www.evolve-h2020.eu

# EBDVF 2021

**The Evolve Project : A Convergence of Machine Learning and HPC to address Big Data Challenges**

**Evolve Partners**

**EBDVF 2021, November 29, 2021**

# Agenda

- ❑ Introduction to Evolve (DDN/Jean-Thomas Acquaviva, 15')
- ❑ Hardware Platform (ATOS/Huy-Nam Nguyen, 15')
- ❑ Software Development (FORTH/Angelos Bilas, 15')
- ❑ Applications (NEUROCOM/Vassilis Spitadakis, 15')
- ❑ Proof-of-Concept (VIF/Alexander Stocker, 15')
- ❑ Q&A (Evolve, 15')

# EBDVF 2021
## The Evolve Hardware Platform

**Evolve Partners**

**EBDVF 2021, November 29, 2021**

# Objectives

❑ **Motivations and Objectives**

- o Evolution of System Architecture for integration of Hw Accelerators
  - Flexible, scalable and interoperable exploitation of Hw accelerators
  - Acceleration of Computing and Data Transfer
  - Enabling HPC-BD/AI Convergence
- o Validation on a wide spectrum of applications
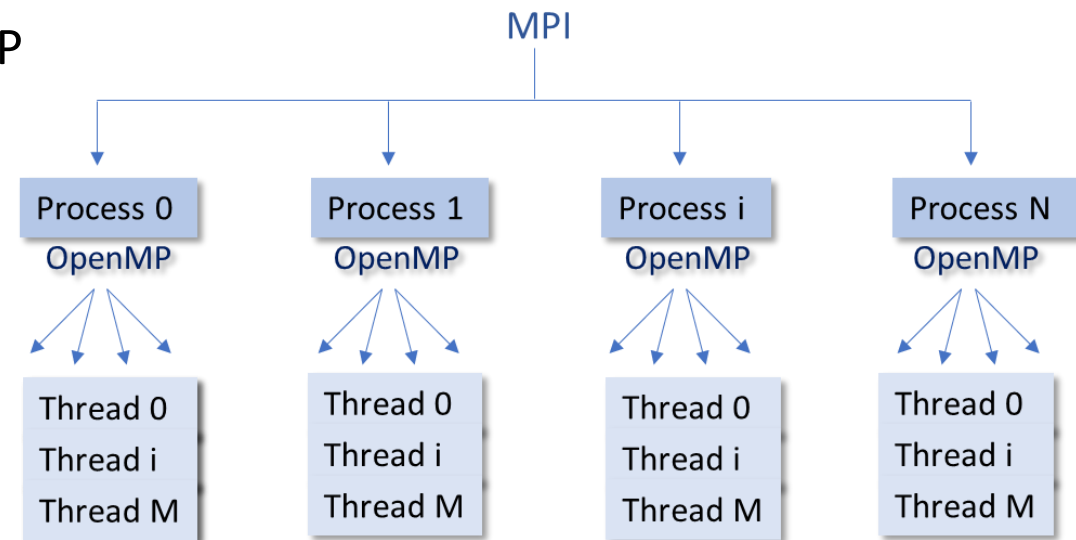
❑**Main Technical Challenges**

- o Scalable approach to heterogeneous computing technologies
- o Combination/Convergence of acceleration software stacks
- o Convergence of HPC, BigData and AI
- o Expansion of HPC's scope to the Cloud/Edge Computing
- o HPC Features: Availability/Reliability/Accessibility/Security

# Computing Technologies

❑ **CPUs** (Intel/Broadwell-Haswell-SKX)

- o Homogeneous Numa Multi-core
- o Performance = #AddMult*#FVect*#Cores*Frequency
- o Frequency, #Cores, Cache size, Integrated GPU
- o Cache Coherence / Memory Consistency Protocol
- o Compute Node size, Node Controller, Numa Factor

