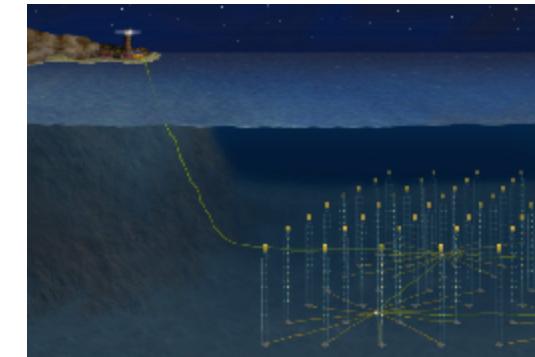
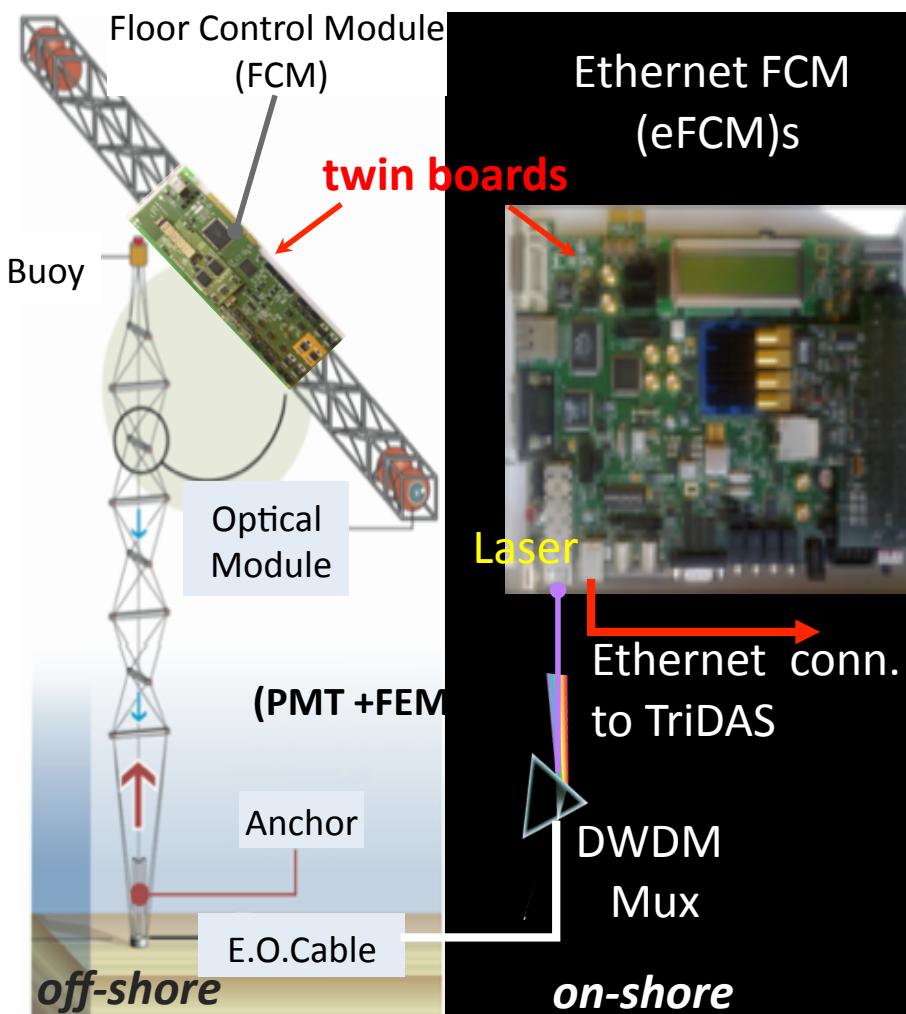


# Stato della scheda NaNet<sup>3</sup>

Alessandro Lonardo  
INFN - Sezione di Roma

# KM3Net-IT Experiment

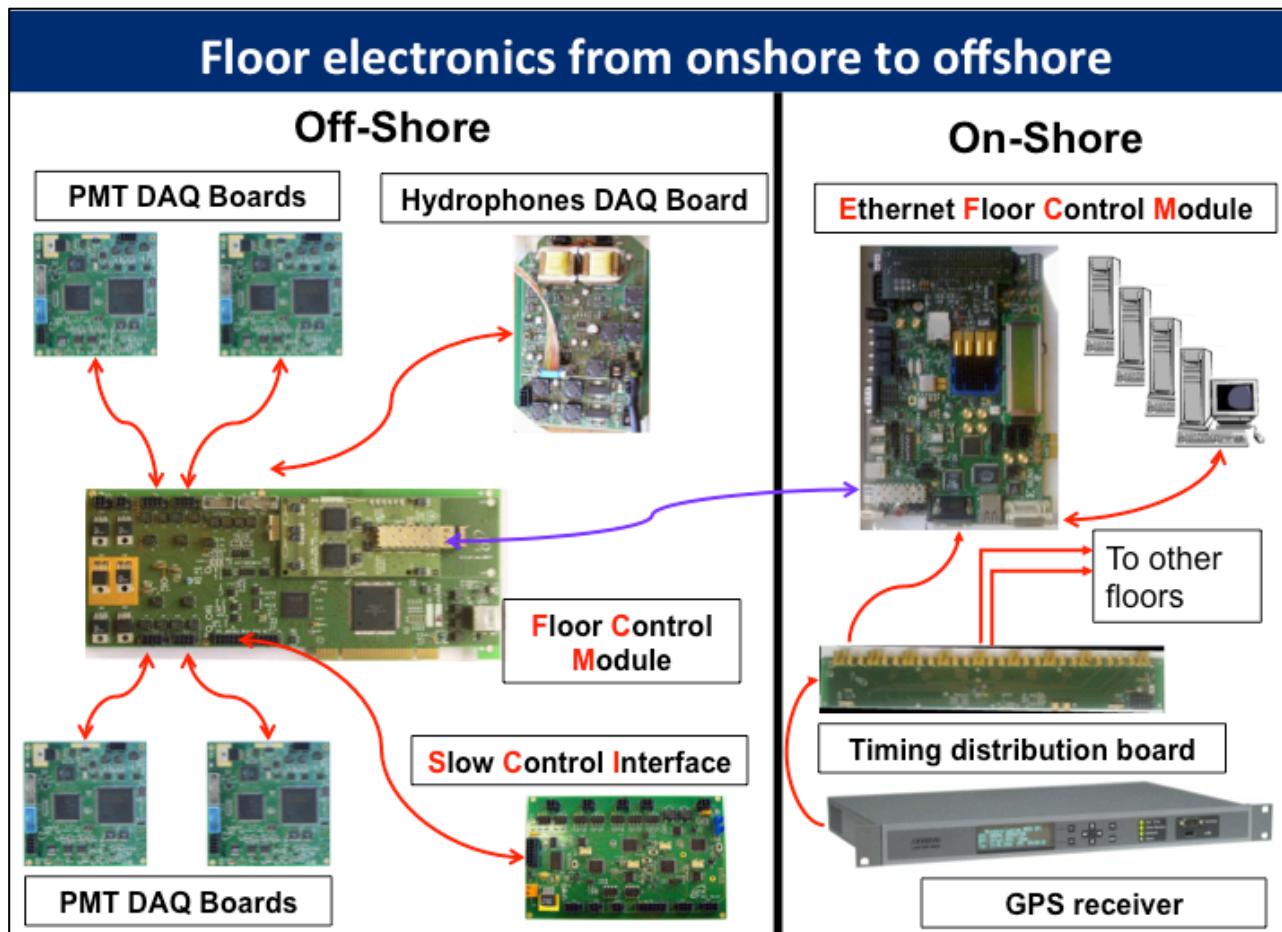
- **KM3Net-IT:** an European deep-sea research infrastructure, hosting a neutrino telescope with a volume of cubic kilometer at the bottom of the Mediterranean Sea.



- “Phase-2 tower” composed of:
  - 8 floors
  - 8 m bars,
  - vertical dist. = 40 m, H<sub>tot</sub> = 450 m
  - 32 OM, 18 hydrophones
  - oceanographic instrumentation
  - 2 twinned (off-shore/on-shore) FCM boards per floor.

# KM3Net-IT experiment: read-out

- Current **read-out system** employs a huge number, NO state-of-the-art components
  - 2.5 Gbps (800 Mbps payload) optical link
  - 2 twinned FCM boards per floor
  - Many PCs for HW read-out hosting



F. Simeone  
IEEE/NSS 2011

- When scaling to KM3 many cost/size/power/reliability issues!!!

