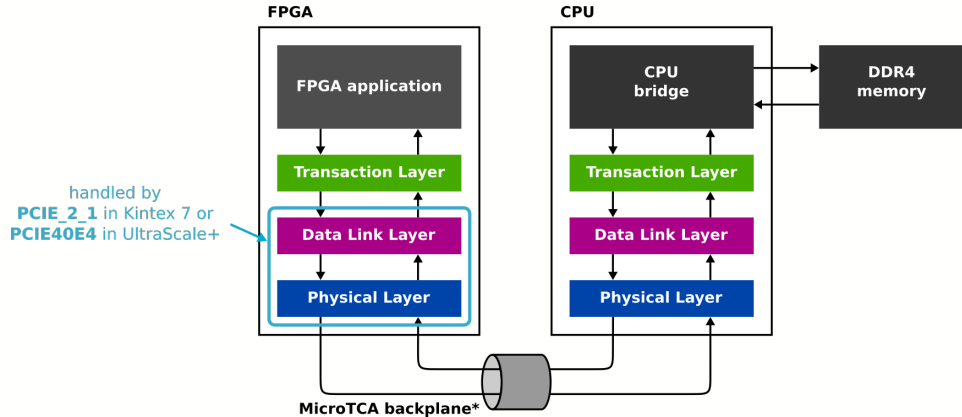


FPGA implementation

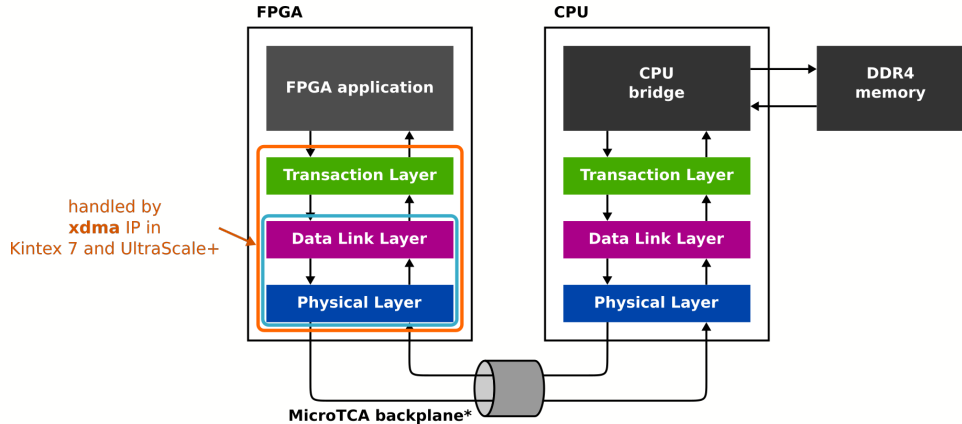
Modern FPGAs from Xilinx and Intel include hard IP cores to handle lower layers of protocol.



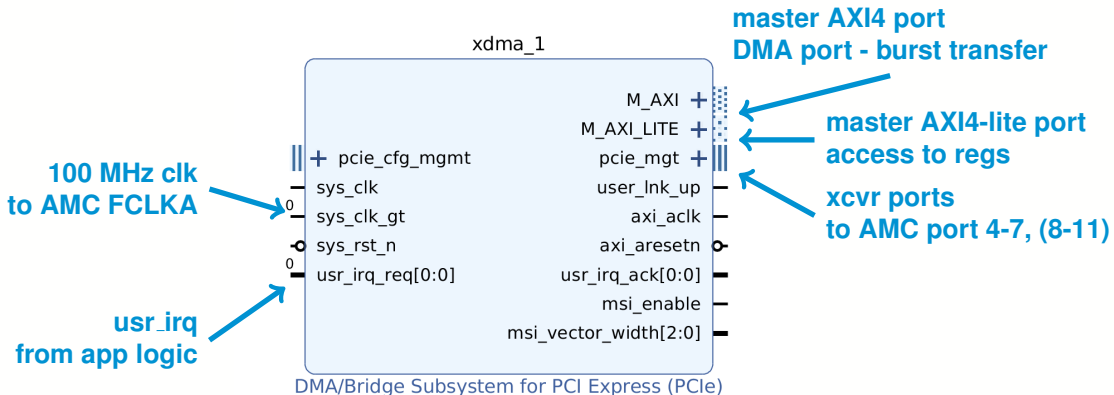




On top of that, both Xilinx and Intel provide IP cores to also handle the Transaction Layer protocol



This is how the Xilinx DMA Subsystem for PCI Express looks in Vivado:



DAMC-TCK7: Vivado project including support for PCIe, DDR3, clock configuration, IBERT on SFP+ and UDP/IPv4 beacon on AMC port 0; available under permissive license (3-clause BSD)

Available on our GitHub page: <https://github.com/MicroTCA-Tech-Lab>

DAMC-TCK7 Board Support Package

This is a Board Support Package in form of a Vivado project for DAMC-TCK7. It serves as a demonstration of board features and as a starting point for future developments.

This project is maintained by MicroTCA Technology Lab at DESY.

Features

PCIe Gen2 x4

PCIe Link Status Menu						
	AMT	AME	AMS	AMK	AME	AME
Link Speed	8.0	8.0	8.0	8.0	8.0	8.0
	8.0	8.0	8.0	8.0	8.0	8.0

DDR3 memory

On-board clock configuration

Project Overview

Component	Part	Package	Part	Package	Part	Package
U1	AMD_ZYNQ7020	XC7Z020-1CLG400	U2	AMD_ZYNQ7020	XC7Z020-1CLG400	

Timing Diagram

Console

```
ERROR: [PLACE-10] Cannot place component 'U1' on board 'xc7z020-1clg400'. Reason: The device 'U1' has a fixed location in the board. The location is 'xc7z020-1clg400'. The component 'U1' is already placed at this location.
```

Block Diagram

Setting up software development environment

1. Set the project dir.
2. Source the Vivado TCL script.
3. Source the Vivado TCL script.
4. Source the Vivado TCL script.

Compiling the project

1. Source the Vivado TCL script.
2. Run a compile, this is going to take some time.

Setting up software development environment

1. Create block diagram and IP wrapper (see "Block diagram").
2. From Vivado open SDC file in "Launch SDC".
3. For **Exported location** select "ExportedLocation". For **Workspace** select "Workspace".
4. Run the tool and save.
5. Import project file in Vivado -> Existing Projects Into Workspace.
6. Select root directory in Vivado.
7. Make sure that both "Add board connector" and "Add board connector" are selected.
8. Press Finish.

Committing changes

Block diagram

From project dir:

```
write_bd_tcl -xc7z020-1clg400 -tcl {xc7z020-1clg400.tcl}
```

Update MicroBlaze bit script

1. Open generator of the "MicroBlaze".
2. Add all to Vivado project.
3. Set parameters "SOURCE_PATH" in MicroBlaze with SOURCE_PATH in system.

