

# Implementations of NVIDIA enterprise solutions in the CEE region

Petr Plodik, M Computers



# Implementations of NVIDIA enterprise solutions in the CEE region

Learn about the IT architectures and designs of supercomputers that M Computers has delivered based on NVIDIA enterprise solutions – Volta and Ampere accelerators and DGX systems in the Central and Eastern Europe (CEE) region. We will describe our experiences, best practices and lessons learnt from real-life implementations.

During the session we will cover 3 approaches:

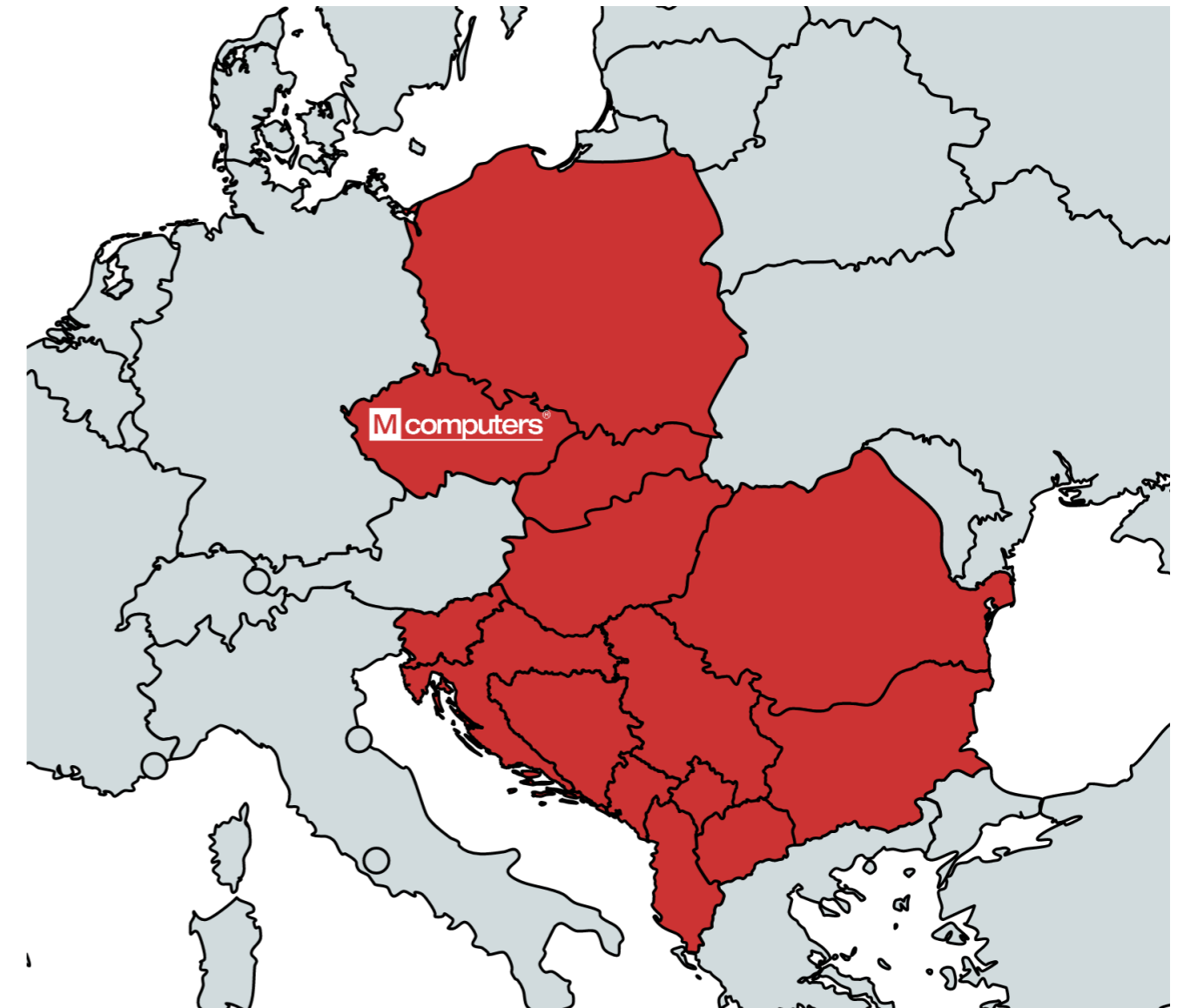
- 1) Single NVIDIA DGX System
- 2) NVIDIA reference architectures
- 3) Tailored solutions with NVIDIA accelerators

Each topic is accompanied by customer references and drawings of each IT architecture. Summary of 7 steps how to successfully implement HPC and AI projects.