# APEnet+: a 3D NIC for HPC with GPUdirect P2P capability

## APEnet+ Card:

- **FPGA based (ALTERA EP4SGX290)**
- PCI Express x16 slot, x8 gen2 signaling
- Fully bidir 3D torus links, 34 Gbps/channel

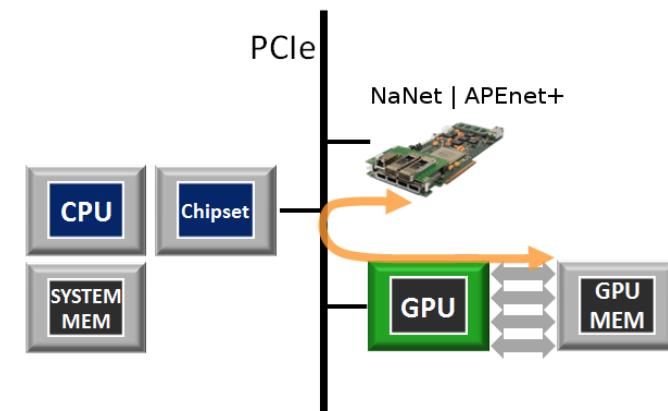
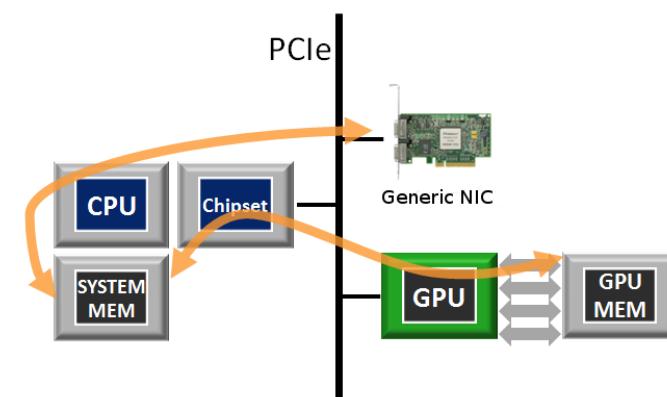
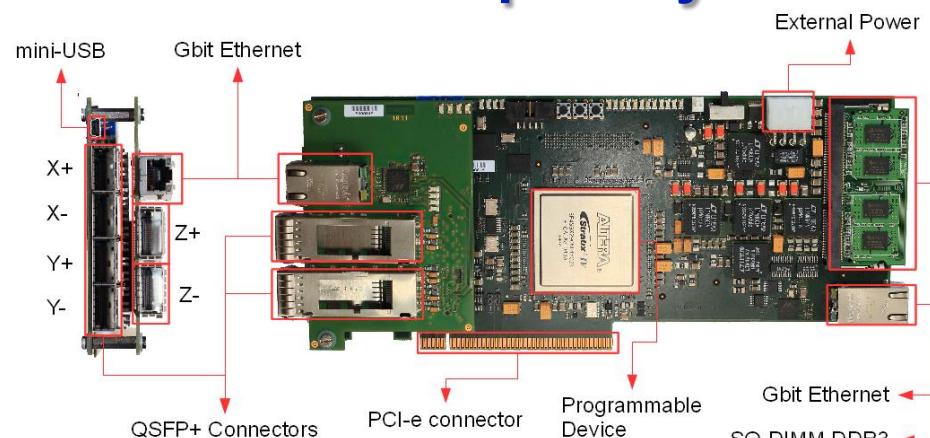
## APEnet+ Logic:

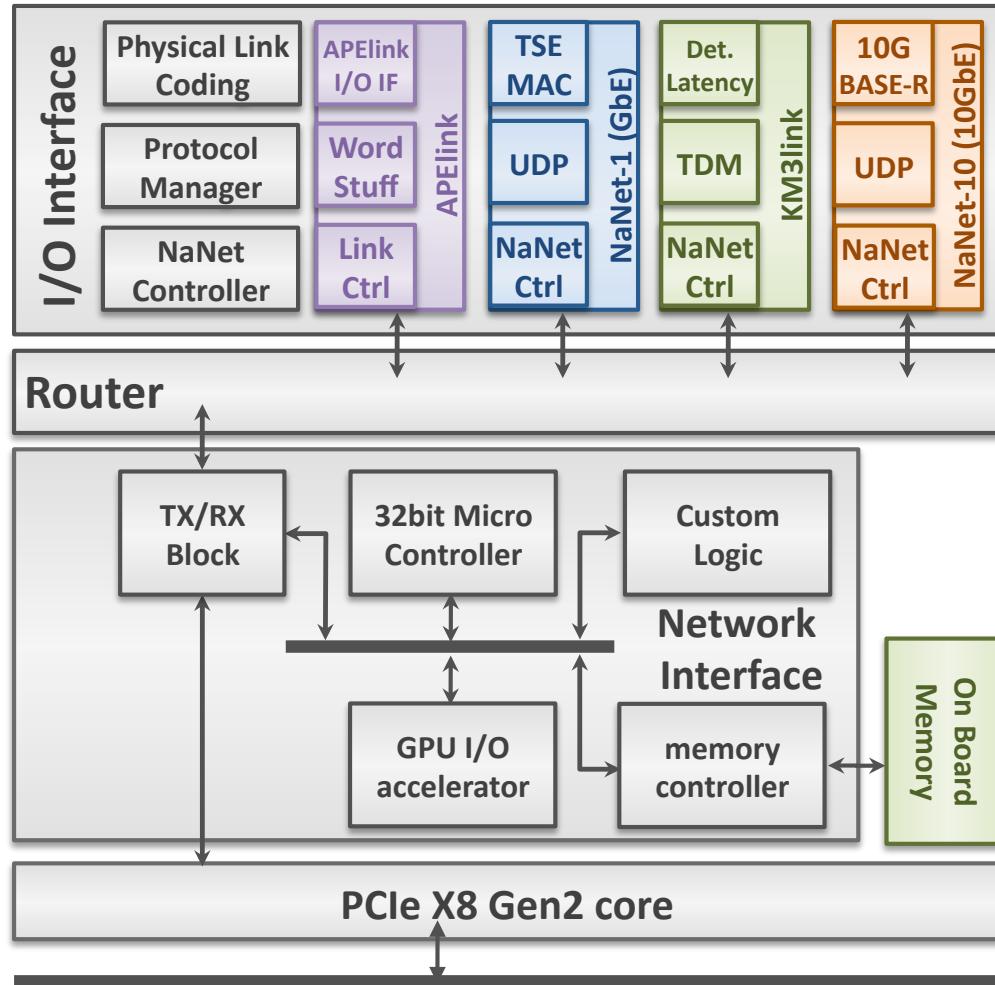
- Network Interface
  - NIOS II 32 bit uP
  - RDMA engine
  - Virtual Memory Management support
  - GPU I/O accelerator.

## GPUDirect P2P

- PCIe P2P protocol between NVIDIA Fermi/Kepler devices and APEnet+
- GPUDirect P2P allows direct data exchange on the PCIe bus with no CPU involvement
 

