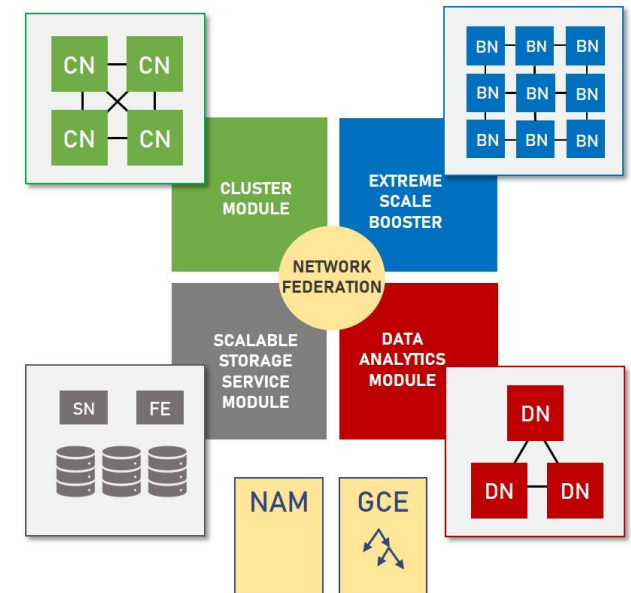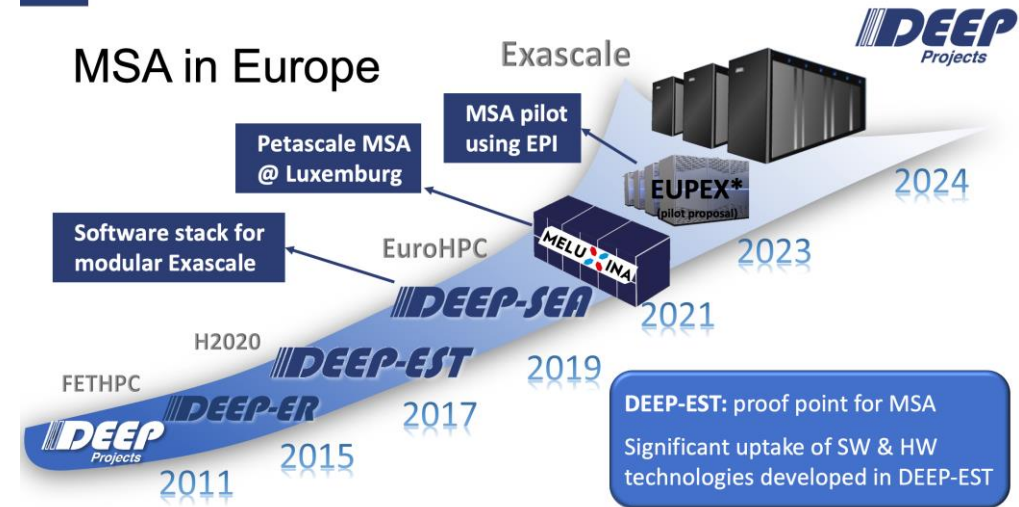**DEEP**
Projects