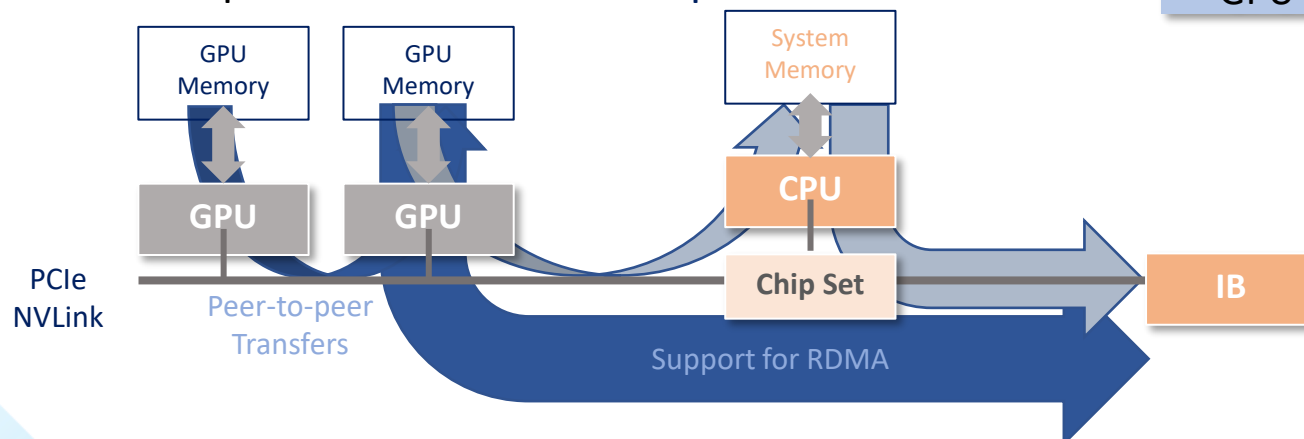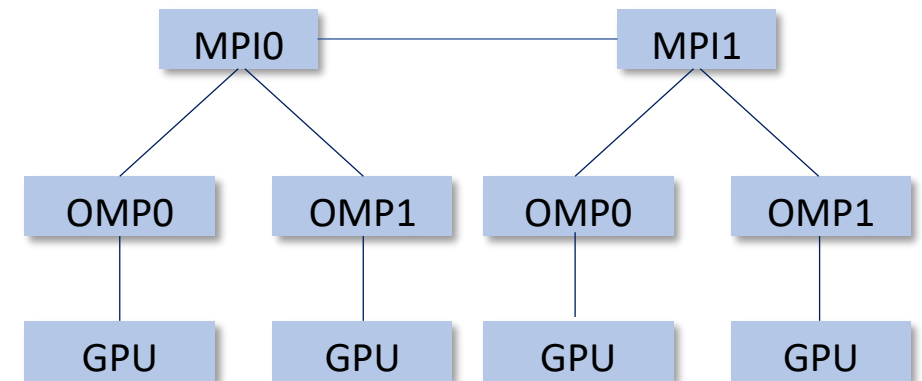❑ **Programming Models** (Hybrid OpenMP+MPI)

- o Combinations from pure MPI to MPI+OpenMP
- o Adequation with system architecture
- o Memory management, Load Balancing
- o Performance vs Memory occupation
- o Criterions : NUMA factor, Message size,
  Synchronization (e.g. locks, race)
- o MPI-3 shared memory programming
- o Accelerator support in OpenMP4.0

MPI

| Process 0 | Process 1 | Process i | Process N |
|-----------|-----------|-----------|-----------|
| OpenMP | OpenMP | OpenMP | OpenMP |

| Thread 0 | Thread 0 | Thread 0 | Thread 0 |
|----------|----------|----------|----------|
| Thread i | Thread i | Thread i | Thread i |
| Thread M | Thread M | Thread M | Thread M |

# Computing Technologies (cont'd)

❑ **GPUs** (Intel embedded GPU, Nvidia Tesla K20, P40, V100)

