

HPC Network Architecture: Multi-FPGA Systems and Upcoming Projects

Bare metal application

KERNEL SPACE

0-4kB PACKETS

4 ports

Switch

64B/66B

PEROUTER

256bit@225MHz

(similar to IB verbs)

APE custom user library (LIBQCM)

BAR + DMA + REG

PCIe X16 GEN3 CORE

MPI application

APEnetX MPI BTL

APEnetX NI device driver

Xilinx device driver library (LIBQDMA) with custom SW added

H2C Stream

Engine

Routing logic

Arbiter

Descriptor

H2C DESC

2-4kB PACKETS

H2C module

APElink

Aurora

64B/66B

Direct device driver

I/O application

MASTER BRIDGE

IMMEDIATE

INFN APE Lab Team

USER SPACE

USER SPACE

DATA STRUCTURES

APEni

512bit@250MHz

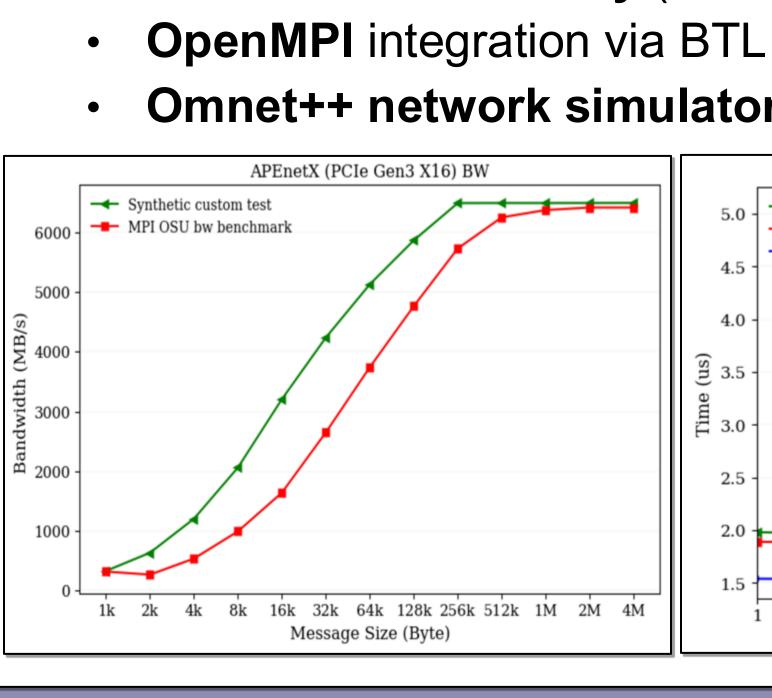
Network IP

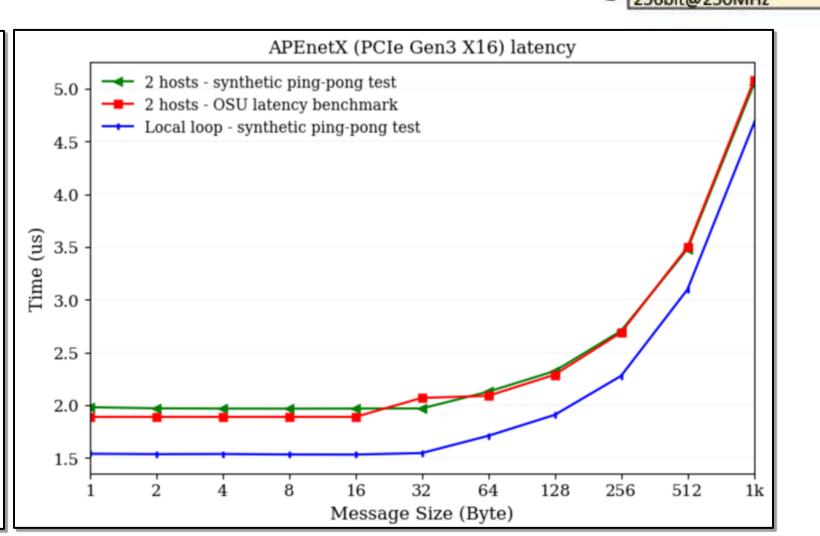
CUDA application

APEnetX: INFN Network Interface Card

- Direct network topology for low latency and highthroughput
- Use of FPGA (programmable component) for architecture exploration and fast&safe desigN
- PCIe x16 Gen3/Gen4 interface
- N-Dim **toroidal** mesh topology
- RDMA semantics + GPUdirect support
- Custom software stack:
 - Optimized Linux device driver
 - Low-level user library (similar to verbs)