# Towards a Modular Supercomputer Architecture: The DEEP-EST Project

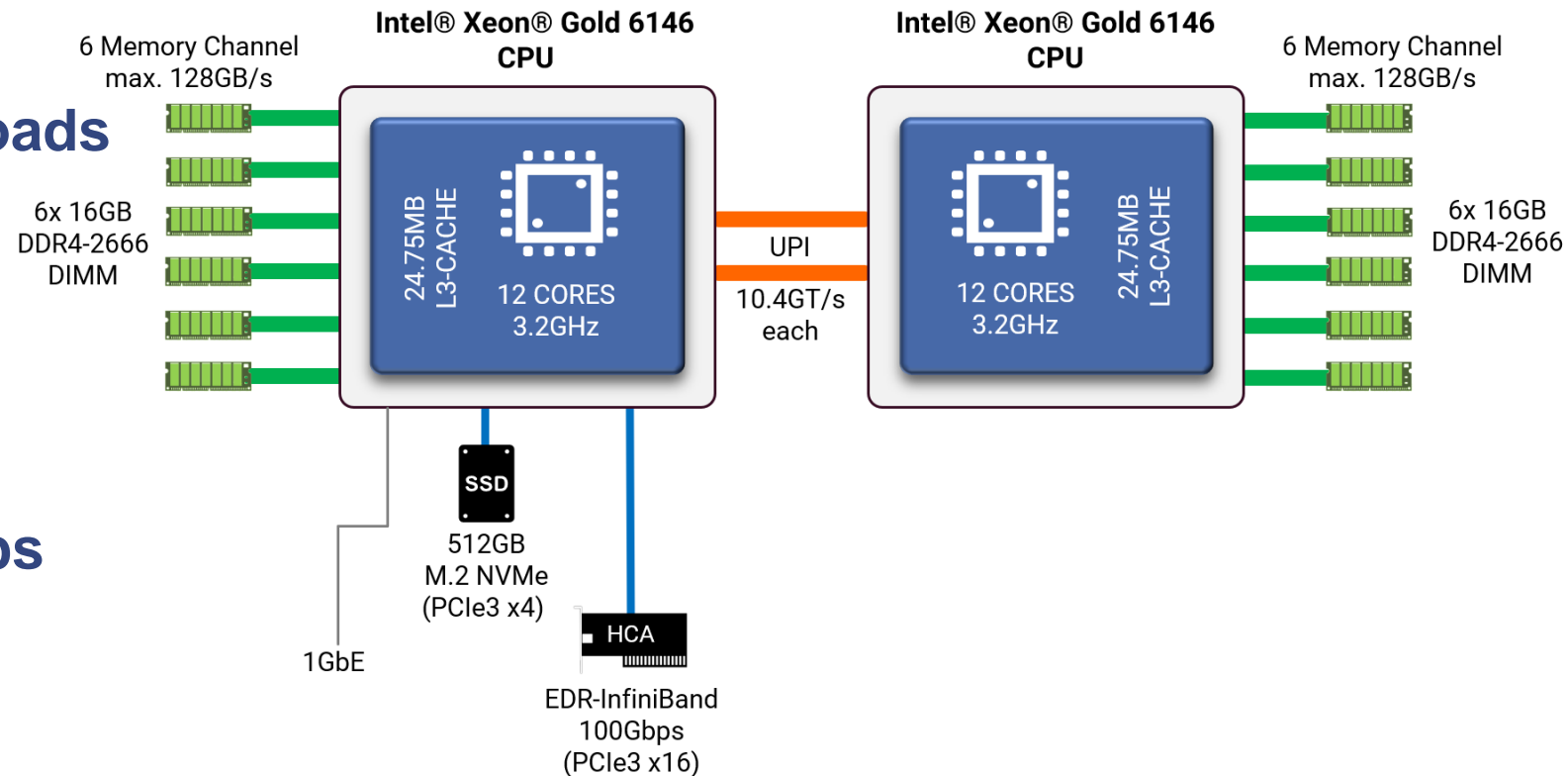*Viktor Khristenko (CERN), Maria Girone (CERN)*

# DEEP-EST Modular Supercomputer

- Prototype for the Modular Heterogeneous HPC system