→ *Latency reduction for small messages*





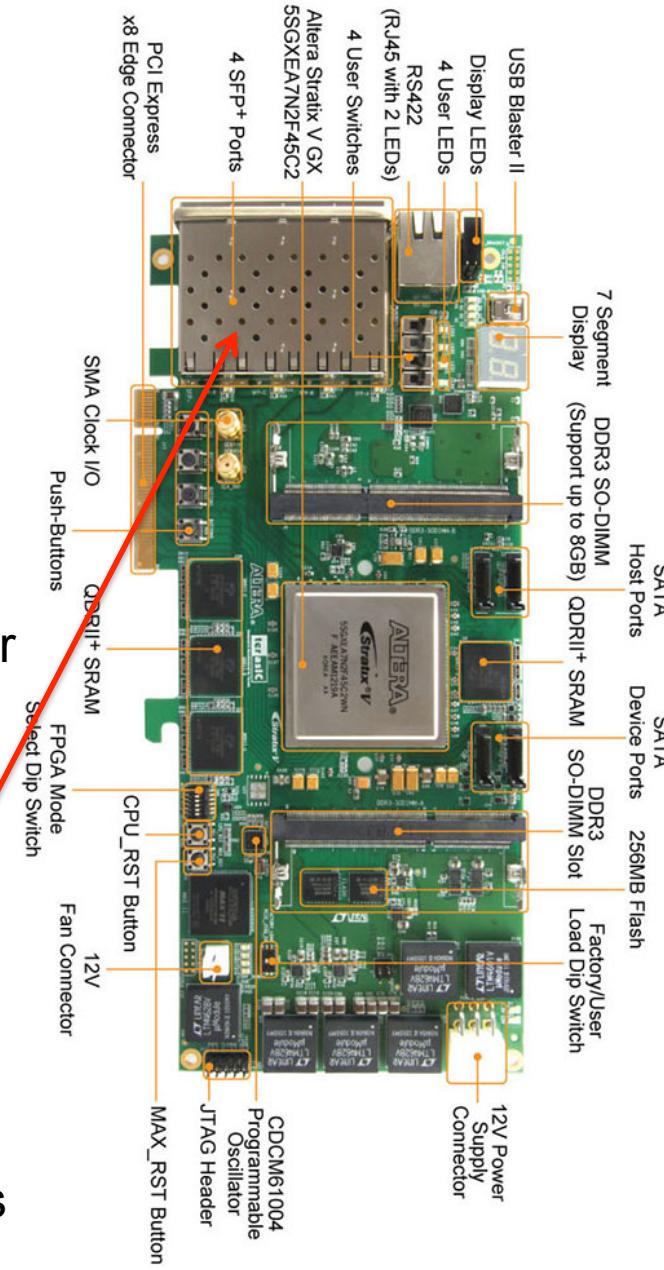
## NaNet

- **Multistandard support of I/O interface**
  - Off-the-Shelf: 1-GbE, 10-GbE (work in progress).
  - Custom: APElink (34gbps/QSFP), KM3link (work in progress).
- GPU I/O accelerator
  - GPUDirect v2 (custom NVP2P)
  - GPUDirect RDMA
- CPU offload – NO OS Jitter
  - Transport layer off-loading module
  - Virtual Memory Management
- NaNet Controller: protocol adaptation between on-board and off-board protocol
- Altera NIOS II microcontroller
- Custom Logic (application-specific task, e.g. compression)
- DDR3 memory controller
- PCIe Gen2 X8 host interface

# Why NaNet<sup>3</sup>?

## NaNet<sup>3</sup> specifications for future enhanced read-out system

- more (4) channels/PCIe board
  - i.e. less PC, less read-out boards,... ✓
- Link speed-up to 10Gb/s
  - to get 4:1 multiplexing ratio use an optional “Octolink” aggregating board ...NO
- GPU Direct
  - for support to future enhanced read-out and trigger system GPU-based... ✓
- “Fixed Latency” clock distribution for under-water events time-stamping
  - OK for FCM (Xilinx) ✓
  - **Altera Stratix IV/V** ✓



## Final design on TERASIC DE5-NET

- Altera Stratix V based dev board
- PCIe Gen2 x8 (Gen3 support)
- 4 independent SFP+ ports (up to 10Gb/s)
- Deterministic Latency mode for transceivers

# NaNet<sup>3</sup>: State of the art

- **NaNet<sup>3</sup>link (3 up to now):**

- Physical Layer: Altera Deterministic Latency Transceivers (8B10B encoding scheme)
- Data Layer: Time Division Multiplexing (TDM) data transmission protocol.
  - RX path: payload of different off-shore devices, multiplexed on continuous data stream at fixed time slot
    - (PCIe DMA transaction to CPU/GPU memory)
  - TX path: limited data rate per FCM
    - (PCIe TARGET transaction from CPU/GPU memory)

- **NaNet Ctrl:**

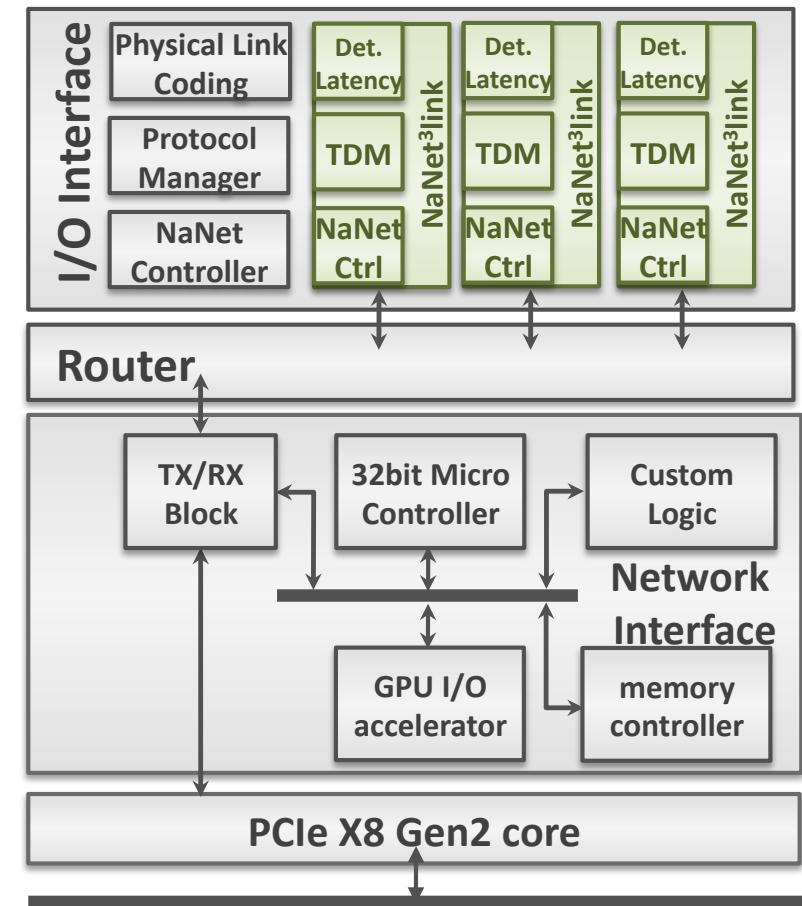
- protocol translation; encapsulate TDM data stream in APEnet protocol
- Virtual Memory management (CPU/ $\mu$ C offloading)

- **Virtual-to-Physical Translation**

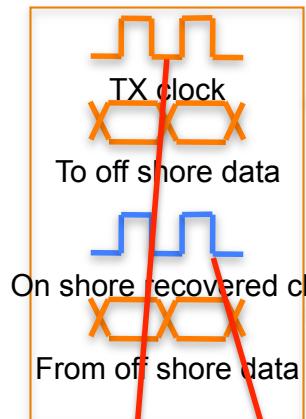
- Nios Implementation
- HW acceleration: Translation Lookaside Buffer (TLB) based on associative memory

- **PCIe X8 Gen2**

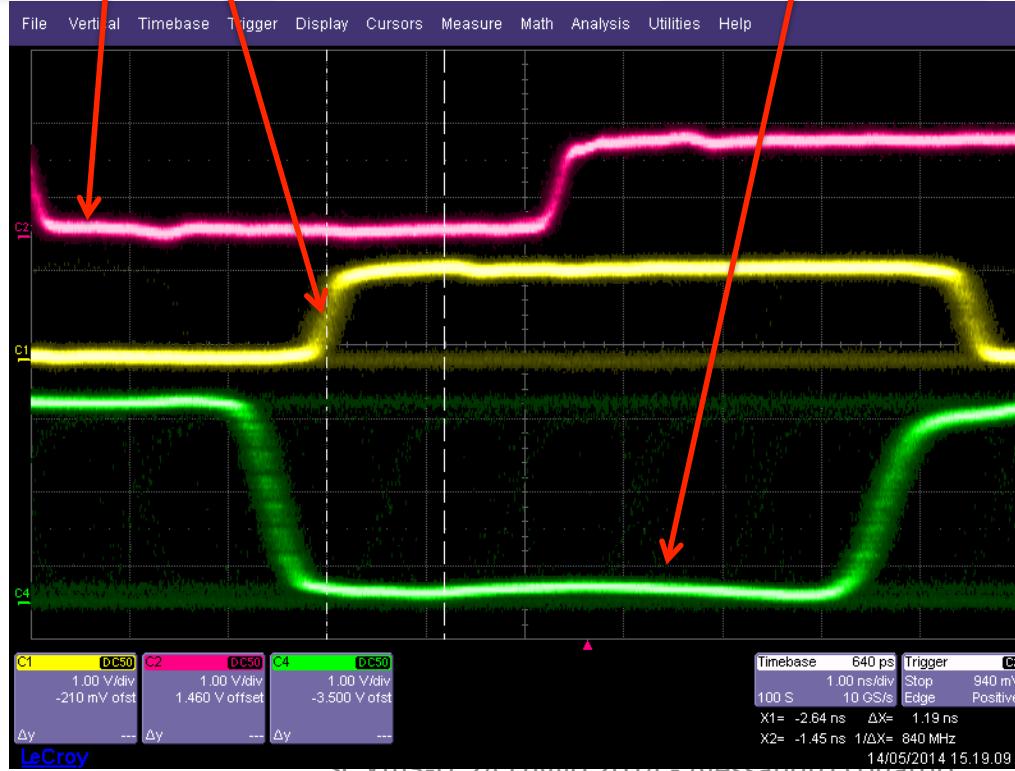
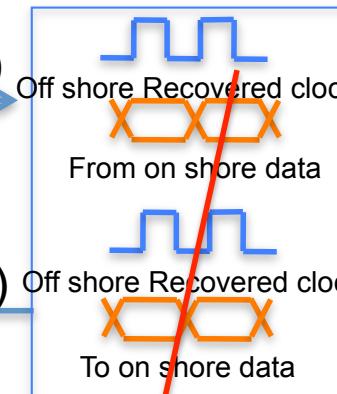
- CPU/GPU Memory Write Bandwidth  $\sim 2.4$  GB/s