Project settings

If you need to update "MicroBlaze" bit:

```
write_project_bit -xc7z020-1clg400 -tcl {xc7z020-1clg400.tcl}
```

Please be careful as there are some hard-coded modification in "MicroBlaze".

DAMC-TCK7: Vivado project including support for PCIe, DDR3, clock configuration, IBERT on SFP+ and UDP/IPv4 beacon on AMC port 0; available under permissive license (3-clause BSD)

Available on our GitHub page: <https://github.com/MicroTCA-Tech-Lab>

The image shows a collage of screenshots from the DAMC-TCK7 Board Support Package documentation and the Vivado IDE. On the left, the documentation includes a photo of the board, a list of features, and a PCIe Link Status Monitor table. The table shows link speeds for lanes 0-6, with lane 0 at 5.0 Gb/s and lanes 1-6 at 2.5 Gb/s. In the center, there is a Vivado console window showing the 'on-board check' results, which confirm the board configuration. On the right, there is a Vivado block design diagram and a 'Scripts' window listing project setup steps. A large orange text overlay reads: 'Stay tuned for DAMC-FMC2ZUP (Zynq UltraScale+ MPSoC)'.

PCIe Link Status Monitor						
	LANE0	LANE1	LANE2	LANE3	LANE4	LANE5
Link Speed	5.0 Gb/s	2.5 Gb/s	2.5 Gb/s	2.5 Gb/s	2.5 Gb/s	2.5 Gb/s

Stay tuned for
DAMC-FMC2ZUP
(Zynq UltraScale+
MPSoC)

- ▶ Xilinx provides a Linux driver for their DMA subsystem for PCIe
- ▶ Available on GitHub:

https://github.com/Xilinx/dma_ip_drivers/tree/master/XDMA/

- ▶ Previously available in Xilinx Answer Record #65444:

<https://www.xilinx.com/support/answers/65444.html>

- ▶ MicroTCA Tech Lab maintains an internal fork with some improvements → discussion with our legal department on licensing ongoing

On the other side (towards the user space), Xilinx DMA driver provides several char devices:

```
$ ll /dev | grep xdma
drwxr-xr-x  3 root root          60 Jun 20 17:19 xdma
crw-rw-rw-  1 root root    238, 36 Jun 20 17:19 xdma0_c2h_0
crw-rw-rw-  1 root root    238,  1 Jun 20 17:19 xdma0_control
crw-rw-rw-  1 root root    238, 10 Jun 20 17:19 xdma0_events_0
crw-rw-rw-  1 root root    238, 11 Jun 20 17:19 xdma0_events_1
[...]
crw-rw-rw-  1 root root    238, 24 Jun 20 17:19 xdma0_events_14
crw-rw-rw-  1 root root    238, 25 Jun 20 17:19 xdma0_events_15
crw-rw-rw-  1 root root    238, 32 Jun 20 17:19 xdma0_h2c_0
crw-rw-rw-  1 root root    238,  0 Jun 20 17:19 xdma0_user
crw-rw-rw-  1 root root    238,  2 Jun 20 17:19 xdma0_xvc
```

On the other side (towards the user space), Xilinx DMA driver provides several char devices:

```
$ ll /dev | grep xdma
drwxr-xr-x  3 root root          60 Jun 20 17:19 xdma
crw-rw-rw-  1 root root    238, 36 Jun 20 17:19 xdma0_c2h_0
crw-rw-rw-  1 root root    238,  1 Jun 20 17:19 xdma0_control
crw-rw-rw-  1 root root    238, 10 Jun 20 17:19 xdma0_events_0
crw-rw-rw-  1 root root    238, 11 Jun 20 17:19 xdma0_events_1
[...]
crw-rw-rw-  1 root root    238, 24 Jun 20 17:19 xdma0_events_14
crw-rw-rw-  1 root root    238, 25 Jun 20 17:19 xdma0_events_15
crw-rw-rw-  1 root root    238, 32 Jun 20 17:19 xdma0_h2c_0
crw-rw-rw-  1 root root    238,  0 Jun 20 17:19 xdma0_user
crw-rw-rw-  1 root root    238,  2 Jun 20 17:19 xdma0_xvc
```

c2h

AXI4 port (DMA from device)

h2c
AXI4 port (DMA to device)

user
AXI4-Lite port

events_N
interrupts

Examples

Example 1: CPU freezes on memcpy()

Using GPIO to trigger on software events

```

0x7ffff717f827 <__memcpy_sse2_unaligned+9>: lea    (%rdi,%rdi,1),%rcx
(gdb) x/5! $pc
=> 0x7ffff717f820 <__memcpy_sse2_unaligned>:  mov    %rsi,%rax
0x7ffff717f823 <__memcpy_sse2_unaligned+3>:  lea   (%rdi,%rdi,1),%rcx
0x7ffff717f827 <__memcpy_sse2_unaligned+7>:  sub   %rdi,%rax
0x7ffff717f82a <__memcpy_sse2_unaligned+10>: sub   %rdx,%rax
0x7ffff717f82d <__memcpy_sse2_unaligned+13>: cmp   %rcx,%rax
(gdb) x/10! $pc
=> 0x7ffff717f820 <__memcpy_sse2_unaligned>:  mov    %rsi,%rax
0x7ffff717f823 <__memcpy_sse2_unaligned+3>:  lea   (%rdi,%rdi,1),%rcx
0x7ffff717f827 <__memcpy_sse2_unaligned+7>:  sub   %rdi,%rax
0x7ffff717f82a <__memcpy_sse2_unaligned+10>: sub   %rdx,%rax
0x7ffff717f82d <__memcpy_sse2_unaligned+13>: cmp   %rcx,%rax
0x7ffff717f830 <__memcpy_sse2_unaligned+16>:  jb    0x7ffff717f93d <__memcpy_sse2_unaligned+285>
0x7ffff717f836 <__memcpy_sse2_unaligned+22>:  cnp   %s0,%rdx
0x7ffff717f83a <__memcpy_sse2_unaligned+26>:  jbe   0x7ffff717f9cb <__memcpy_sse2_unaligned+427>
0x7ffff717f840 <__memcpy_sse2_unaligned+32>:  movdqu(%rsi),%xmm8
0x7ffff717f845 <__memcpy_sse2_unaligned+37>:  cnp   %s0x20,%rdx

```

Example 2: CPU crashes on DMA (after some time)