- Convergence of HPC and HPDA worlds

- Variety of hardware to enable wide range of applications

- Software Hardware co-design driven by 6 applications (reconstruction in CMS included)
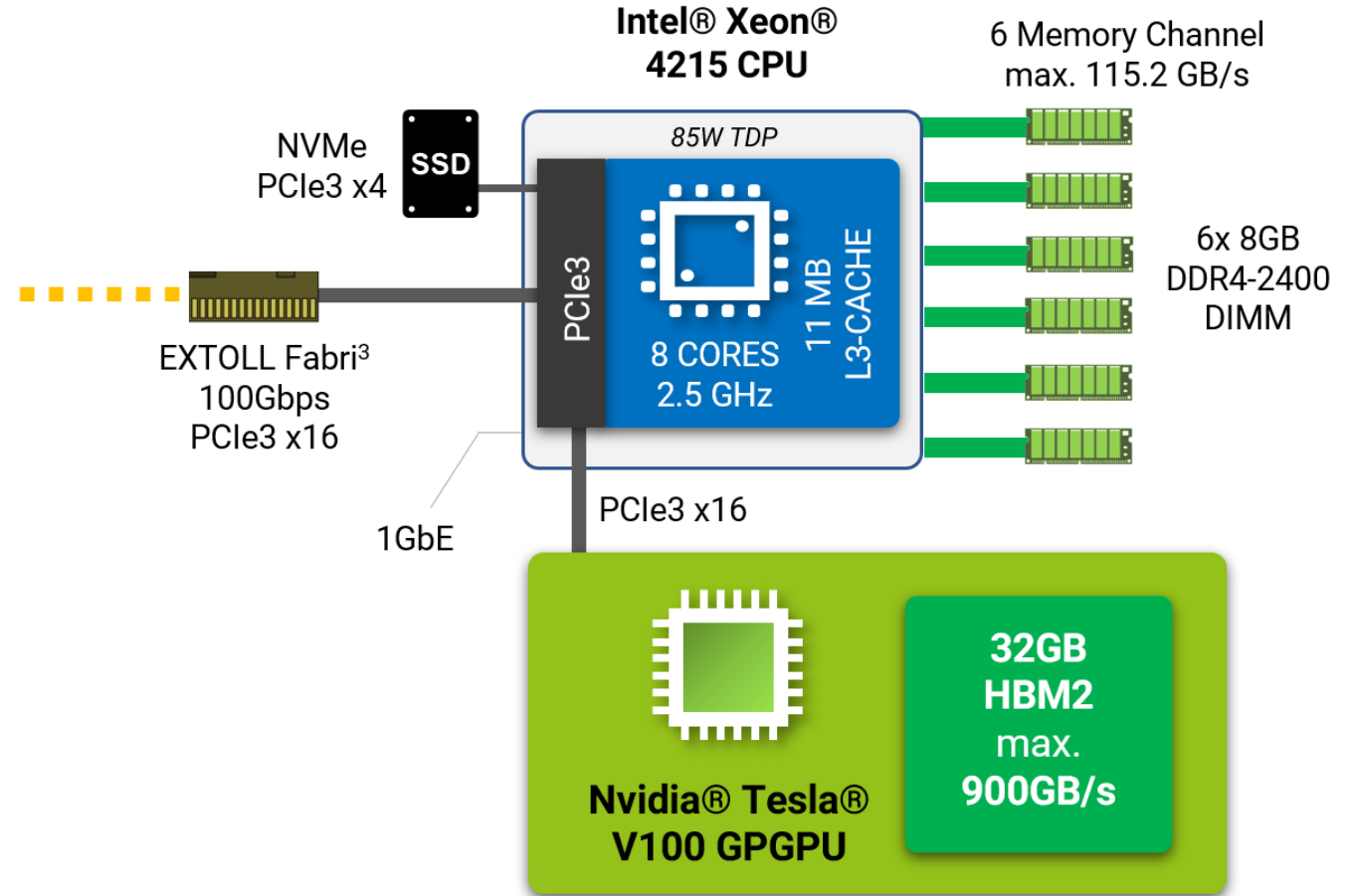
# Cluster Module