- o  Performance = #AddMult*#SM*#C$_{SM}$*Frequence
- o  #Streaming Multiprocessor/#cores, Thread parallelism, Frequency
- o  Memory Type, Size, Bandwidth and repartition
- o  Programming Models : TF/MatLab, OpenMP/OpenAcc, CUDA/OpenCL
- o  Interaction with CPUs
  - Hybrid MPI + OpenMP/CUDA-OpenCL
  - GPU scalability
    - o  Mapping MPI tasks to GPUs
    - o  Scaling code on multiple GPUs
    - o  Monitoring GPUs in a heterogeneous cluster
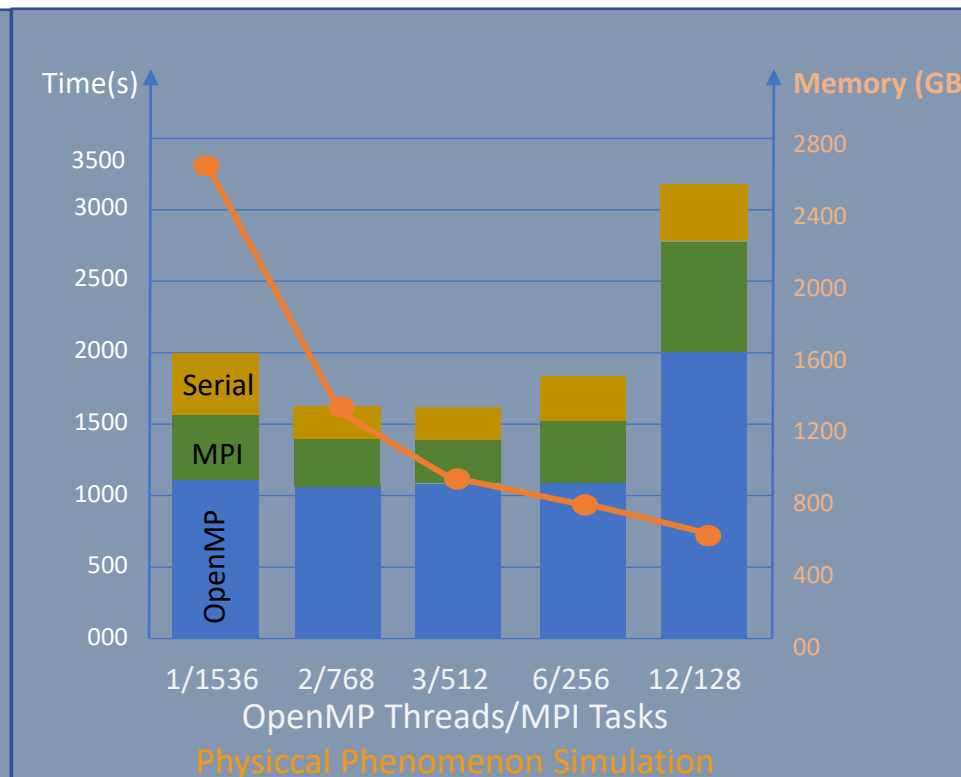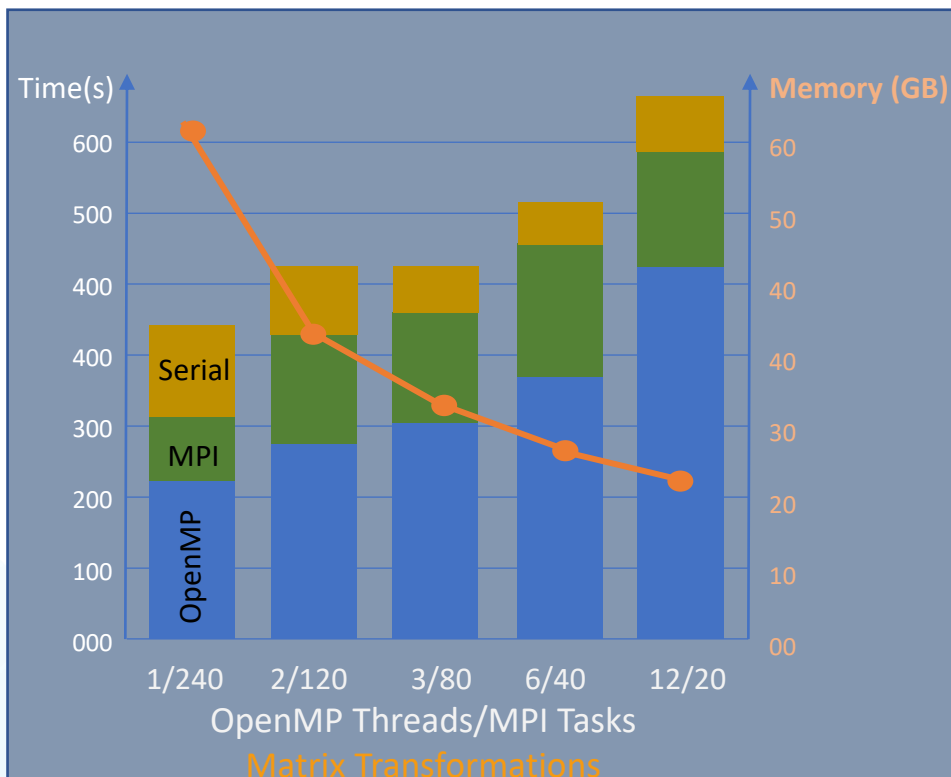  - Multiple CUDA-aware MPI implementation