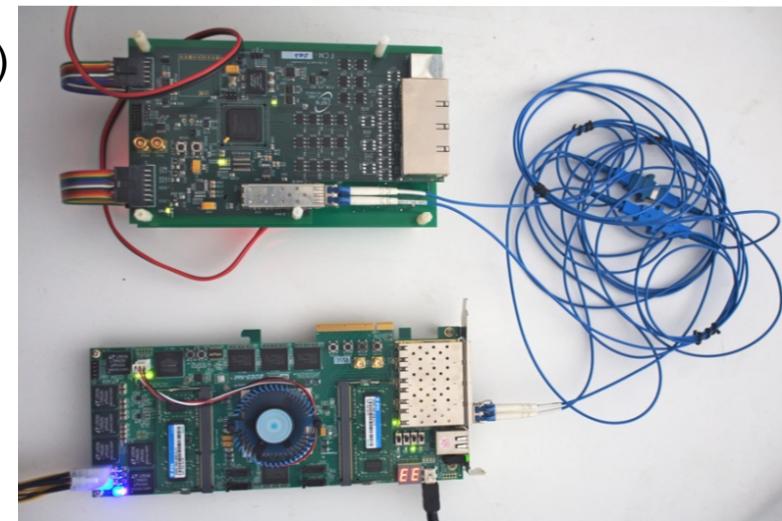
## NaNet3 On-shore (StratixV)



## FCM Off-Shore (Virtex5)



# Deterministic latency link inter-operability



## Preliminary but encouraging results

- Testbed: FCM vs Terasic DE5-Net
  - Custom hw mode for FCM Transceivers (Xilinx)
  - Latency deterministic mode for Stratix V Transceiver
  - 2mt copper and 2 mt long fiber
- Test:
  - 12 hours of periodic (~s) Tx clock reset to verify pll locking and rx word alignment

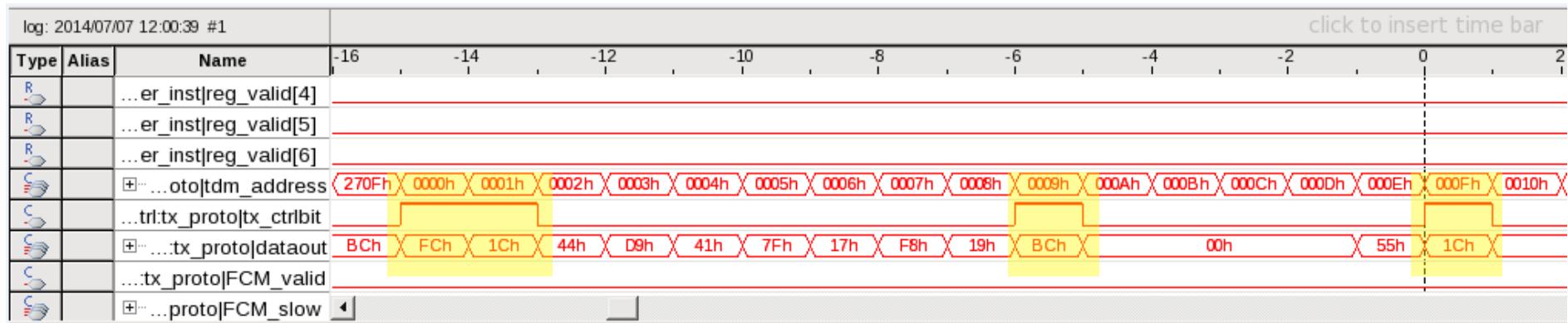
# Time Division Multiplexing (TDM) support

## NaNet I/O interface Customization:

- TDM implementation: compliant with KM3Net-IT specification, echo test OK!

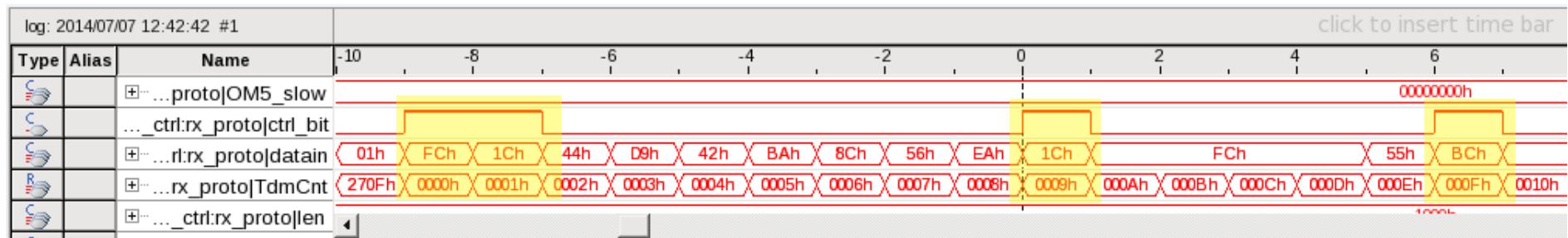
	K Code	Kout	D
IDLE_CODE	K28.5	1	0xBC
FRAME_CODE	K28.7	1	0xFC
DATA_CODE	K28.0	1	0x1C

- NaNet<sup>3</sup> TX



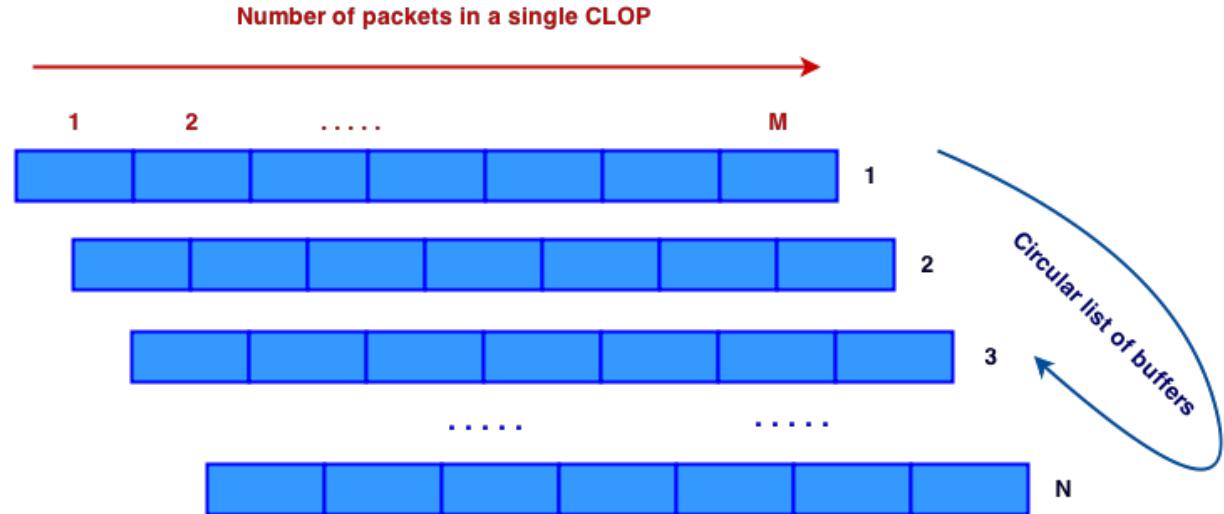
- NaNet<sup>3</sup> RX:

Data flow: NaNet<sup>3</sup> → FCM board → NaNet<sup>3</sup>

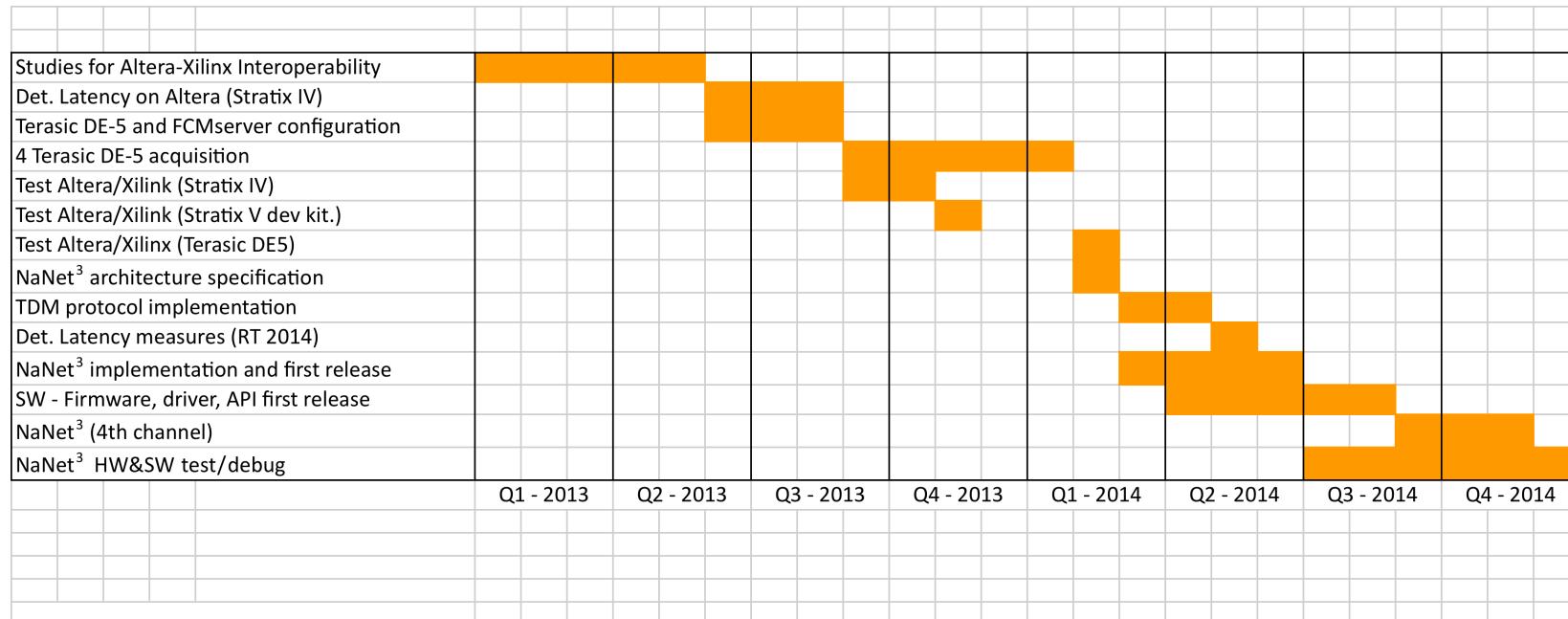


# NaNet<sup>3</sup> Software

- **Linux Kernel Driver**
  - Status/Configuration registers.
  - TX registers interface.
- **NIOS II Firmware**
  - New BSP for NaNet3 board.
  - Initialization of NaNet3 channels.
  - Management of 4 concurrent data streams
- **Application Library**
  - `nan3_t nan3_open(int card_id); int nan3_close(nan3_t nan);`
  - `int nan3_register_clop(nan3_t nan, nan3 Chan_id_t chan, u32 pkts, u32 items, u8 is_gpu);`
  - `int nan3_wait_event(nan3_t nan, nan3_Event_recv_t* event);`
  - `u32 nan3_write_slow_data(nan3_t nan, u8 channel, u8 deviceld, u32 data);`
- **Application Template**



# NaNet<sup>3</sup>: Roadmap



- Dec 2012: Piero Vicini talk on NaNet in Catania.
- Aug 2013: Deterministic Latency on Altera Stratix IV dev kit.
- Oct 2013: Interoperability Altera/Xilinx.
- Feb 2014: Deterministic Latency on Terasic DE5 (Stratix V).
- **Mid Feb 2014: NaNet<sup>3</sup> architecture specification.**
- Jul 2014: first release of NaNet<sup>3</sup> of complete HW design (3 channels).
- Aug 2014: first release of system software and application library.

## Budget & People

- NaNet HW and SW developments financed by FP7 EURETILE project (exploitation activities): 8 PM in 2013, 9 PM in 2014.
- EURETILE ending on September 2014.
- People involved
  - HW
    - Andrea Biagioni
    - Ottorino Frezza
    - Francesca Lo Cicero
    - Piero Vicini (Staff)
  - SW
    - Alessandro Lonardo (Staff)
    - Michele Martinelli
- Integrated remaining man power on HW side : 0.5 MM before the end of the project.
- People available to continue activities after September (if supported):
  - Integration of 4<sup>th</sup> channel.
  - Test and debug
  - ...



# BACKUP SLIDES

## NaNet<sup>3</sup> testbeds



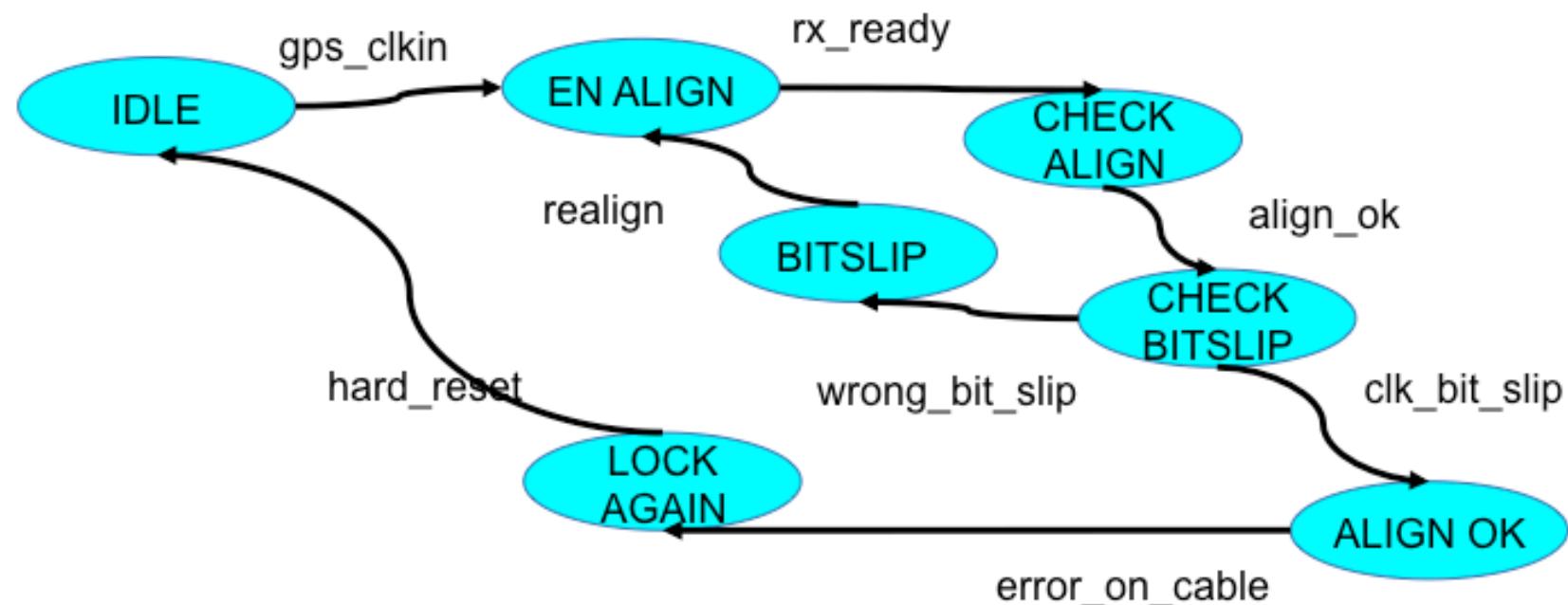
## Deterministic latency link inter-operability

<b>NaNet<sup>3</sup> on-shore</b>	Tx	Send Magic word	Sending Magic word	Sending Magic word	Sending data
	Rx	Under reset	Under reset	Under reset	Clock recovered with Det. Lat.
<b>FCM offshore</b>	Tx	Idle	Idle	Sending Magic word	Sending data
	Rx	Waiting for CDR	Clock recovered	Clock recovered	Clock recovered