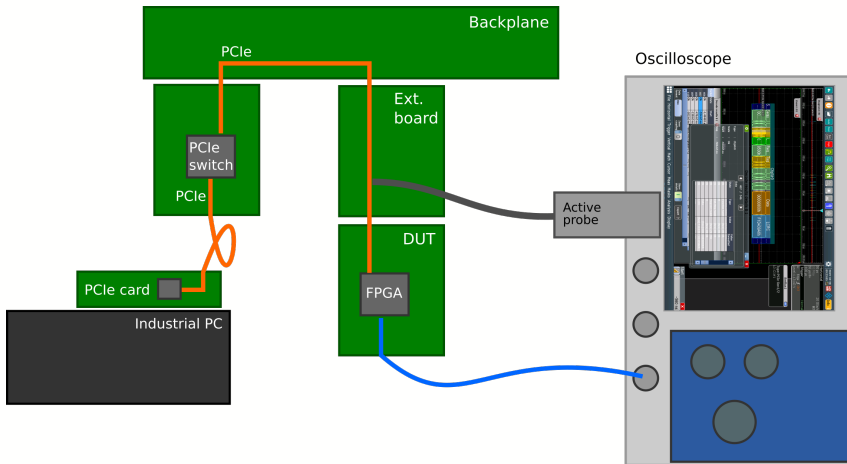
Using FPGA to help with the trigger

```

1923.983788 pcleport 0000:b9:00.0: AER: [0] RXERR
1923.983801 pcleport 0000:b9:00.0: AER: arc_layer=Physical Layer, aer_agent=Accelerator
1923.983565 (3)[Hardware Error]: Hardware error from APEI Generic Hardware Error Source: 0
1923.983568 (3)[Hardware Error]: It has been corrected by h/w and requires no further action
1923.983570 (3)[Hardware Error]: event severity: corrected
1923.983572 (3)[Hardware Error]: Error 0, type: corrected
1923.983574 (3)[Hardware Error]: section_type: PCIe error
1923.983576 (3)[Hardware Error]: port_type: 9, upstream_switch_port
1923.983577 (3)[Hardware Error]: version: 3.0
1923.983579 (3)[Hardware Error]: command: 0x0147, status: 0x0010
1923.983581 (3)[Hardware Error]: device_id: 0000:b9:00.0
1923.983582 (3)[Hardware Error]: slot: 0
1923.983584 (3)[Hardware Error]: secondary_bus: 0x8a
1923.983585 (3)[Hardware Error]: vendor_id: 0x10b5, device_id: 0x8780
1923.983587 (3)[Hardware Error]: class_code: 000406
1923.983588 (3)[Hardware Error]: bridge: secondary_status: 0x0000, control: 0x0013
1923.983823 pcleport 0000:b9:00.0: AER: aer_status: 0x0000a001, aer_mask: 0x00007000
1923.992489 pcleport 0000:b9:00.0: AER: [0] RXERR
1923.992702 pcleport 0000:b9:00.0: AER: [15] Receiver
1924.007000 pcleport 0000:b9:00.0: AER: arc_layer=Physical Layer, aer_agent=Accelerator ID
1985.423801 (4)[Hardware Error]: Hardware error from APEI Generic Hardware Error Source: 0
1985.453804 (4)[Hardware Error]: It has been corrected by h/w and requires no further action

```

In both examples a similar setup is used: active probe monitors PCIe differential pair, and another signal is connected to LVDS or LVCMOS output for trigger.



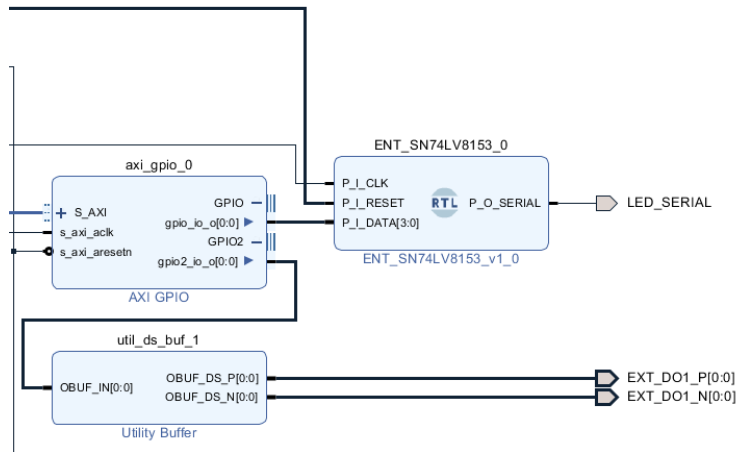
Introduction:

- ▶ To copy large amounts of data one usually uses DMA
- ▶ But sometimes in our applications we have to copy just a couple of registers (e.g. 32 bytes in total) → setting up DMA makes no sense
- ▶ In this case it would be convenient (from SW point of view) to use `memcpy()`

Issue:

- ▶ Performing `memcpy()` on more than 4 bytes freezes the CPU, reboot is required to recover (luckily we have IPMI)
- ▶ Reading registers one-by-one (4 bytes at the time) works
- ▶ `memcpy()` also works on Kintex UltraScale, but does not work on Kintex 7-series (with the same IP - xdma)

- ▶ ...
- ▶ The most convenient way to trigger on such event is to **set GPIO pin high** before starting the relevant SW test case

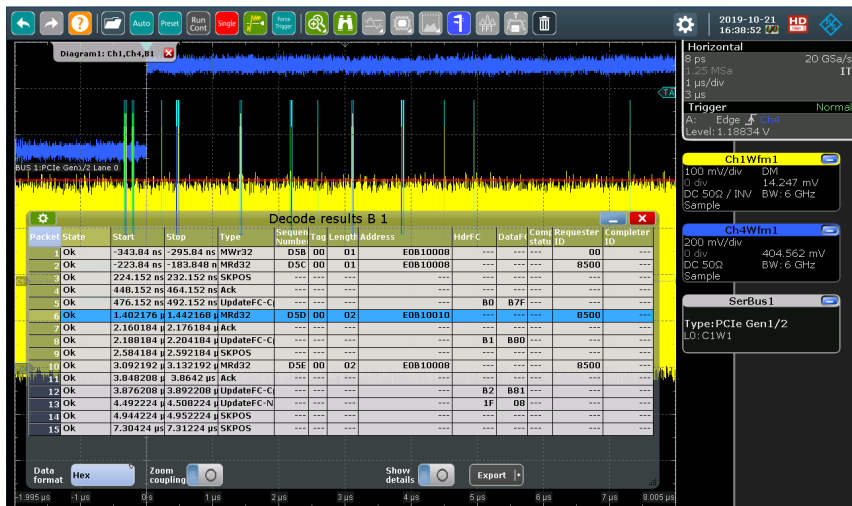


We first set GPIO pin high to trigger the oscilloscope, then perform a 32-bit read (for check) and then perform `memcpy()`.