- **Overall 50 nodes**

- **Aimed at CPU-bound workloads**

- **To/from ESB**
  - **Infiniband/Extoll Bridge**

- **To/from DAM**
  - **Inifiband/Ethernet 40Gbps Bridge**

6 Memory Channel
max. 128GB/s

6x 16GB
DDR4-2666
DIMM

**Intel® Xeon® Gold 6146 CPU**

24.75MB L3-CACHE

12 CORES 3.2GHz

UPI

10.4GT/s each

**Intel® Xeon® Gold 6146 CPU**

12 CORES 3.2GHz

24.75MB L3-CACHE

6 Memory Channel
max. 128GB/s

6x 16GB
DDR4-2666
DIMM

SSD

512GB
M.2 NVMe
(PCIe3 x4)
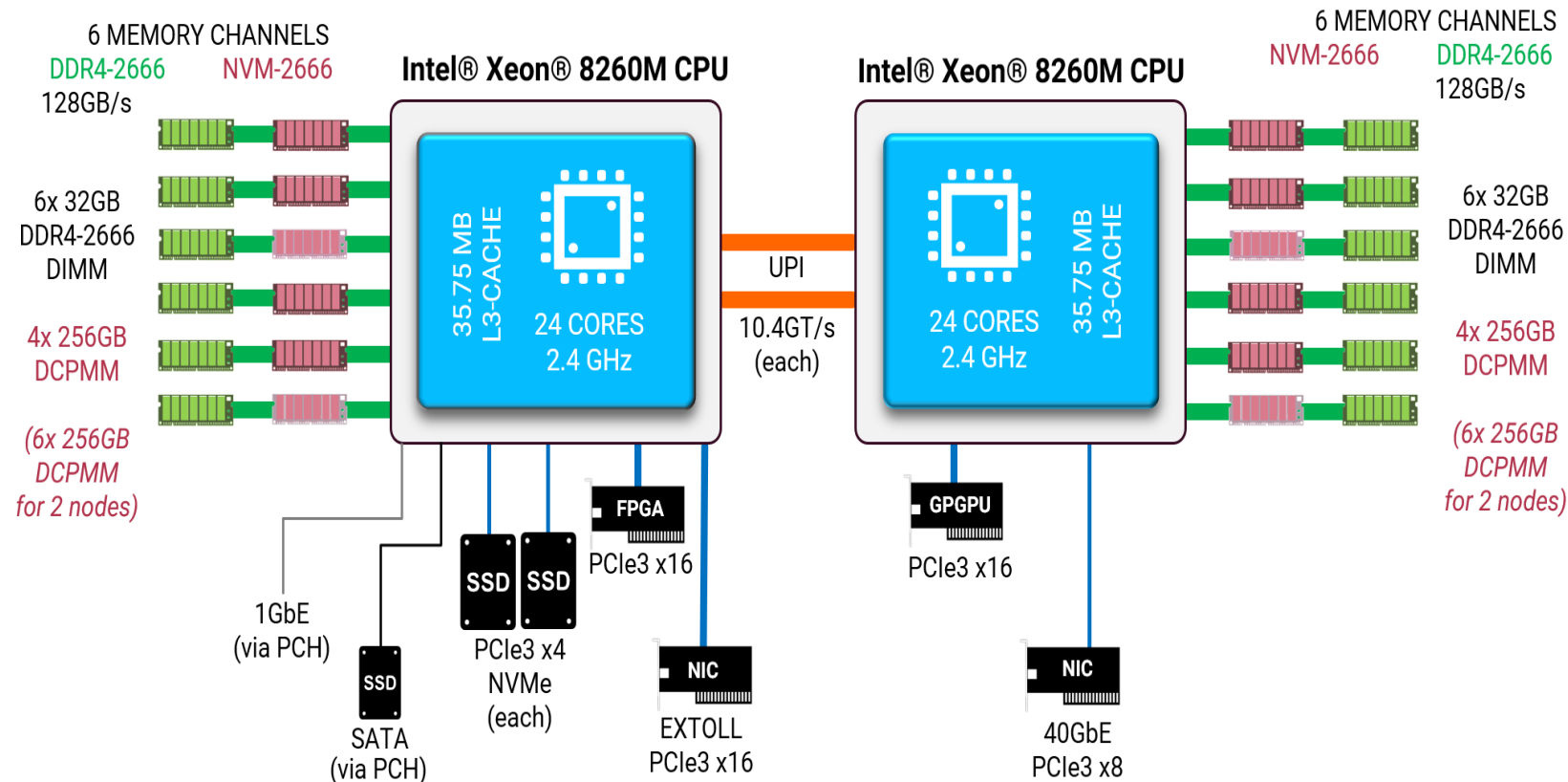
HCA

1GbE

EDR-InfiniBand
100Gbps
(PCIe3 x16)

# Extreme Scale Booster

- **Overall 75 nodes**

- **GPU-based, Nvidia V100**

- **Extoll / Infiniband Network Fabric**
  - **~ Quarter of nodes uses Extoll**



Intel® Xeon® 4215 CPU

85W TDP

8 CORES 2.5 GHz

11 MB L3-CACHE

PCIe3

NVMe PCIe3 x4

SSD

EXTOLL Fabri[3] 100Gbps PCIe3 x16

1GbE

PCIe3 x16

6 Memory Channel max. 115.2 GB/s

6x 8GB DDR4-2400 DIMM

Nvidia® Tesla® V100 GPGPU
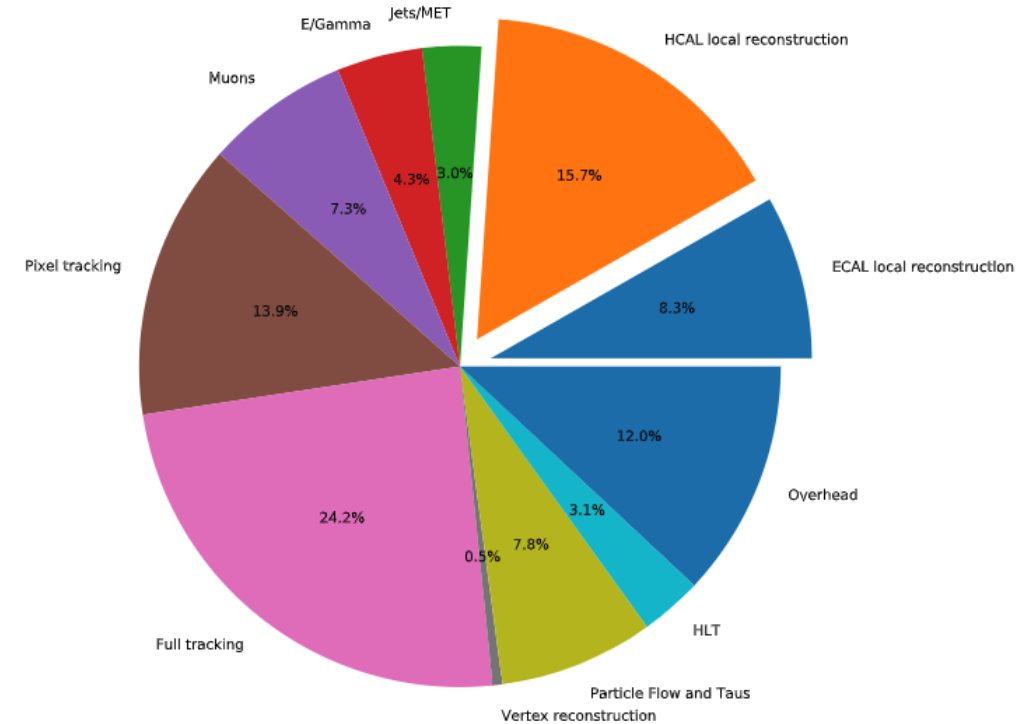
32GB HBM2 max. 900GB/s