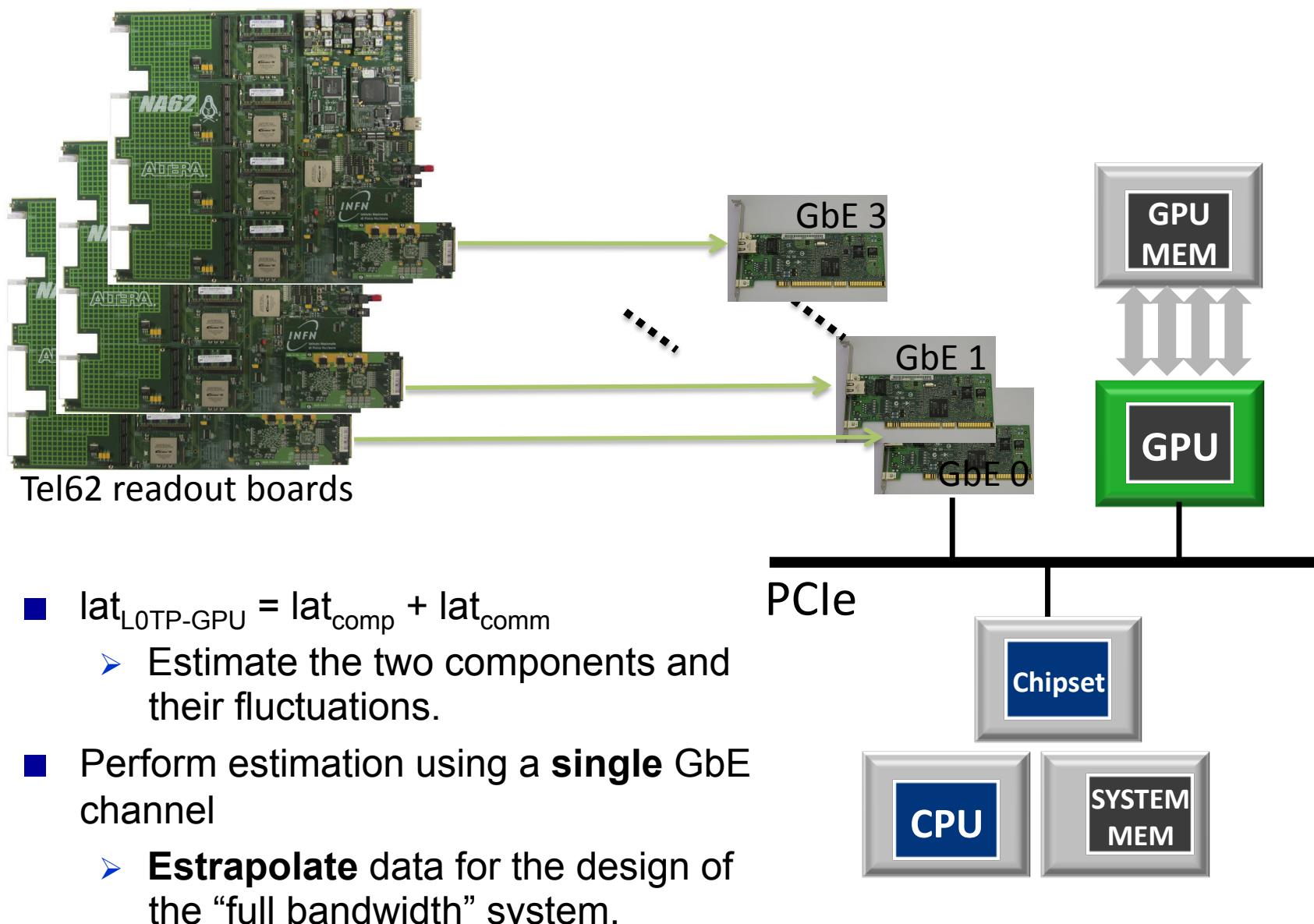
**TIME** →

Start GPS Clock	Off-shore CDR OK	Tx Offshore	On- shore CDR OK
-----------------	------------------	-------------	------------------

## DETERMINISTIC LATENCY ALIGNMENT STATE MACHINE



## GPU-Based RICH Level 0 TP System Latency Estimation



## NaNet: a FPGA-based NIC with GPUDirect RDMA capability

### ■ Challenge:

- Lower communication latency and its fluctuations.

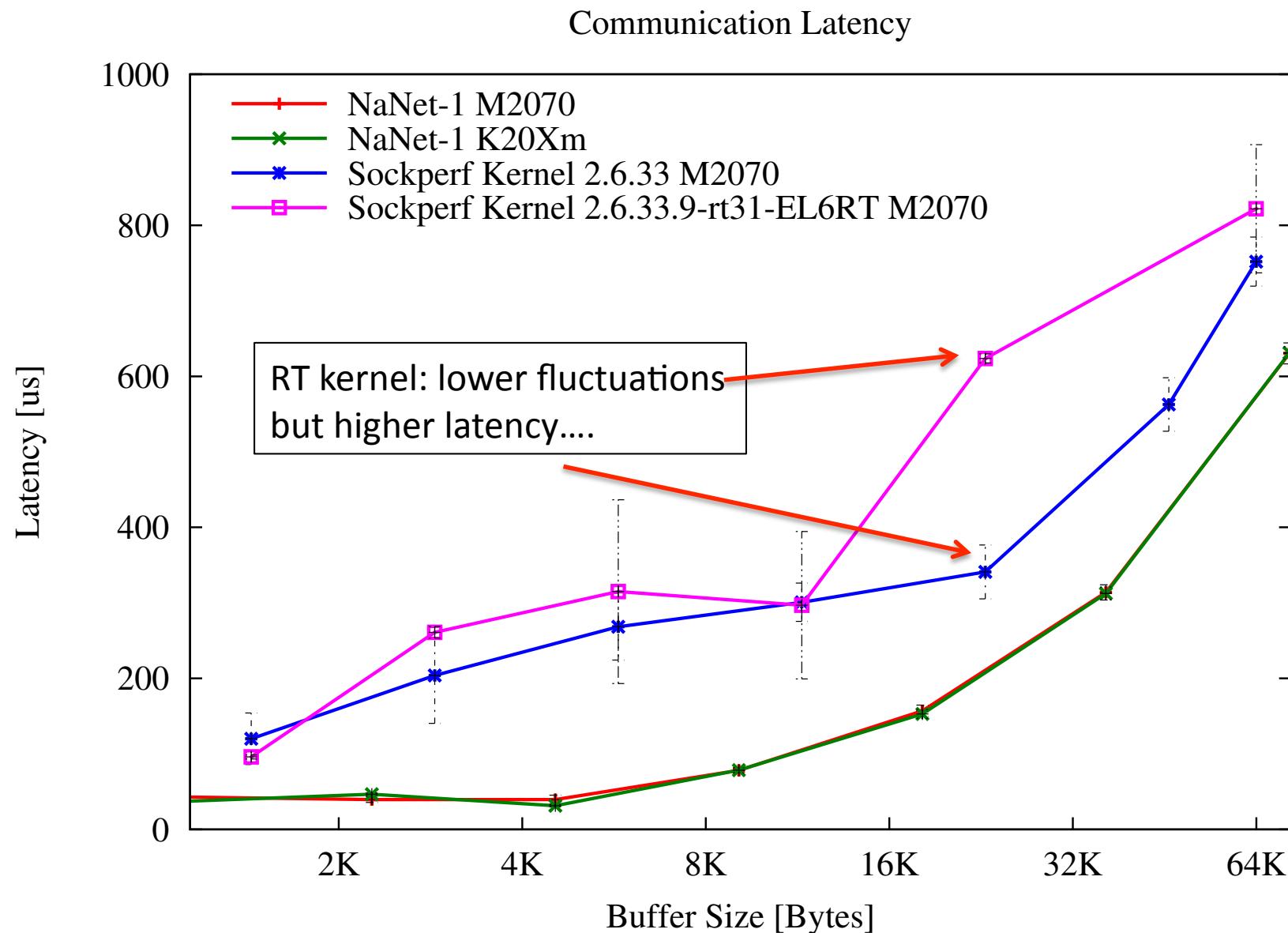
### ■ How?

1. Injecting directly data from the NIC into the GPU memory **without intermediate buffering**.
2. **Offloading** the CPU from network stack protocol management, avoiding OS jitter effects.