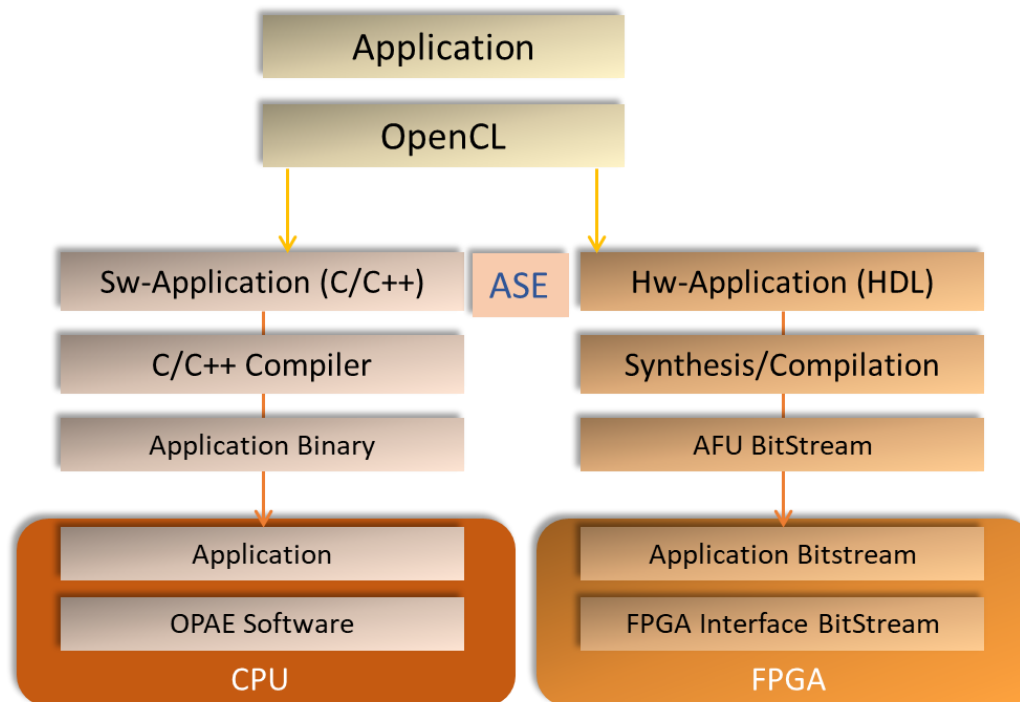# Computing Technologies (cont'd)

❑ Hybrid Programming Models

- ○ Combination of OpenMP + MPI
- ○ CPU vs OpenAcc vs CUDA
- ○ CUDA-aware MPI

| FFT Time (s) | CPU | GPU OpenACC | GPU CUDA |
|---|---|---|---|
| Exec | 53.426 | 3.767 | 1.08 |
| Transfer | 0 | 3.342 | 4.719 |
| Total | 53.426 | 7.109 | 5.799 |
| SpeedUp | 1 | 7.5 | 9.2 |



Matrix Transformations



Physiccal Phenomenon Simulation
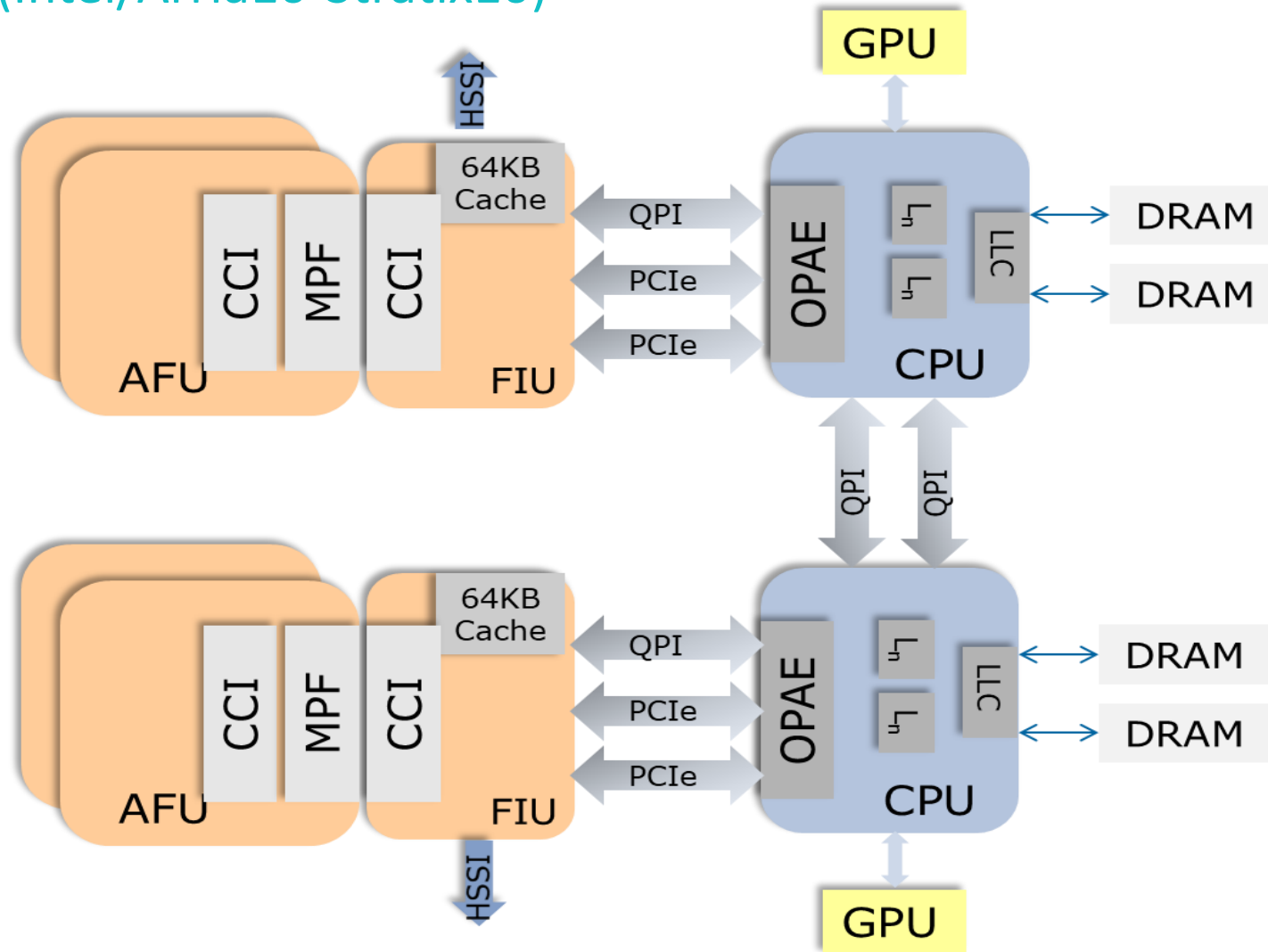
# Computing Technologies (cont'd)