# Data Analytics Module

- 16 nodes

- 2 accelerators per node
  - 1 GPU - Nvidia V100
  - 1 FPGA - Intel Stratix 10

- Memory
  - 2-3TBs Intel Optane Memory
  - 384GB DDR4
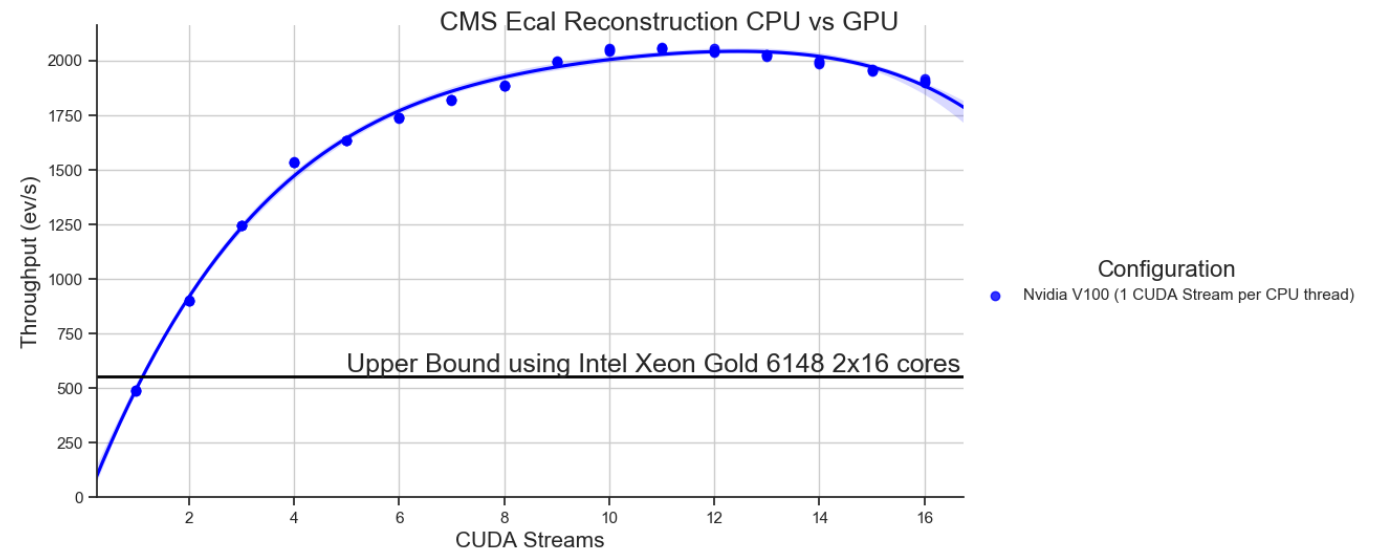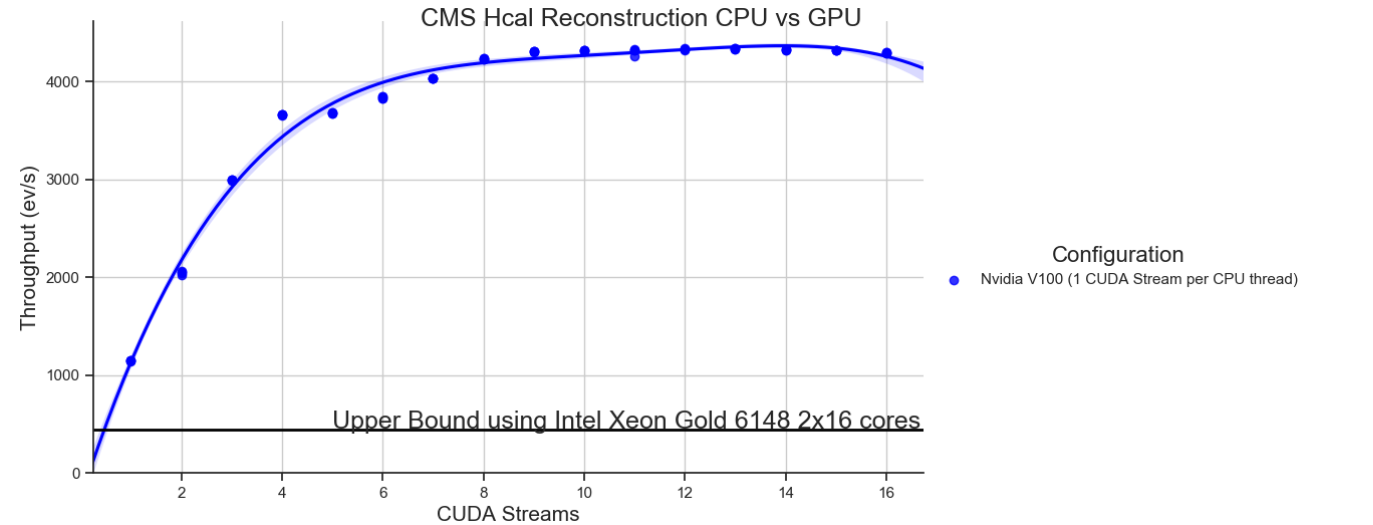
# DEEP-EST: Heterogenous data processing

- Heterogenous Execution for CMSSW within Patatrack Effort

- Porting of HCAL / ECAL Local Energy Reconstruction

- Current Calorimeters consume ~15-20% of the total HLT time

    – Both (Ecal/Hcal) utilize the same algorithm (the core part) for energy regression