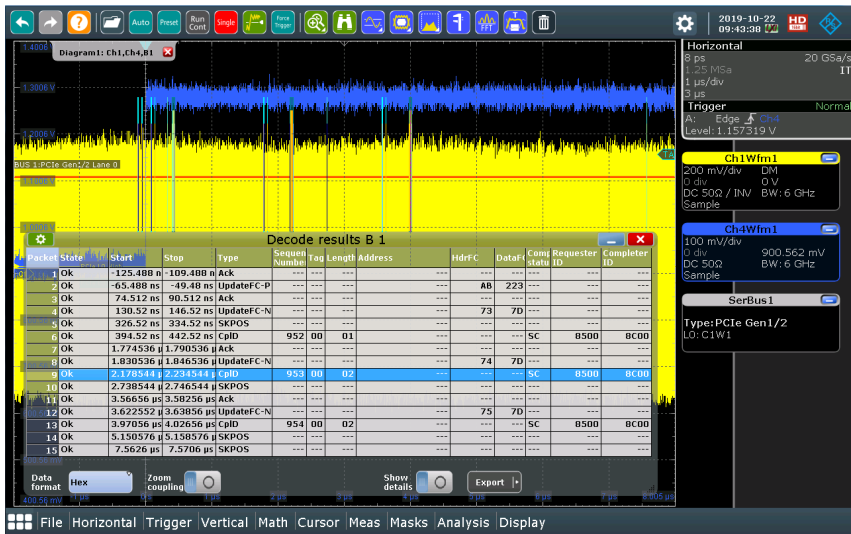
```
// virt_addr from mmap(), points to a memory in FPGA
gpio_control(GPIO0, 1);
uint32_t tmp = *((uint32_t*)virt_addr);
char buf[16];
memcpy(buf, (void*)virt_addr, sizeof(buf));
```


Example 1: measurement

TX side on SIS8300-KU (with Kintex UltraScale)

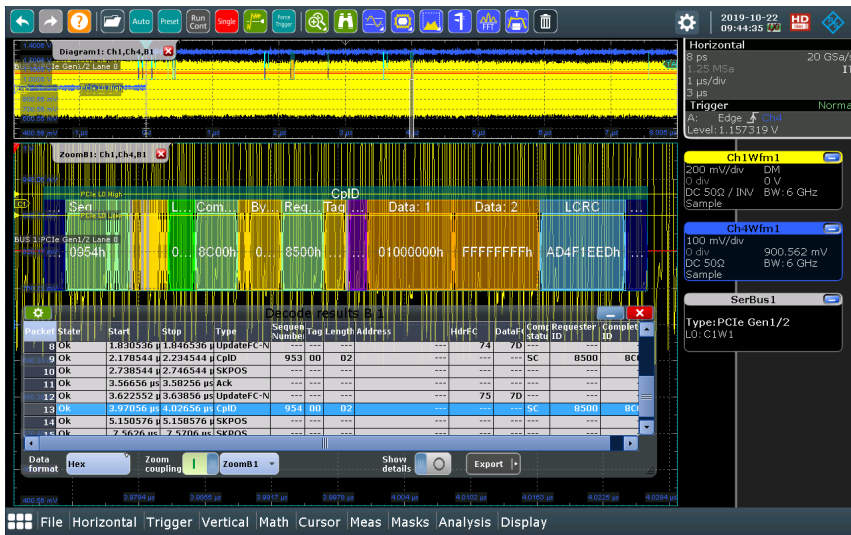


RX side on SIS8300-KU (with Kintex UltraScale)



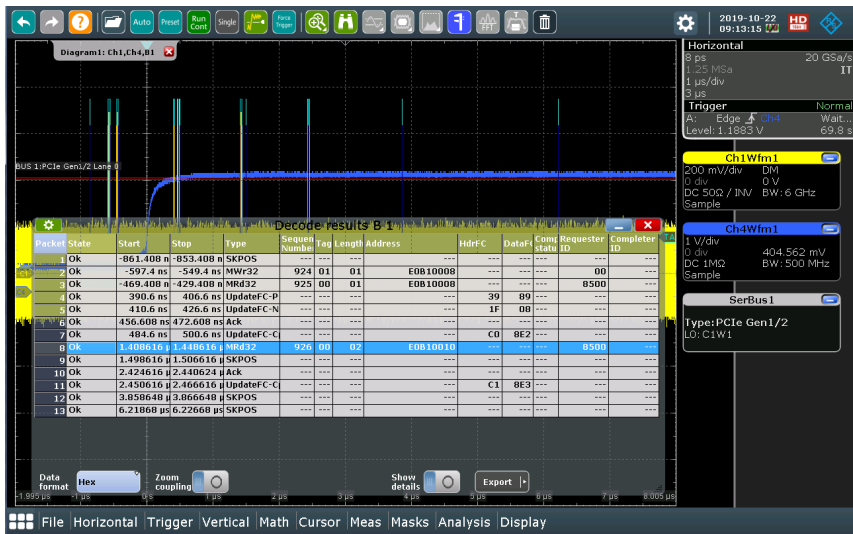
Example 1: measurement

RX side on SIS8300-KU (with Kintex UltraScale)



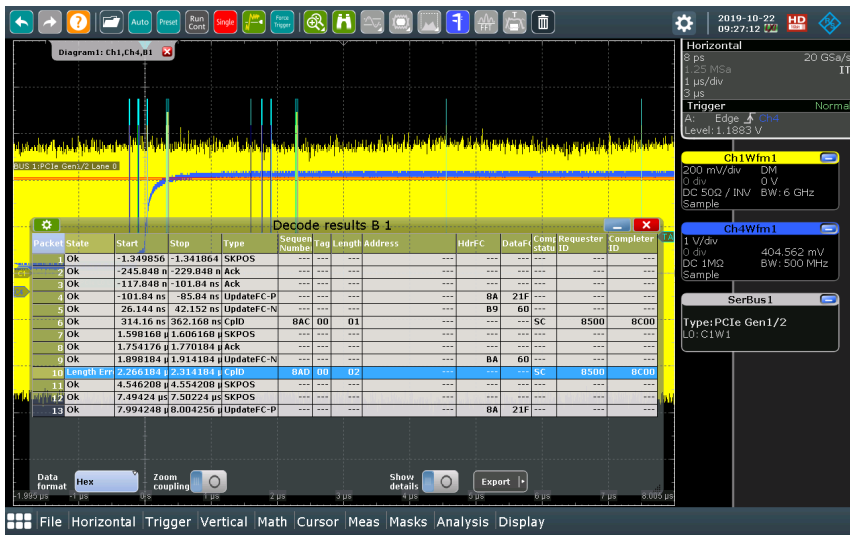
Example 1: measurement

TX side on DAMC-TCK7 (with Kintex 7)



Example 1: measurement

RX side on DAMC-TCK7 (with Kintex 7)



Example 1: measurement

RX side on DAMC-TCK7 (with Kintex 7) - **length error**



- ▶ Triggering with GPIO pin is a convenient way to "isolate" relevant events
- ▶ Difference between 7-series and UltraScale xdma
- ▶ Work-around for 7-series xdma: `mempcy()` based on 32-bit reads

It was observed that after some time (usually less than an hour) the DMA stops working - every time we attempt to perform a DMA transfer the CPU crashes. It was soon notice that