❑ **FPGA** (Intel/Arria10-Stratix10)
- o #CLBs/LEs/LUTs, RAMs, DSPs, Interconnect, SERDES/Transceivers
- o FPG-based Programming Models : HDL (Vhdl/SVerilog), OpenCL
- o Discrete vs Integrated implementation
- o Challenge/Evolution trend : Clocking, Multi-FPGA

# Computing Technologies (cont'd)

❑ **FPGA** (Intel/Arria10-Stratix10)

# Computing Technologies (cont'd)

❑ **Acceleration of Computation**

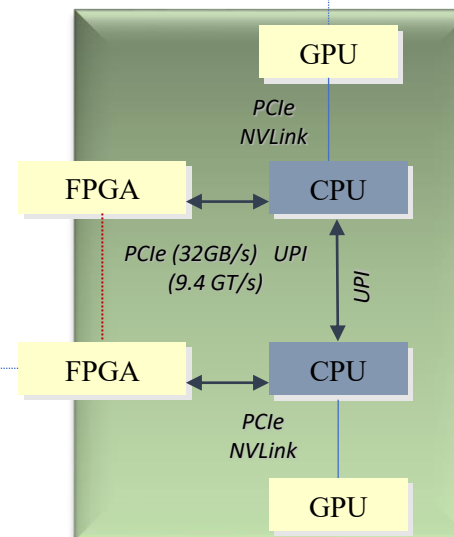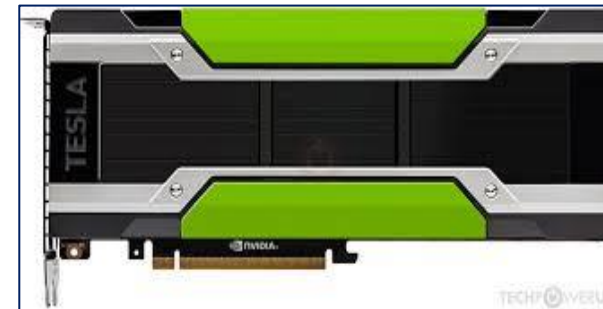Nvidia/P40 : 3840 Cores, 24 GB GDDR5, PCIe 3.0
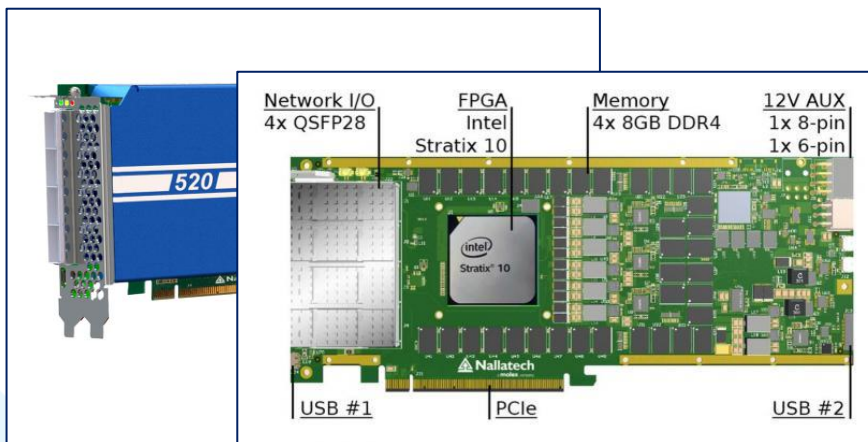
Nvidia/V100 : 5120 Cores, 32 GB HBM2, Pcie Gen3