 - Omnet++ network simulator module

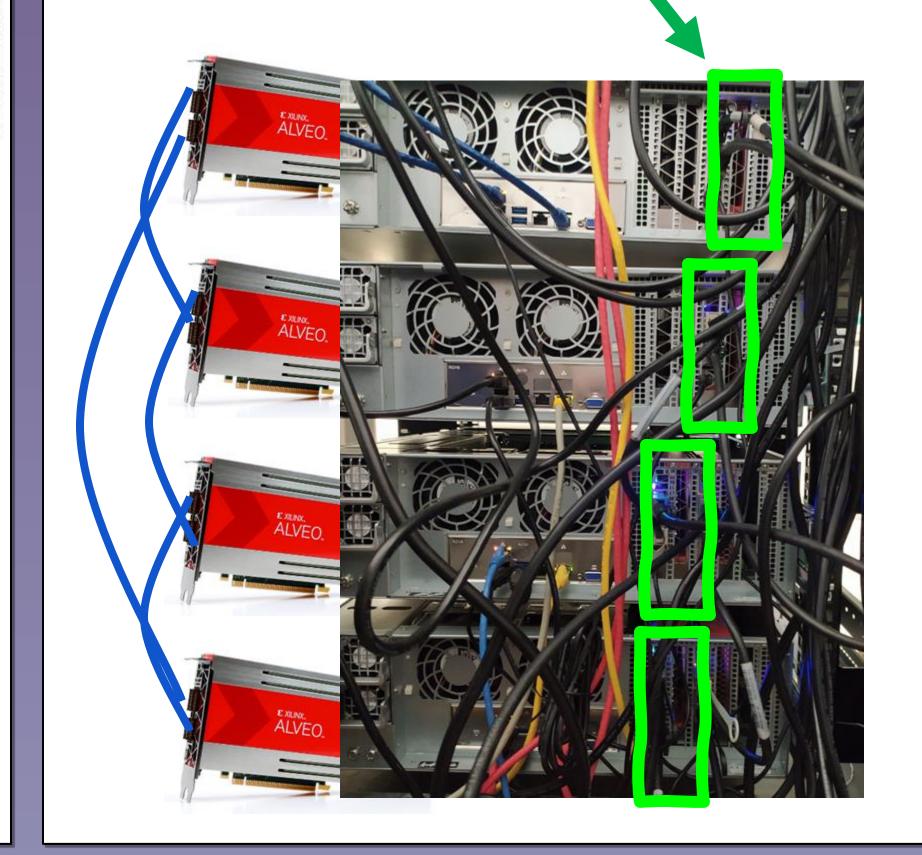




FPGA- based testbed

- **AMD Xilinx FPGA U200**
- 4x Supermicro servers
- Sapphire Rapids CPUs
- Gen3 PCIe

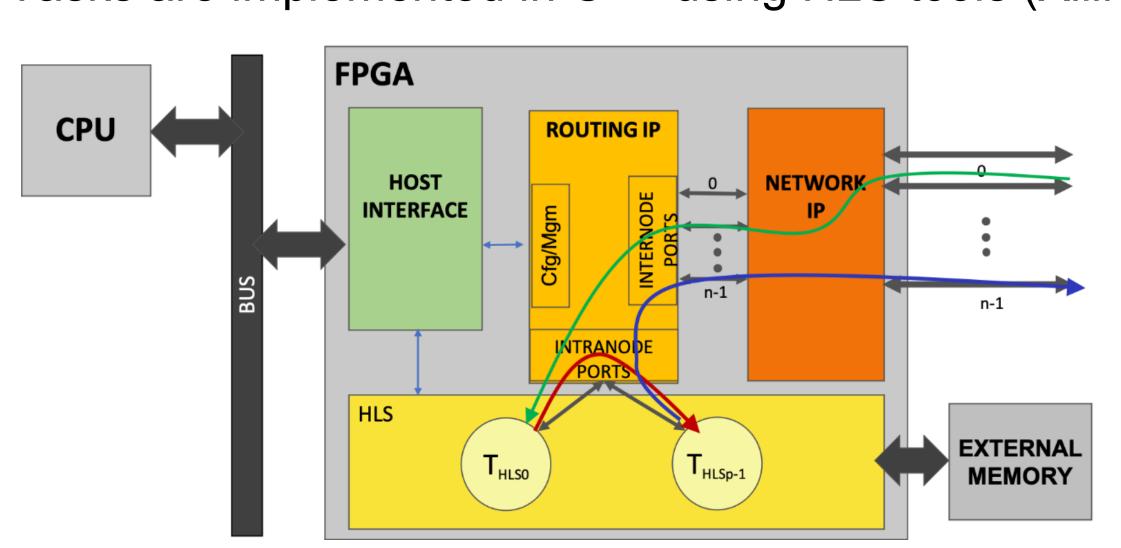
APEnetX prototypes connected point-to-point with QSFP28 cables