# Results: CMS Hcal/Ecal only

- http://opendata.cern.ch/record/12303

- 20K events. Replicate twice

- Hcal -> speed of 7-8x
  – Using Nvidia V100 GPU
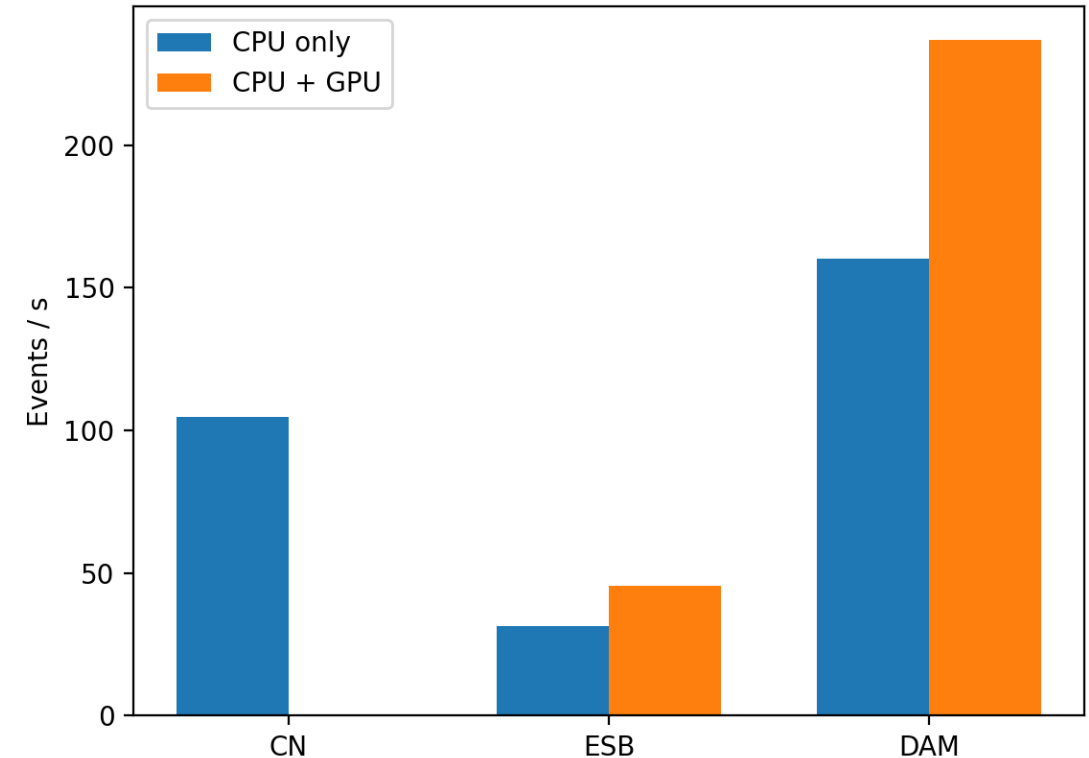
- Ecal -> speed of 3-4x
  – Using Nvidia V100 GPU

Intel Xeon Gold 6148 used for the baseline



CMS Hcal Reconstruction CPU vs GPU

Configuration
Nvidia V100 (1 CUDA Stream per CPU thread)



CMS Ecal Reconstruction CPU vs GPU

Configuration
Nvidia V100 (1 CUDA Stream per CPU thread)

# Results: CMS HLT-like Run3
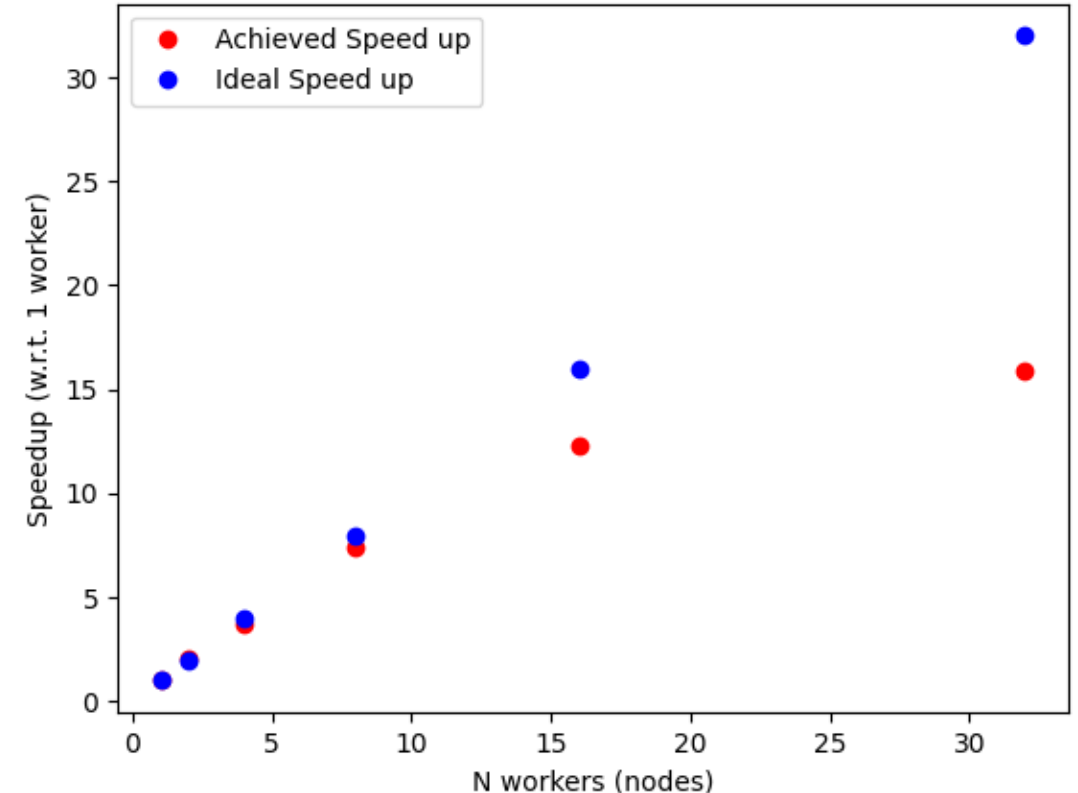
- http://opendata.cern.ch/record/12303 - Open Data used