Nallatech/520N

    Intel/Stratix 10 GX 2800

    4x QSFP28s for 400Gbps

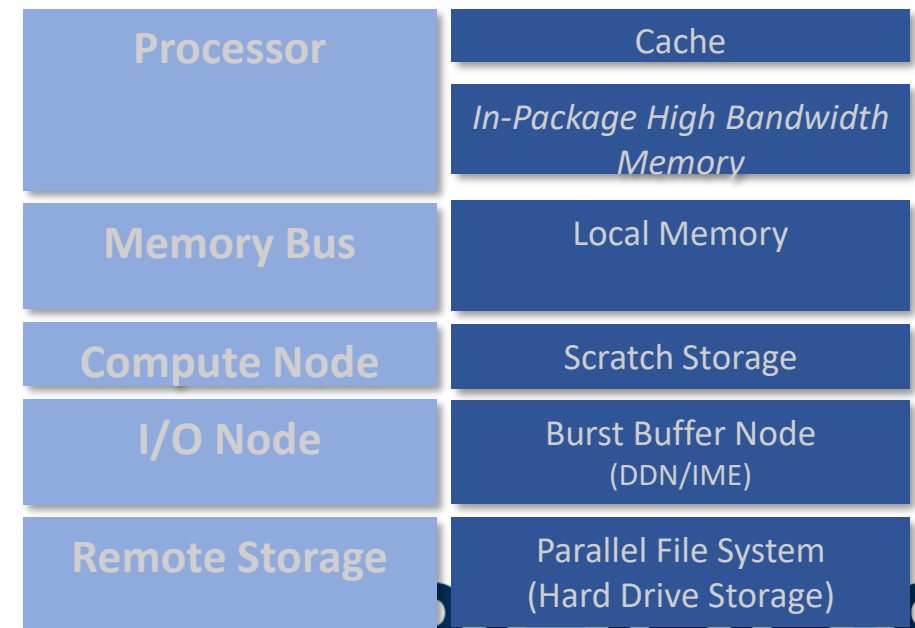    BittWare-optimized OpenCL BSP for Intel SDK