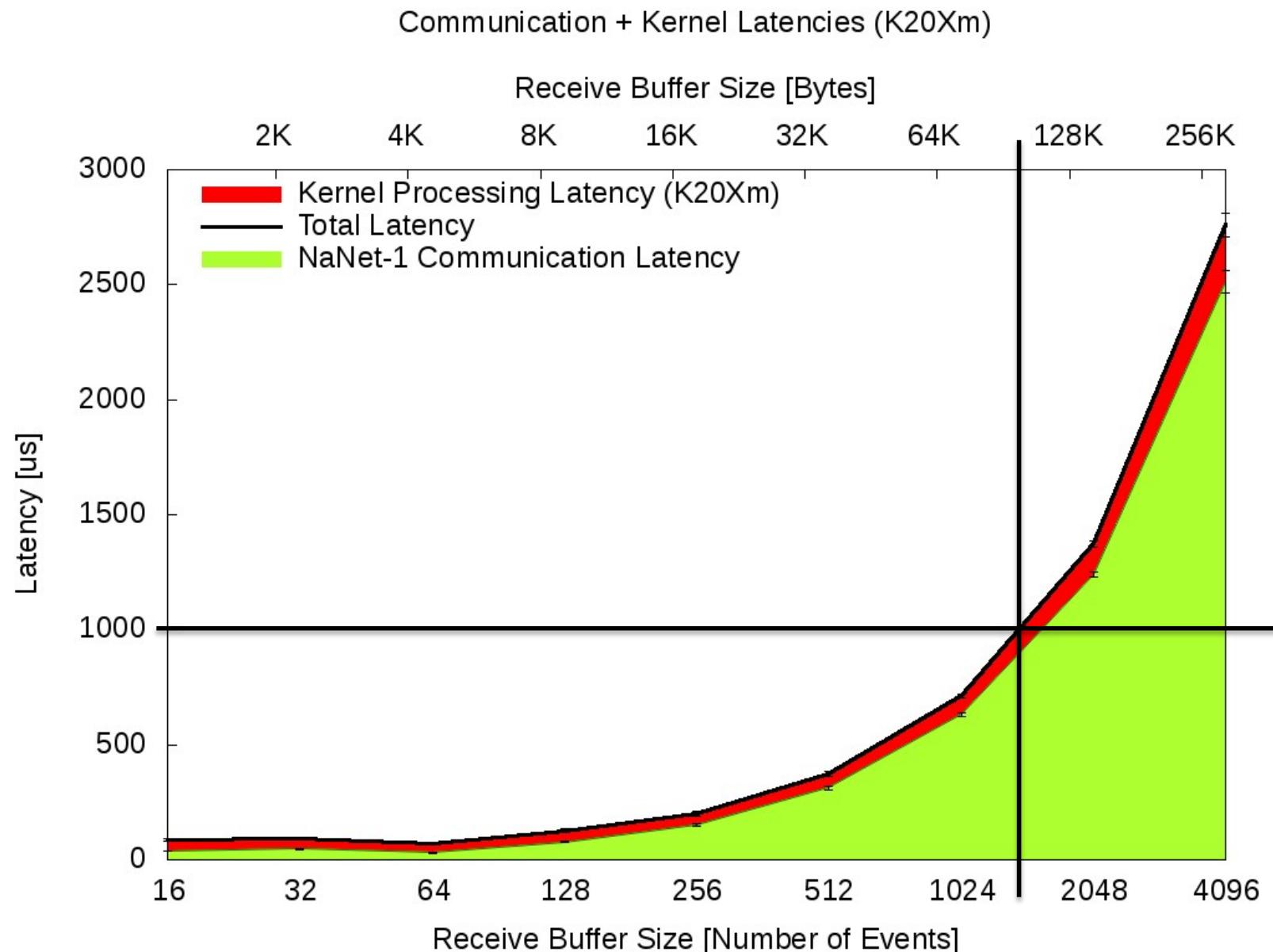
### ■ NaNet solution:

- Re-use the **APEnet+** design, implementing **PCIe** host interface and **GPUDirect** .
- Add a network stack protocol management offloading engine to the logic (**UDP Offloading Engine**).
- Use FPGA resources to perform **processing on data stream** (e.g. reformat data on-the-fly in a GPU-friendly fashion)

# NaNet-1 Communication Latency



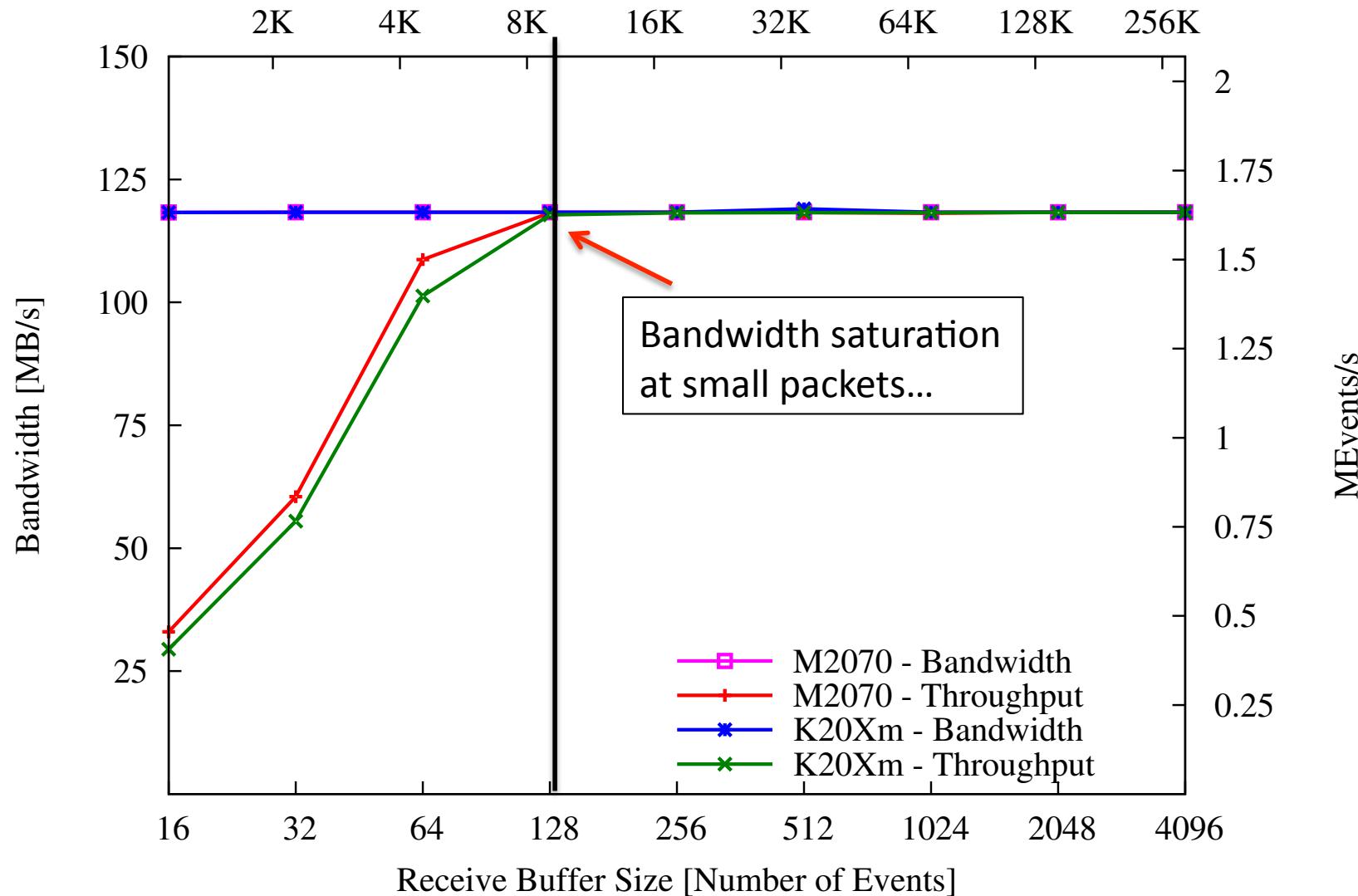
# Total latency of the GPU-Based RICH L0-TP using NaNet-1 (Kepler K20X)