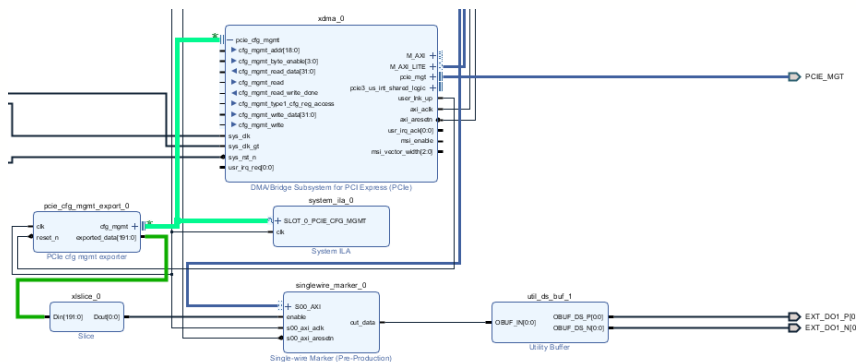
everything works OK until first error report appears in `dmesg`. This error report itself is nothing unusual - HW handled the error on its own.

```
1555.357788] pciport 0000:89:00.0: AER: [15] HeaderOF
1555.365082] pciport 0000:89:00.0: AER: aer_layer=Physical Layer, aer_agent=Receiver ID
1923.983565] {3}[Hardware Error]: Hardware error from APEI Generic Hardware Error Source: 0
1923.983568] {3}[Hardware Error]: It has been corrected by h/w and requires no further action
1923.983570] {3}[Hardware Error]: event severity: corrected
1923.983572] {3}[Hardware Error]: Error 0, type: corrected
1923.983574] {3}[Hardware Error]: section_type: PCIe error
1923.983576] {3}[Hardware Error]: port_type: 5, upstream switch port
1923.983577] {3}[Hardware Error]: version: 3.0
1923.983579] {3}[Hardware Error]: command: 0x0147, status: 0x0010
1923.983581] {3}[Hardware Error]: device_id: 0000:89:00.0
1923.983582] {3}[Hardware Error]: slot: 0
1923.983584] {3}[Hardware Error]: secondary_bus: 0x8a
1923.983585] {3}[Hardware Error]: vendor_id: 0x10b5, device_id: 0x8780
1923.983587] {3}[Hardware Error]: class_code: 000406
1923.983588] {3}[Hardware Error]: bridge: secondary_status: 0x0000, control: 0x0013
1923.983823] pciport 0000:89:00.0: AER: aer_status: 0x0000a001, aer_mask: 0x00002000
1923.992489] pciport 0000:89:00.0: AER: [ 0] RxErr
1923.999792] pciport 0000:89:00.0: AER: [15] HeaderOF
1924.007086] pciport 0000:89:00.0: AER: aer_layer=Physical Layer, aer_agent=Receiver ID
1985.423861] {4}[Hardware Error]: Hardware error from APEI Generic Hardware Error Source: 0
1985.423864] {4}[Hardware Error]: It has been corrected by h/w and requires no further action
```

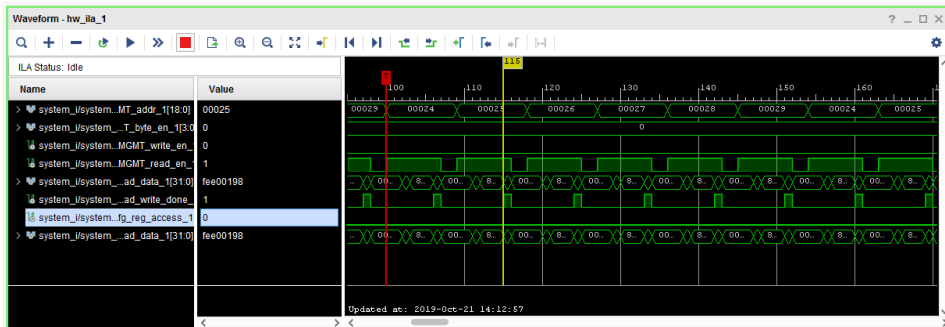

We noticed that MSI vector gets corrupted (we performed `diff` on `lspci` output before and after first "Hardware Error" message).

```
$ sudo lspci -s 8c:00 -vv
8c:00.0 Serial controller: Xilinx Corporation Device 8021 (prog-if 01 [16450])
  Subsystem: Xilinx Corporation Device 0007
  [...]
  NUMA node: 1
  Region 0: Memory at e0b00000 (32-bit, non-prefetchable) [size=1M]
  Region 1: Memory at e0c00000 (32-bit, non-prefetchable) [size=64K]
  Capabilities: [80] Power Management version 3
    Flags: PMEClk- DSI- D1- D2- AuxCurrent=0mA PME(D0-,D1-,D2-,D3hot-,D3cold-)
    Status: D0 NoSoftRst+ PME-Enable- DSel=0 DScale=0 PME-
  Capabilities: [90] MSI: Enable+ Count=1/1 Maskable- 64bit+
    Address: 00080000fee00198 Data: 0000
  [...]
  Capabilities: [100 v2] Advanced Error Reporting
    UESSta: DLP- SDES- TLP- FCP- CmplTTO- CmpltAbrt- UnxCmplt- RxOF- MalfTL -
    UEMsk: DLP- SDES- TLP- FCP- CmplTTO- CmpltAbrt- UnxCmplt- RxOF- MalfTL -
    UESvrt: DLP+ SDES+ TLP- FCP+ CmplTTO- CmpltAbrt- UnxCmplt- RxOF+ MalfTL +
    CESTa: RxErr- BadTLP- BadDLLP- Rollover- Timeout- AdvNonFatalErr-
    CEMsk: RxErr- BadTLP- BadDLLP- Rollover- Timeout- AdvNonFatalErr-
    AERCap: First Error Pointer: 00, ECRGenCap- ECRGenEn- ECRChkCap- ECR
    MultHdrRecCap- MultHdrRecEn- TLPPfxPres- HdrLogCap-
    HeaderLog: 04000001 85000003 8c010000 00000000
  Kernel driver in use: xdma
```

We prepared a test setup to trigger the scope on this event - an IP (PCIe config mgmt exporter) continuously copies the values from the PCIe endpoint configuration area to a vector. This vector is then used to drive LVDS output on the front panel.



Here we observe the connection between **PCIe config mgmt exporter** and PCIe IP core - we see that **exporter** is cyclically reading relevant addresses.



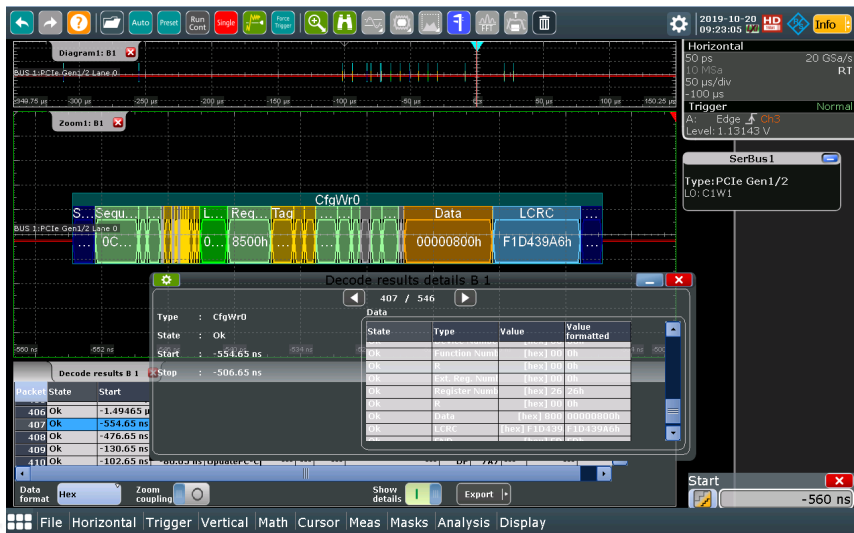
Example 2: capture overview

We arm the trigger and wait. After some time, the oscilloscope triggers:



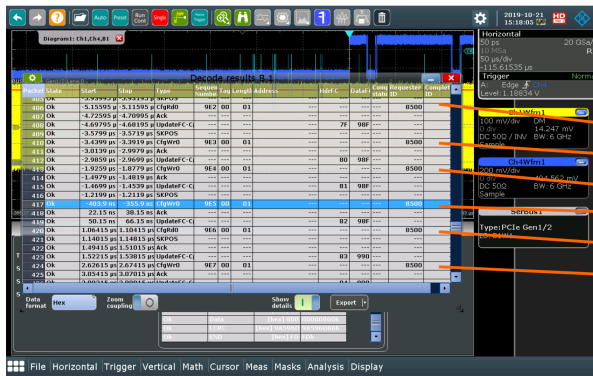
Example 2: the packet

This is the packet which corrupts MSI vector! It means that the source of this issue is on the software side - probably in Linux driver.



Example 2: analysis

We can analyze the configuration packets and compare them with software - this allows us pin-pointing the source of the erroneous packet.



- rd @ 0x0 [-15.8 us]
- rd @ 0x4 [-14.3 us]
- rd @ 0x4 [-12.7 us]
- rd @ 0x34 [-11.2 us]
- rd @ 0x80 [- 9.7 us]
- rd @ 0x90 [- 8.2 us]
- rd @ 0xC0 [- 6.6 us]
- rd @ 0x100 [- 5.1 us]
- wr @ 0x104 [- 3.5 us]
- wr @ 0x110 [- 1.9 us]
- wr @ 0x98 [- 0.4 us]**
- rd @ 0x108 [1.1 us]
- wr @ 0x108 [2.6 us]
- rd @ 0x114 [4.1 us]
- wr @ 0x114 [5.6 us]
- rd @ 0xC [7.2 us]
- rd @ 0x104 [8.8 us]
- rd @ 0x104 [10.4 us]



- ▶ Using a clever triggering method we were able to observe the relevant packets
- ▶ Protocol decode enables us to understand software better/easier
- ▶ Long buffers in oscilloscope (10 M samples) allows capturing large SW routines
- ▶ This debug is on-going, but we are very close to resolving the problem

Summary

- ▶ L. Petrosyan, "MicroTCA and PCIe HotSwap under Linux",
<http://webcast.desy.de/?m=201312&lang=en>,
at MTCA Workshop for Industry and Research 2013, Hamburg
- ▶ J. Marjanovic, "Direct Memory Access with FPGA",
https://github.com/MicroTCA-Tech-Lab/damc-tck7-fpga-bsp/blob/demo-mtcaws2018/docs/demo/mtcaws2018_marjanovic_dma.pdf
at MTCA Workshop for Industry and Research 2018, Hamburg
(similar to this one, but more hands-on - Vivado project and SW available)
- ▶ J. Corbet, A. Rubini, G. Kroah-Hartman, "Linux Device Drivers, 3rd Edition",
(quite old at this point, 4th Edition long over-due)
- ▶ M. Jackson, R. Budruk, "PCI Express Technology 3.0"

- ▶ PCI Express is a protocol of choice for many institutes and vendors
- ▶ Provides low latency and high throughput
- ▶ FPGA vendors provide a lot of IP, drivers, ...
- ▶ MicroTCA Tech Lab provides BSP with PCIe for some of the boards, more BSPs (e.g. for Zynq UltraScale MPSoc) to follow
- ▶ There are a lot of tools to troubleshoot PCIe-based systems



TRANSFER MTCA TO RESEARCH AND INDUSTRY

- ▶ Custom developments
- ▶ High-end test & measurement services
- ▶ System configuration & integration
- ▶ LLRF design

Marketing.
Services & Support.
Tech-Shop.

**Thank you
Vielen Dank**

<https://techlab.desy.de>

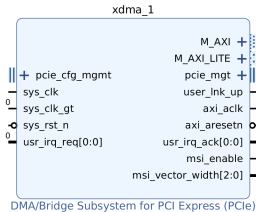
Deutsches Elektronen-Synchrotron DESY
A Research Centre of the Helmholtz Association
Notkestr. 85, 22607 Hamburg, Germany



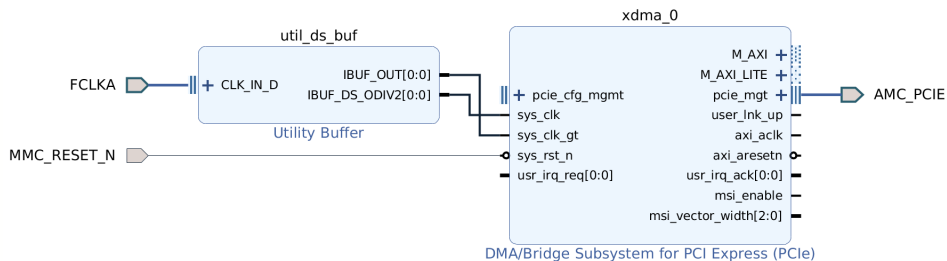
Backup slides

FPGA side

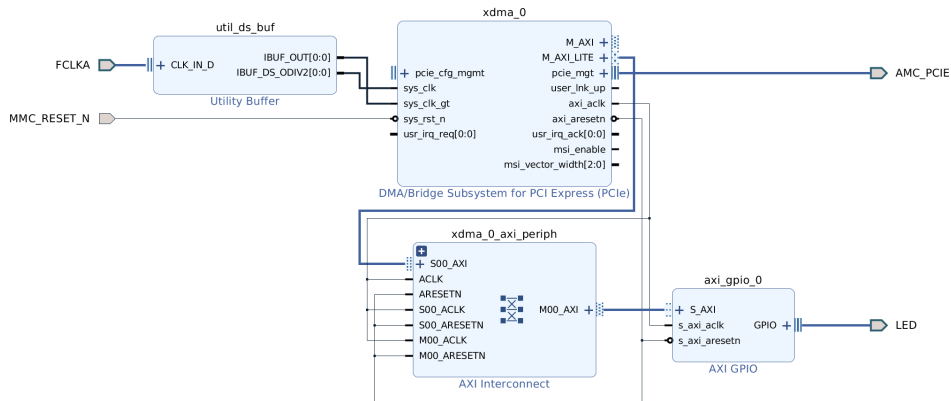
1. We open Vivado, start a new project, create a Block Diagram and place an xdma IP:



2. We then connect external connections:



3. We then proceed to add a GPIO, to blink the LEDs



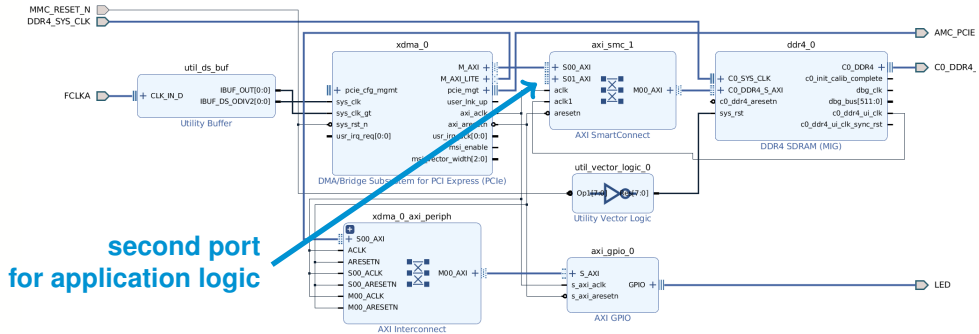
4. Let's not forget to assign the addresses

The screenshot shows the BLOCK DESIGN software interface. The main window is titled "BLOCK DESIGN" and contains two panes. The left pane, "Sources", shows a tree view of the design hierarchy for "design_1", including "External Interfaces" (AMC_PCIE, FCLKA, LED) and "Interface Connections" (axi_gpio_0_GPIO, diff_clock_rt1_0_1, xdma_0_axi_periph_M00_AXI, xdma_0_M_AXI_LITE, xdma_0_pcie_mgt). The right pane, "Address Editor", displays a table with the following data:

Cell	Slave Interface	Base Name	Offset Address	Range	High Address
▼ xdma_0					
		M_AXI (64 address bits : 16E)			
▼		M_AXI_LITE (32 address bits : 4G)			
	S_AXI	Reg	0x0001_0000	4K	- 0x0001_0FFF

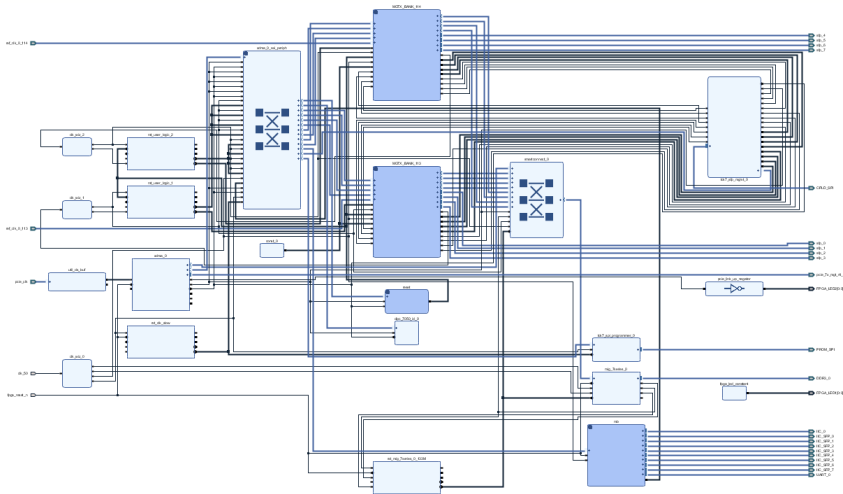
A blue tip is overlaid on the image: "Mini tip: place at a know address (e.g. at 0x0) ident register (e.g. 0xde30b02d)".

5. We add a DDR4 controller (not a trivial step, a lot of configuration) and connect the DMA port of PCIe to DDR4 controller:

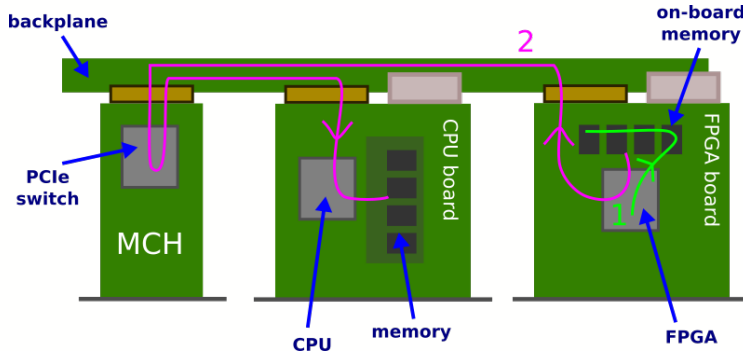


6. some magic (aka months of work on application logic)

7. This is how a typical PCIe-based system looks

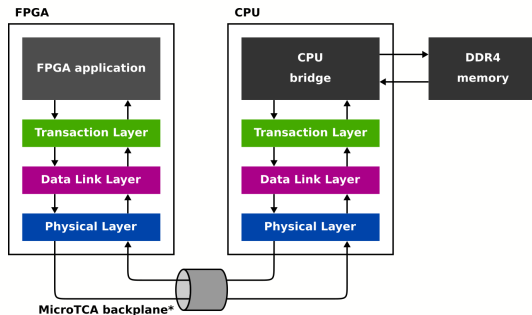


- ▶ Modern FPGAs provide a lot of hard IP
- ▶ Typical flow involves configuring IPs and connecting them together
- ▶ Custom logic required for application-specific part
- ▶ Two-step transfer: first to on-board memory, then to CPU



Software side

From "FPGA side" we have solved the left part, now it is time to look at the right part:



- ▶ Typically there is one (or more) CPU AMCs in MicroTCA systems
- ▶ Usually (in experimental physics) the CPU runs GNU/Linux
- ▶ DESY runs Ubuntu 16.04 LTS, Ubuntu 18.04 LTS is also in use
- ▶ PCIe hotplug was modernized in Linux kernel 4.19[1] → available in Ubuntu 20.04 LTS

[1] <https://lwn.net/Articles/767885/>

- ▶ Linux has a subsystem responsible for managing PCI/PCIe device driver
- ▶ Hot-swap module (pciehp) is also part of Linux kernel
- ▶ Device-specific driver register by providing `pci_driver` structure to `pci_register_driver()` function

<https://elixir.bootlin.com/linux/v4.15/source/include/linux/pci.h#L744>

```
struct pci_driver {
    struct list_head node;
    const char *name;
    const struct pci_device_id *id_table; /* must be non-NULL for probe to be called */
    int (*probe) (struct pci_dev *dev, const struct pci_device_id *id); /* New device
        inserted */
    void (*remove) (struct pci_dev *dev); /* Device removed (NULL if not a hot-plug
        capable driver) */
    int (*suspend) (struct pci_dev *dev, pm_message_t state); /* Device suspended */
    int (*suspend_late) (struct pci_dev *dev, pm_message_t state);
    int (*resume_early) (struct pci_dev *dev);
    int (*resume) (struct pci_dev *dev); /* Device woken up */
    void (*shutdown) (struct pci_dev *dev);
    int (*sriov_configure) (struct pci_dev *dev, int num_vfs); /* PF pdev */
    const struct pci_error_handlers *err_handler;
    const struct attribute_group **groups;
    struct device_driver driver;
    struct pci_dynids dynids;
};
```