# Memory/Storage Hierarchy

❑ **Memory Sizing and Distribution**

- ○ Memory requirements of production HPC Applications
  - ○ HPL/HPCG benchmarks : **2 GB (resp. 1GB) / x86 (resp. PowerPC) core**
  - ○ UEABS benchmarks : $10^2$ MB/core
  - ○ Increase with #cores : $10^6$ cores => 10GB/core
- ○ Reduction of Application to/from Network latency
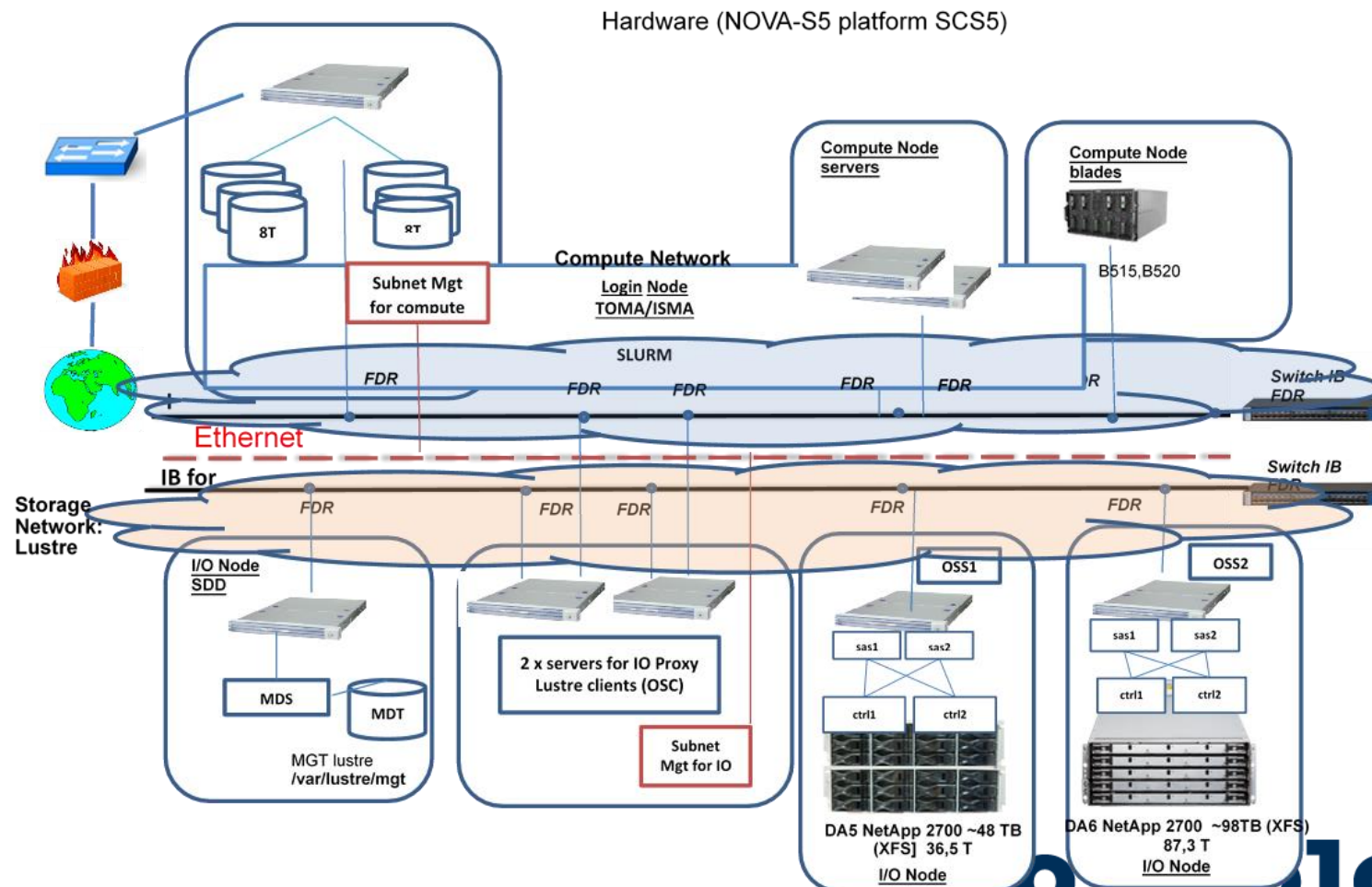
○ **Storage**

- ○ DDN's IME® scale-out, software-defined, flash storage platform
  - ○ Streamlines the data path for application I/O
  - ○ Realizes flash-cache economics with the storage
  - ○ 8 IME120 systems. For a total of 48TB flash available
- ○ Global vs Local storage
  - ○ 60 disks (1.8Gb each): 98TB
  - ○ 24 disks (1.8Gb each): 47TB
  - ○ Local storage : 1TB/Compute node
- ○ Revisited/New Features
  - ○ System/Users CheckPointing/Restart Procedures
  - ○ Local data Compression