- CMS HLT-like Run3 configuration
  - Includes Patatrack GPU developments

- **50% more out of nodes with Nvidia GPUs (V100 here)**

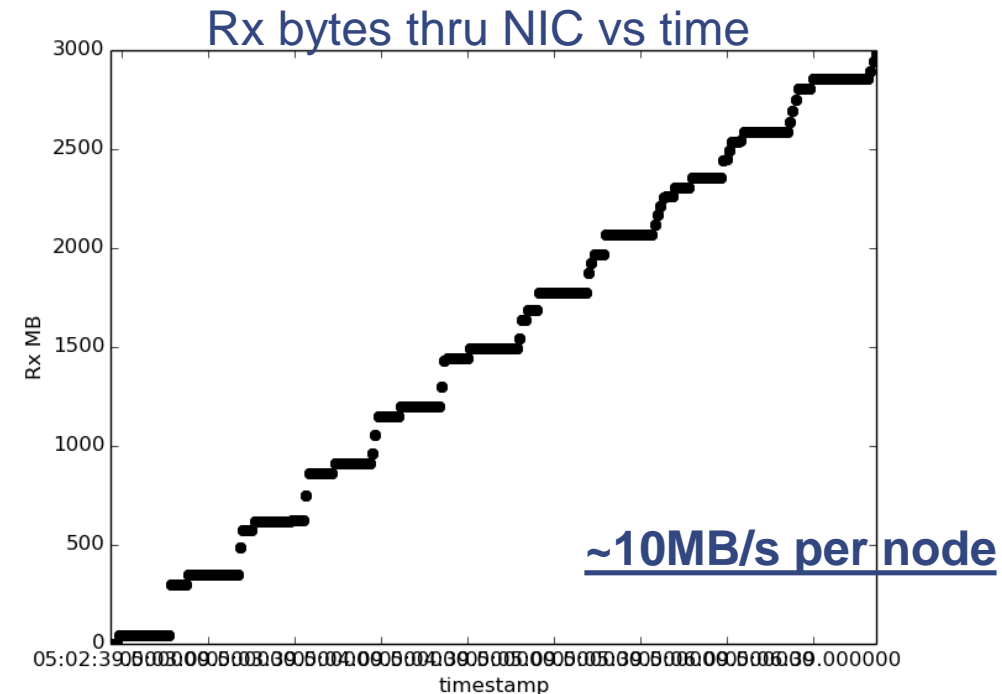Throughput by node type. CMS HLT Run3 configuration with Open Data

# Results: Testing AI Workflows

- Goal: Understand scalability of distributed DL training on HPC
  - specifically for HEP workloads

- Used
  - JEDI-net (arxiv: 1908.05318) model
  - NNLO package for distributed MPI
    - *ParastationMPI* on DEEP-EST prototype
    - *With GPUDirect support*
  - Data from zenodo

- Using up to 32 DEEP-EST ESB nodes
  - 1 Nvidia V100 GPU per node
  - 100Gbps Infiniband interconnect

- **Good scalability observed**

# Data Access / HPC:

- How will HEP process Exabytes of data?!

- **Would HPC Storage Systems be enough**
  - Or Exascale HEP workflows will require extra

- CMS MINIAOD2NANO workflow
  - Running up to 64 nodes
    - *~640MB/s aggregate from Ceph (SDSC Popeye)*
    - *This is on average (HEP i/o is bursty)*

- Hypothetically, Exascale HPC O(1M) cores total -> O(10K) nodes -> O(100GB/s) aggregate
  - Although we might never use from a single site

### Rx bytes thru NIC vs time

**~10MB/s per node**

# Conclusions

- The DEEP-EST Project concludes end of March 2021

- The DEEP-EST Project proved to be an invaluable platform for
  - Collaboration with HPC experts from other sciences/centers
  - HEP Tests and Developments towards the usage of Exascale HPCs

- Contribution to the Patatrack effort has been integrated into CMS Experiment's framework and will be used in Run 3.