Xilinx driver only provides `probe()` and `remove()` functions - enough to perform hot-swap.

from https://github.com/Xilinx/dma_ip_drivers/blob/8b8c70b697f049649d5fa99be9c6bc4302d89ac9/XDMA/linux-kernel/xdma/xdma_mod.c#L316:

```
static struct pci_driver pci_driver = {
    .name = DRV_MODULE_NAME,
    .id_table = pci_ids,
    .probe = probe_one,
    .remove = remove_one,
    .err_handler = &xdma_err_handler,
};
```

On the other side (towards the user space), Xilinx DMA driver provides several char devices:

```
$ ll /dev | grep xdma
drwxr-xr-x  3 root root          60 Jun 20 17:19 xdma
crw-rw-rw-  1 root root    238, 36 Jun 20 17:19 xdma0_c2h_0
crw-rw-rw-  1 root root    238,  1 Jun 20 17:19 xdma0_control
crw-rw-rw-  1 root root    238, 10 Jun 20 17:19 xdma0_events_0
crw-rw-rw-  1 root root    238, 11 Jun 20 17:19 xdma0_events_1
[...]
crw-rw-rw-  1 root root    238, 24 Jun 20 17:19 xdma0_events_14
crw-rw-rw-  1 root root    238, 25 Jun 20 17:19 xdma0_events_15
crw-rw-rw-  1 root root    238, 32 Jun 20 17:19 xdma0_h2c_0
crw-rw-rw-  1 root root    238,  0 Jun 20 17:19 xdma0_user
crw-rw-rw-  1 root root    238,  2 Jun 20 17:19 xdma0_xvc
```

On the other side (towards the user space), Xilinx DMA driver provides several char devices:

```
$ ll /dev | grep xdma
drwxr-xr-x  3 root root          60 Jun 20 17:19 xdma
crw-rw-rw-  1 root root    238, 36 Jun 20 17:19 xdma0_c2h_0
crw-rw-rw-  1 root root    238,  1 Jun 20 17:19 xdma0_control
crw-rw-rw-  1 root root    238, 10 Jun 20 17:19 xdma0_events_0
crw-rw-rw-  1 root root    238, 11 Jun 20 17:19 xdma0_events_1
[...]
crw-rw-rw-  1 root root    238, 24 Jun 20 17:19 xdma0_events_14
crw-rw-rw-  1 root root    238, 25 Jun 20 17:19 xdma0_events_15
crw-rw-rw-  1 root root    238, 32 Jun 20 17:19 xdma0_h2c_0
crw-rw-rw-  1 root root    238,  0 Jun 20 17:19 xdma0_user
crw-rw-rw-  1 root root    238,  2 Jun 20 17:19 xdma0_xvc
```

c2h

AXI4 port (DMA from device)

h2c
AXI4 port (DMA to device)

user
AXI4-Lite port

events_N
AXI4 port (DMA from device)

To access the status and control registers (on AXI4-Lite interface), C variant:

```
#define ADDR_STATUS_REG          (0x1020)

int fd = open("/dev/xdma/card0/user", O_RDWR | O_SYNC);
if (fd < 0) {
    perror("open()");
    return EXIT_FAILURE;
}

void* mem = mmap(0, mmap_size, PROT_READ | PROT_WRITE, MAP_SHARED, fd, mmap_offset);
if (mem == MAP_FAILED) {
    perror("mmap()");
    return EXIT_FAILURE;
}

uint32_t status = *(uint32_t*)((char*)mmap_base + ADDR_STATUS_REG);
```

To access the status and control registers (on AXI4-Lite interface), Python variant

```
class HwAccessPcie(object):

    def __init__(self):
        user_filename = "/dev/xdma/card0/user"
        self.fd_user = os.open(user_filename, os.O_RDWR)
        self.mem = mmap.mmap(self.fd_user, 4 * 1024 * 1024)

    def __close__(self):
        self.mem.close()
        os.close(self.fd_user)

    def _rd32(self, addr):
        bs = self.mem[addr:addr + 4]
        return struct.unpack("I", bs)[0]

    def _wr32(self, addr, data):
        bs = struct.pack("I", data)
        self.mem[addr:addr + 4] = bs
```

- ▶ Hot-plug requires that virtual memory is additional kernel arguments:

```
pci=realloc, assign-busses, hpmemsize=32M
```

- ▶ Listing of a hot-plug procedure is shown on the next slide


```
[ 11.719854] xdma: loading out-of-tree module taints kernel.
[ 11.719902] xdma: module verification failed: signature and/or required key missing - tainting kernel
[ 11.720630] xdma:xdma_mod_init: Xilinx XDMA Reference Driver xdma v2018.3.50
[ 11.720631] xdma:xdma_mod_init: desc_blen_max: 0xffffffff/268435455, sgdma_timeout: 10 sec.
[ 12.938964] JAN: plug in AMC and push in the handle
[ 26.347956] pciehp 0000:04:08.0:pcie204: Slot(5): Attention button pressed
[ 26.347978] pciehp 0000:04:08.0:pcie204: Slot(5) Powering on due to button press
[ 26.505506] pciehp 0000:04:08.0:pcie204: Slot(5): Link Up
[ 27.420194] pci 0000:07:00.0: [10ee:7024] type 00 class 0x070001
...
[ 27.420471] pci 0000:07:00.0: BAR 0: assigned [mem 0xa6000000-0xa67fffff]
[ 27.420477] pci 0000:07:00.0: BAR 1: assigned [mem 0xa6800000-0xa680ffff]
[ 27.420482] pcieport 0000:04:08.0: PCI bridge to [bus 07]
[ 27.420485] pcieport 0000:04:08.0: bridge window [io 0x4000-0x4fff]
[ 27.420489] pcieport 0000:04:08.0: bridge window [mem 0xa6000000-0xa9ffffff]
[ 27.420492] pcieport 0000:04:08.0: bridge window [mem 0x96000000-0x97ffffff 64bit pref]
...
[ 27.420944] xdma:map_single_bar: BAR0 at 0xa6000000 mapped at 0x000000004c4894e6, length=83886
[ 27.420977] xdma:map_single_bar: BAR1 at 0xa6800000 mapped at 0x00000000b314fa42, length=65536
[ 27.420980] xdma:mapBars: config bar 1, pos 1.
[ 27.420981] xdma:identifyBars: 2 BARs: config 1, user 0, bypass -1.
[ 27.421038] xdma:probe_one: 0000:07:00.0 xdma0, pdev 0x00000000a8508428, xdev 0x00000000affcab
[ 27.421833] xdma:cdev_xvc_init: xcdev 0x00000000b85138f3, bar 0, offset 0x40000.
```

- ▶ Drivers are provided by FPGA vendors
- ▶ Easy-to-use user-space interface (for register access, DMA and interrupts)
- ▶ Hot-swap facilitates development and improves up-time