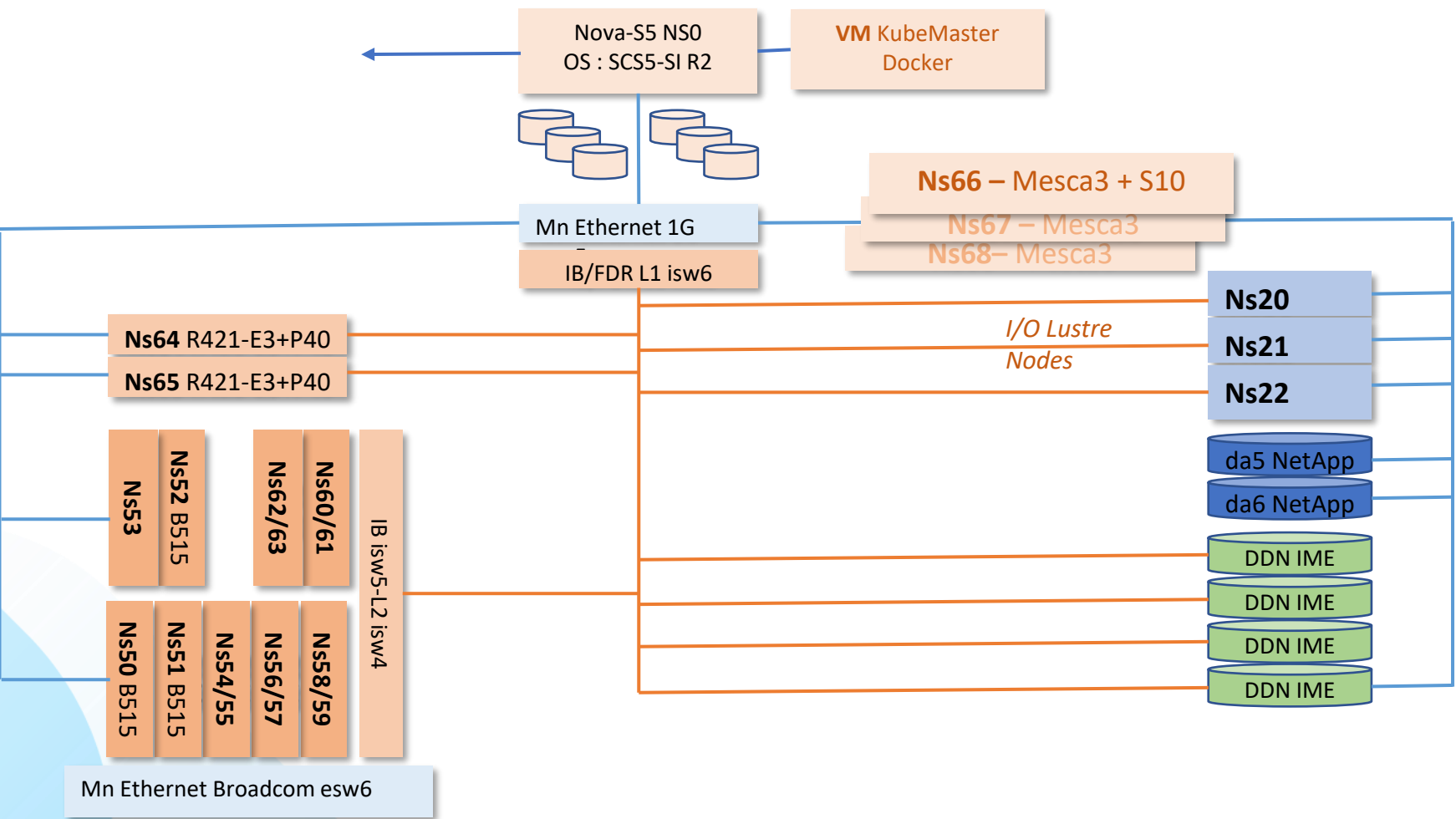| | |
|---|---|
| **Processor** | Cache |
| | *In-Package High Bandwidth Memory* |
| **Memory Bus** | Local Memory |
| **Compute Node** | Scratch Storage |
| **I/O Node** | Burst Buffer Node (DDN/IME) |
| **Remote Storage** | Parallel File System (Hard Drive Storage) |

# Interconnect Fabric

❑ **Interconnect Fabric**

- o IB/FDR (ISR 9024D and ISR 9024D-M) @56 Gb/s
- o Cisco catalyst 3560G 44 ports – private Ethernet
- o Fat-Tree Topology



Hardware (NOVA-S5 platform SCS5)

# The Hw Evolve Platform



Nova-S5 NS0
OS : SCS5-SI R2

**VM** KubeMaster
Docker

**Ns66 –** Mesca3 + S10

**Ns67 –** Mesca3
**Ns68 –** Mesca3

Mn Ethernet 1G

IB/FDR L1 isw6

**Ns64** R421-E3+P40

**Ns65** R421-E3+P40

*I/O Lustre Nodes*

**Ns20**

**Ns21**

**Ns22**

**Ns53**

**Ns52** B515

**Ns62/63**

**Ns60/61**

IB isw5-L2 isw4

**Ns50** B515

**Ns51** B515

**Ns54/55**

**Ns56/57**

**Ns58/59**

da5 NetApp

da6 NetApp

DDN IME

DDN IME

DDN IME

DDN IME

Mn Ethernet Broadcom esw6

# Conclusions

❑ **Investigation of some aspects of the HPC Evolution**

❑ General-purpose Computing Acceleration (GPU, FPGA)

❑ Restructuration of the Memory/Storage Hierarchy

❑ Co-Design

❑ **Future Development** (in the context of other projects)

o AI-Acceleration

o In-situ Computation

o Co-Design

o Evolution towards the HPC-Cloud/Edge Continuum