- Now starting to work on RAISE HPC project

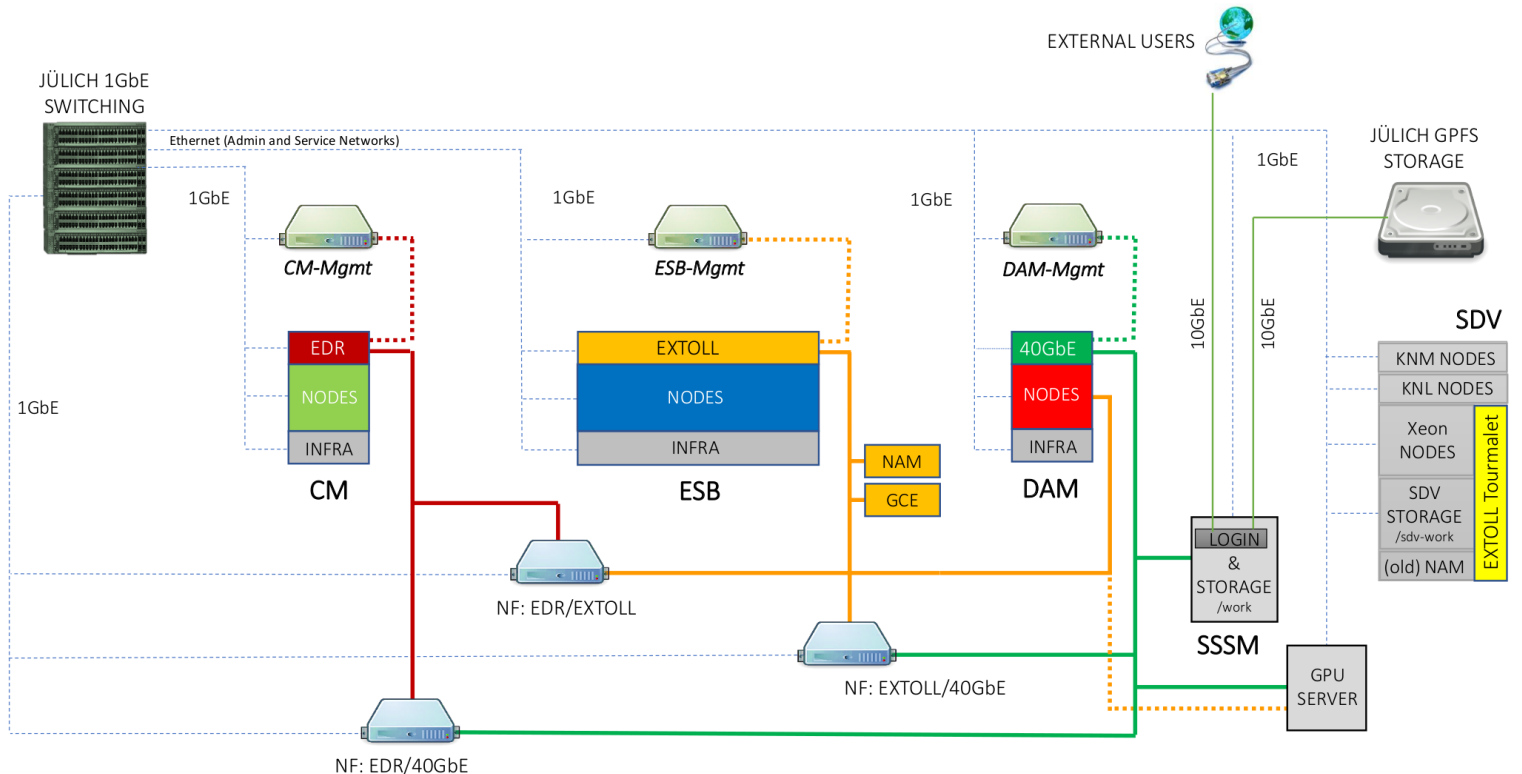**New CoE RAISE: Research on AI- and Simulation-Based Engineering at Exascale**

# Network Federation + Auxiliary

- **Multiple fabrics**
  - **100Gbps Infiniband**
  - **100Gbps Extoll**
  - **40 Gbps Ethernet**
  - **Bridges**

- **Network Attached Memory NAM**
  - **Extoll's FPGA based solution**
  - **128GBs DDR4**
  - **TB(s) SSDs**

- **Global Collective Engine GCE**
  - **Extoll's FPGA based solution**
  - **Accelerate MPI-collective operations**

**DEEP-EST Prototype – Schematic Network Overview**

# Results: Testing Intel OneAPI

- **Ported standalone CMS Ecal Reconstruction algorithm to use Intel OneAPI**

- Results between regular c++ and Intel OneAPI implementations match 100%
  - Tested on a CERN's VM with CPUs only

- Employed Intel OneAPI Compatibility Tool to convert from CUDA-based implementation
  - Almost 0 modifications after the conversion
  - Had to adapt slightly Eigen to be used within OneAPI kernels

- Github repo: https://github.com/vkhristenko/cmsregr-oneapi