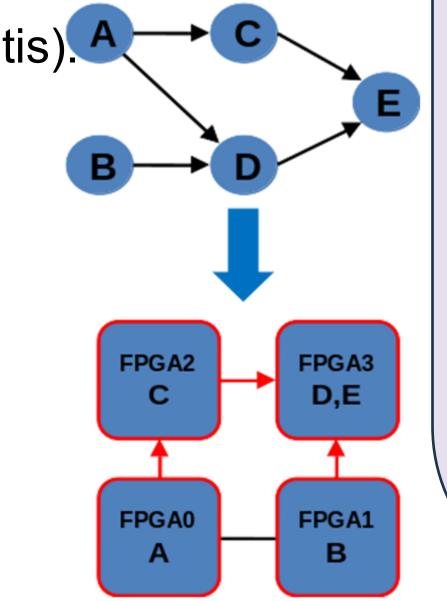


APEIRON: A Framework for High Level Programming of Dataflow Applications on Multi-FPGA Systems

- Enabling to mapping the dataflow graph of the application on the distributed FPGA system and offering runtime support for the execution.
- Allowing users, with no (or little) experience in hardware design tools, to develop their applications on such distributed FPGA-based platforms.
- INFN Communication IP

Tasks are implemented in C++ using HLS tools (Xilinx® Vitis).





APEIRON Physics use case

Development of intelligent Read-Out systems based on Neural networks

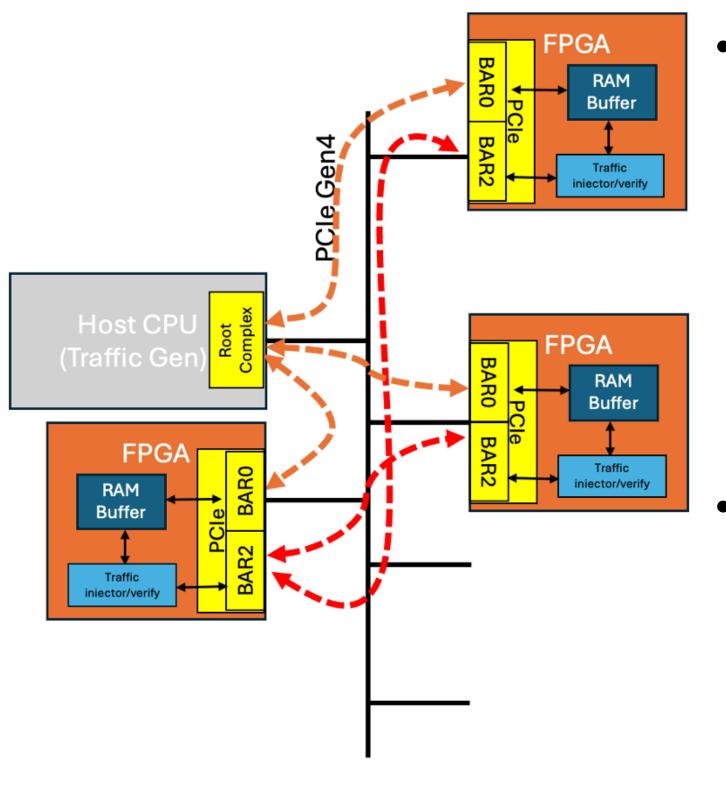
- Data streams from different data sources (detectors or sub-detectors) can be recombined through the processing layers using a low-latency, modular and scalable network infrastructure
- Online processing will be distributed on computing devices (FPGAs for the moment) in *n* subsequent layers.

Real-time inference on FPGAs: why?

- **Customizable I/O**
- **Deterministic latency**
- Useful in **heavy computation**

Project DARE (EuroHPC) **HPC Digital Autonomy with RISC-V in Europe**

The **DARE** consortium aims to establish a clear path for software and hardware development in Europe, leveraging early access to RISC-V hardware emulation and simulation, with the goal of deploying the developed technologies in European HPC systems



- INFN will contribute with Al-Direct **Engine** enabling the deployment of large scale NN models over multiple AIPU (Artificial Intelligence Processor Unit) accelerators to boost performance of applications like Al-accelerated HPC and Generative Al.
- **Al-Direct Specialized hardware** and its companion system software (linux device driver, user library) will be prototyped on a FPGAbased testbed

Project NET4EXA (EuroHPC) **NETwork for EXascale Architectures**

NET4EXA aims to develop a next-generation high-speed interconnect for HPC and AI systems, building on the success of the BXI European HOC Interconnect and the advancements made through research in the RED-SEA project and other previous European RIA initiatives.

INFN contribution

- Innovative mechanisms for congestion control management
- ON-NIC processing for task streaming computing,
- Prototyping new features supporting GPU triggered computing
- INFN scientific applications for benchmarking network

