

## Scheda riassuntiva modulo Future Computing Architecture

### 1) Lezione

Descrizione generale delle architetture HPC e AI e loro componenti di base

- Le metriche HPL, HPCG, Green500, IO500
- I componenti base di un'architettura HPC
- Tassonomia di Flynn
- Concetti generali su architetture vettoriali
- Calcolo delle prestazioni di picco
- Limiti della scalabilità parallela

### 2) Lezione

Architetture di calcolo e loro evoluzione

- Il problema dei consumi delle CPU e limiti fisici
- Concetto di TDP
- L'architettura tipica di un sistema multi GPUs e connettività NVLINK per GPU NVIDIA

### 3) Lezione

Reti a alte prestazioni per architetture HPC e AI e loro evoluzione

- Tipologia di reti per un sistema AI&HPC
- Prestazioni di riferimento
- Concetto di bisection bandwidth
- Over-subscription rate
- Protocollo RDMA e RoCE

### 4) Lezione

Sottosistemi storage a alte prestazioni e loro evoluzione

- Flusso dati dal nodo al sistema storage e colli di bottiglia
- Architetture sistemi storage e loro caratteristiche principali
- Da tecnologia flash a nastro

## 5) Lezione

### Architetture storage a alte prestazioni

- Architettura RAID: 0, 1, 3, 5, 6
- Concetti di RAID distribuito “declustered RAID”
- Architettura Ceph
- Algoritmo di “erasure code”
- Concettigenerali di architettura scale-up e scale-out
- Concetti di block e object storage
- Benchmark IOPS
- Concetti generali su filesystem parallelo

## 6) Lezione

### Problematiche di efficientamento energetico per sistemi HPC a grande scala (architetture pre e exascale)

- TDP: perche' GPU migliore di CPU a parita' di prestazioni (GFs/Watts)
- Differenti sistemi di raffreddamento e corrispondenti intervalli di carico termico
- Perche' acqua meglio di aria
- Tecnologie di raffreddamento a confronto
- Il concetto del PUE e una sua stima di massima
- I parametri PUE, ITUE e ERE
- Nei sistemi a DLC l'importanza di modulare la temperatura acqua ingresso e la portata
- TCO: come si applica tale stima nelle gare Europee. Si prenda un'applicazione come esempio tra quelle nella slide 69 per valutare il TCO.Energy.App

## 7) Lezione

### Accenni sulle architetture innovative in ambito AI&HPC

- L'importanza di utilizzare architetture con elevata memoria condivisa per problemi di training di reti neurali complesse
- Architetture scale-up con molte GPUs per nodo (NVL72) o scale-out con poche GPUs per nodo connesse via NVLNC (NVL4)

- Valori indicativi di BW per NVLINK
- Concetti generali su architetture disaggregate
- Alcuni concetti di architetture multi-cores con elevata capacita di memoria condivisa per inferenza: Esempio Cerebras2

## 8) Lezione

Accenni al disegno e alla progettazione di un'architettura HPC

- Un breve esempio di come scegliere il sistema di raffreddamento ottimale, dato il carico termico complessivo del sistema, e una stima di massima del costo legato al consumo elettrico complessivo