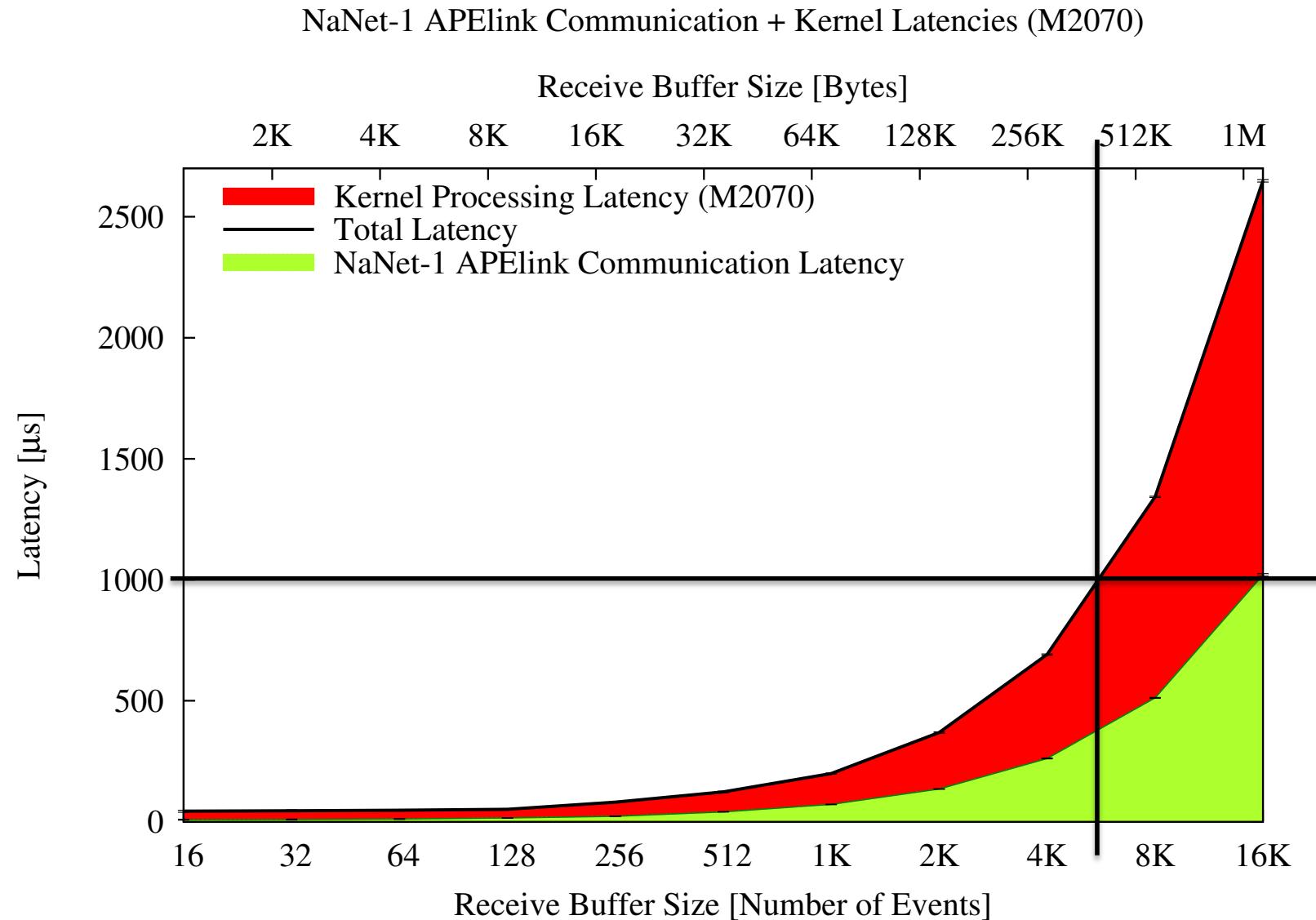
# NaNet-1 1 GbE Bandwidth & Throughput

NaNet-1 GbE Performance

Receive Buffer Size [Bytes]

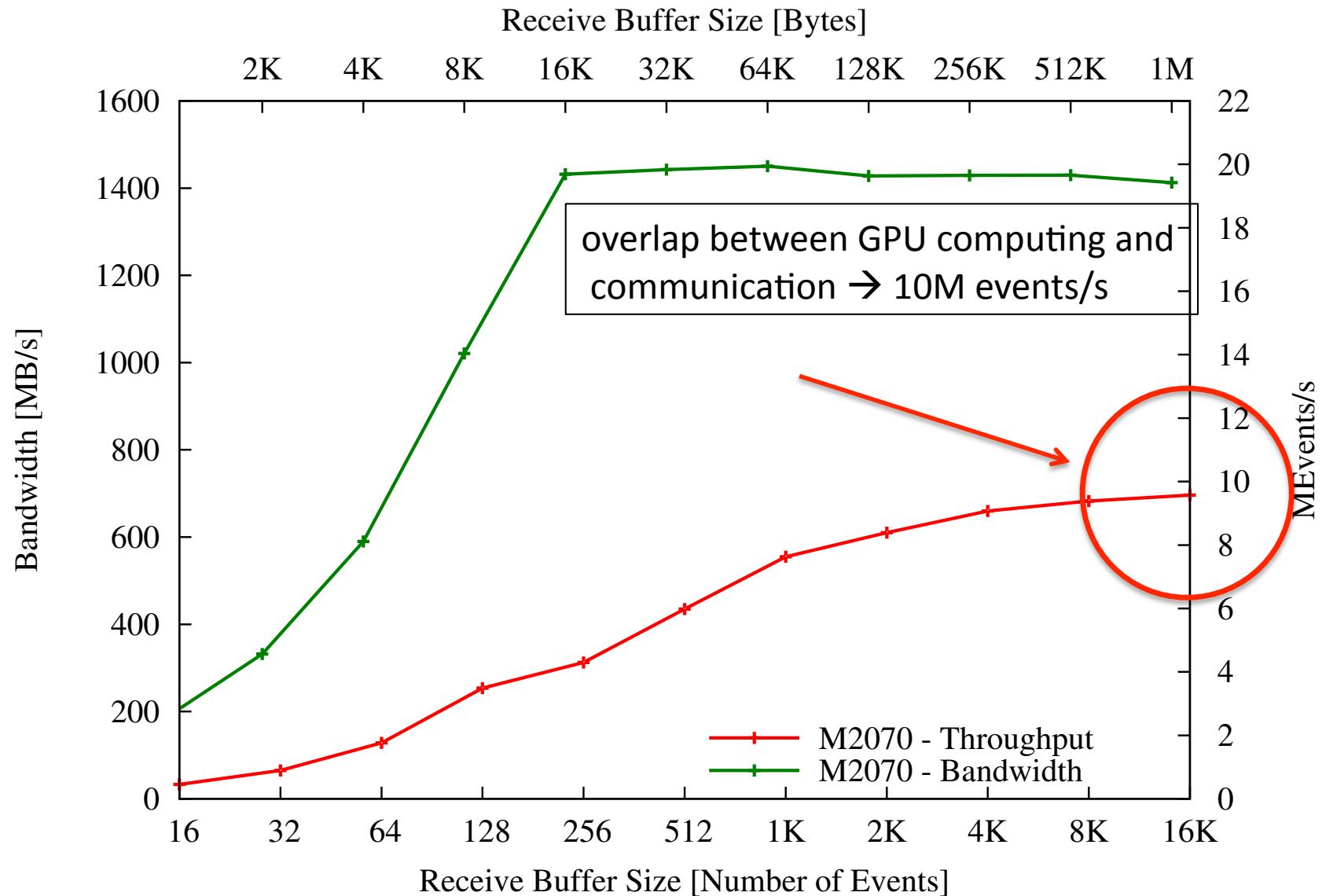


# Total latency of the GPU-Based RICH L0-TP using NaNet-1 (APElink)

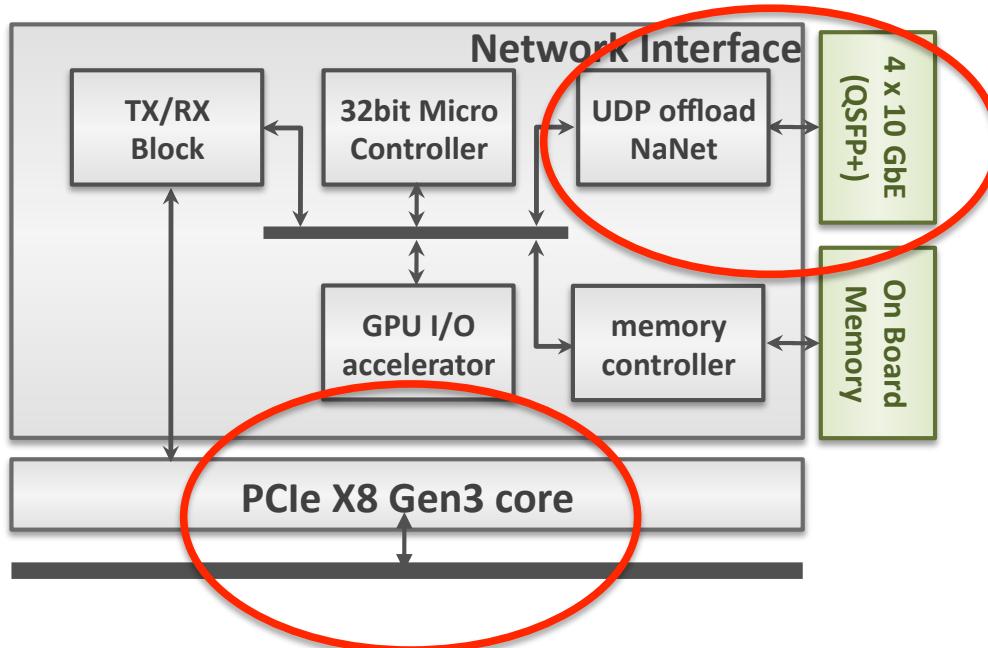


# NaNet-1 APElink Bandwidth & Throughput

## NaNet-1 APElink Performance



## NaNet-10 and last generation FPGAs (Stratix V)



- Implemented on the Altera Stratix V dev board
  - PCIe gen2 x8 but developing PCIe Gen3 (8 GB/s)
  - Faster embedded Altera transceivers (up to 14.1 Gbps)
  - hardened 10GBASE-R PCS features to support 10 Gbps Ethernet (10GbE)



QSFP+ to 4 SFP+ cable