## M Computers profile

- NVIDIA ELITE partner based in the Czech Republic
- Coverage: Czech Republic, Slovakia, rest of CEE
- System integrator of HPC and AI infrastructure
- Focused on the design and delivery of:
  - HPC and AI supercomputers
  - high performance and large-scale storage
  - networking
  - data center integration
  - full software stack, management, scheduling, containers, ...
- Consultations on HPC and AI topics.
- Part of NVIDIA test drive program



# 1. Single NVIDIA DGX System

Products:

- DGX Station, DGX A100,
- or older DGX-1, DGX-2

All-in-one solution for HPC and AI applications:

- compute resources (CPU and GPU)
- fast local storage (SSD/NVMe)
- connectivity to network
- DGX software stack

Applications ready to run minutes after delivery



NVIDIA DGX Station



NVIDIA DGX A100



# 1. Single NVIDIA DGX System

NVIDIA DGX Software Stack

NVIDIA NGC Cloud (hub for GPU-optimized software)

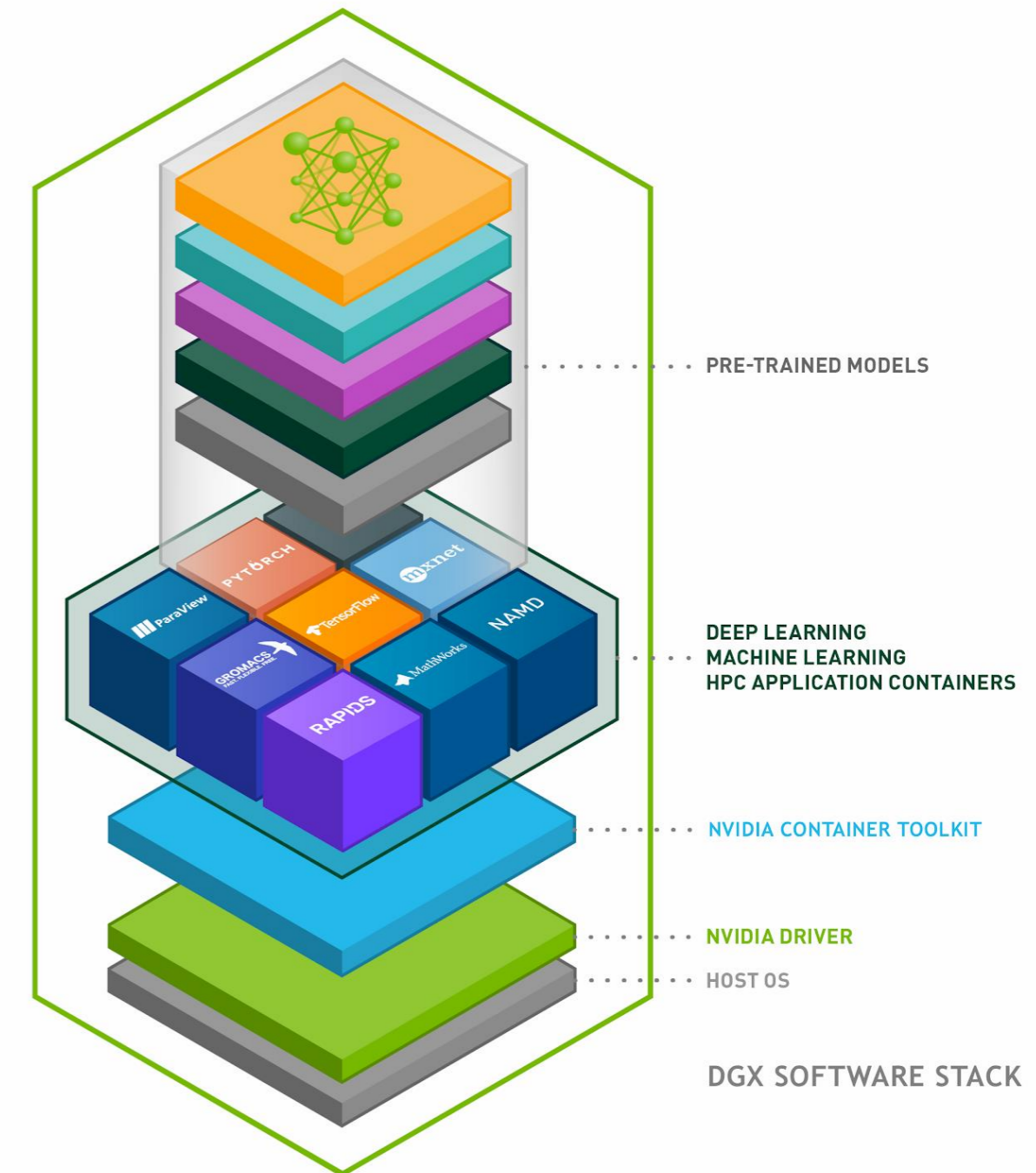
A great solution for:

- testing of HPC and AI applications
- department / workgroup workload
- required performance is within a server

DGX A100 could be divided into 56 GP GPU partitions

NVIDIA Support covers hardware and software stack

Easy migration to larger infrastructure or cloud



# 1. NVIDIA NGC

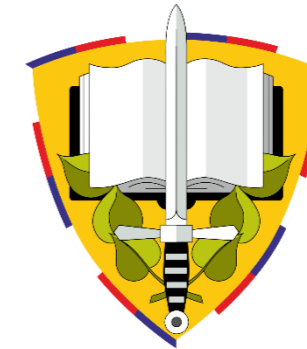
The screenshot displays the NVIDIA NGC Catalog interface. At the top, the NVIDIA logo and 'NGC CATALOG' are visible on the left, and 'Welcome Guest' is on the right. Below the header, there are navigation tabs for 'ALL CONTENT TYPES', 'COLLECTIONS', 'CONTAINERS', and 'HELM CHARTS'. A search bar is present with the text 'Search all content types'. Below the search bar, there is a grid of eight collection cards, each with a representative image, a title, a subtitle, a brief description, and a 'View Labels' link.

Collection Name	Category	Description
Translation	Natural Language Processing	A collection of easy to use, highly optimized Deep Learning Models for Machine Translation. Deep Learning Examples provides Data Scientist and Sof...
Image Segmentation	Image Classification	A collection of easy to use, highly optimized Deep Learning Models for Image Segmentation. Deep Learning Examples provides Data Scientist and Software Engi...
Image Classification	Image Classification	A collection of easy to use, highly optimized Deep Learning Models for Image Classification. Deep Learning Examples provides Data Scientist and Software Engi...
NeMo	Deep Learning	NeMo is a toolkit for creating Conversational AI applications. NeMo toolkit makes it possible for researchers to easily compose complex neural network ...
NGC - Getting Started	Beginner	Looking to get started with containers and models on NGC? This is the place to start..
Clara Deploy Operators	Healthcare	The Clara Deploy Operators Collection includes all of the reference operators that encapsulate the logic, AI algorithms, and utils to build reusable AI application pipel...
Clara Deploy Pipelines	Healthcare	The Clara Deploy Pipelines Collection includes all of the available reference pipelines for medical imaging modalities, including MRI, CT, X-Ray, Pathology, Endo...
Clara Deploy Platform	Healthcare	The Clara Deploy Platform Collection includes the bootstrap and Command Line Interface (CLI). These tools are used to install the main core services that allow y...

<https://ngc.nvidia.com/>



# 1. Single NVIDIA DGX System



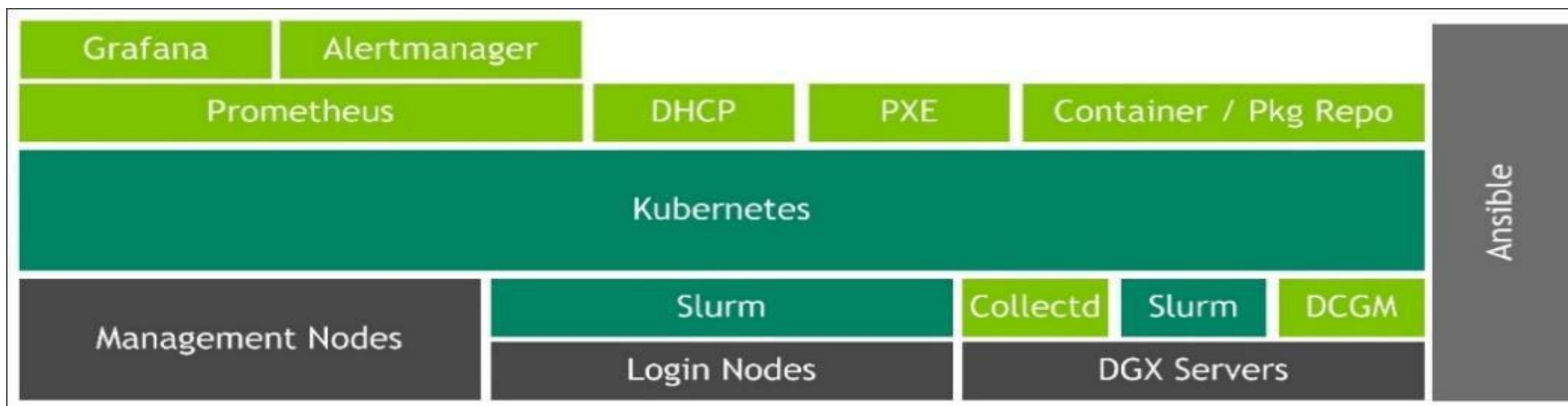
Our references in the CEE region:

- NVIDIA DGX-2 at **IT4Innovations** – National Supercomputing Center, Ostrava, CZ
  - development & testing of new large-scale algorithms for HPC and AI applications
- NVIDIA DGX Station at the **University of Defense**, Hradec Králové
  - HPC applications
- 5x NVIDIA DGX-1 MAXQ at **CIIRC, Czech Technical University**
  - AI algorithms for robotics applications



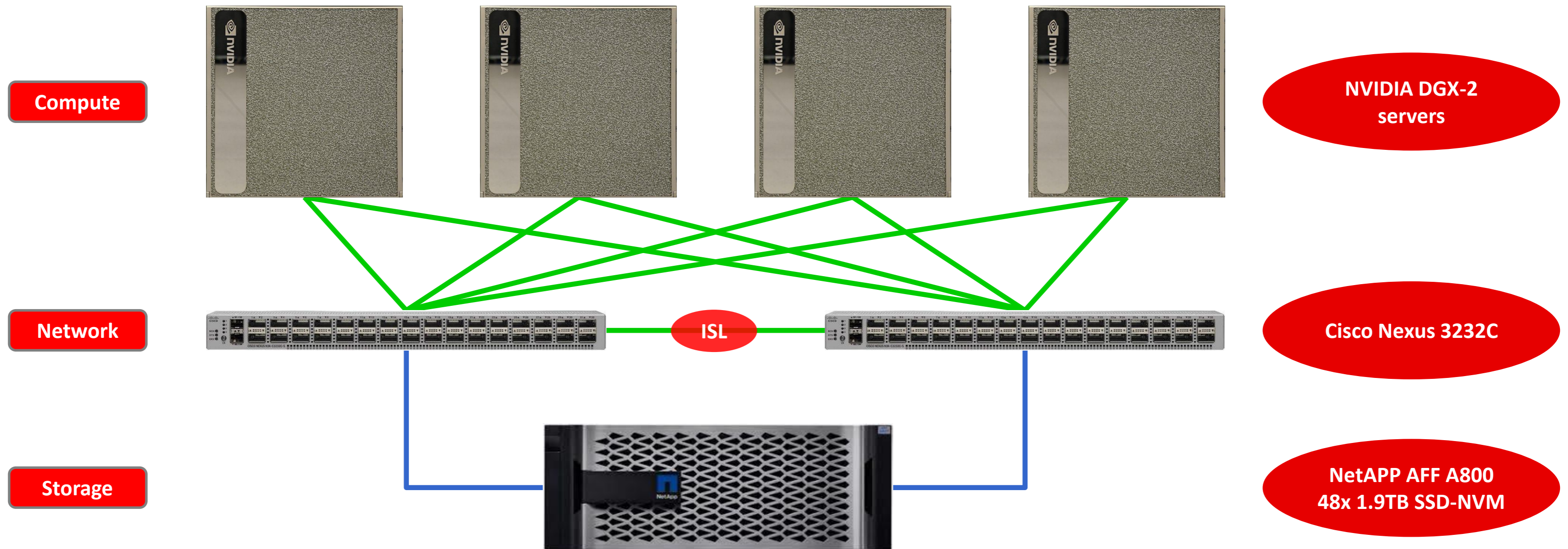
## 2. NVIDIA DGX reference architecture

- Best practices how to build HPC and AI infrastructure
- Compute nodes, fabrics, storage, system and management software stack
- ONTAP-AI = reference architecture based on NVIDIA and NetApp products
  - DGX systems, NetApp AFF A800, Cisco or Mellanox switches
- Includes data center considerations (power, cooling)
- Performance expectations of AI applications/frameworks
- Management stack is covered by NVIDIA support



## 2. NVIDIA DGX reference architecture

We have delivered two projects with NVIDIA reference architectures (DGX-1 and DGX-2)



## 2. NVIDIA DGX reference architecture



Why do we choose NetApp for NVIDIA reference architecture?

There are several storage vendors for NVIDIA reference architectures – NetApp, DELL EMC, IBM, Pure and DDN.

Based on market research we selected the NetApp solution as the best option for projects:

- NetApp is represented by ALEF Distribution in CEE region
- fast responses and support from both ALEF and NetApp teams
- 100Gb storage connectivity available from the beginning (NetApp AFF A800)
- seamless integration of Windows apps requested by some clients
- flexible and fast during price negotiation

During project implementation we appreciated:

- the technical support of HPC/AI solutions and consultations during implementation
- great support from the NVIDIA team – tuning scripts, DeepOps tool



### **3. Tailored solutions with NVIDIA accelerators**

Tuned IT architecture for specific requirements

Flexibility of choice of hardware vendors, special platforms – e.g. Supermicro, Gigabyte

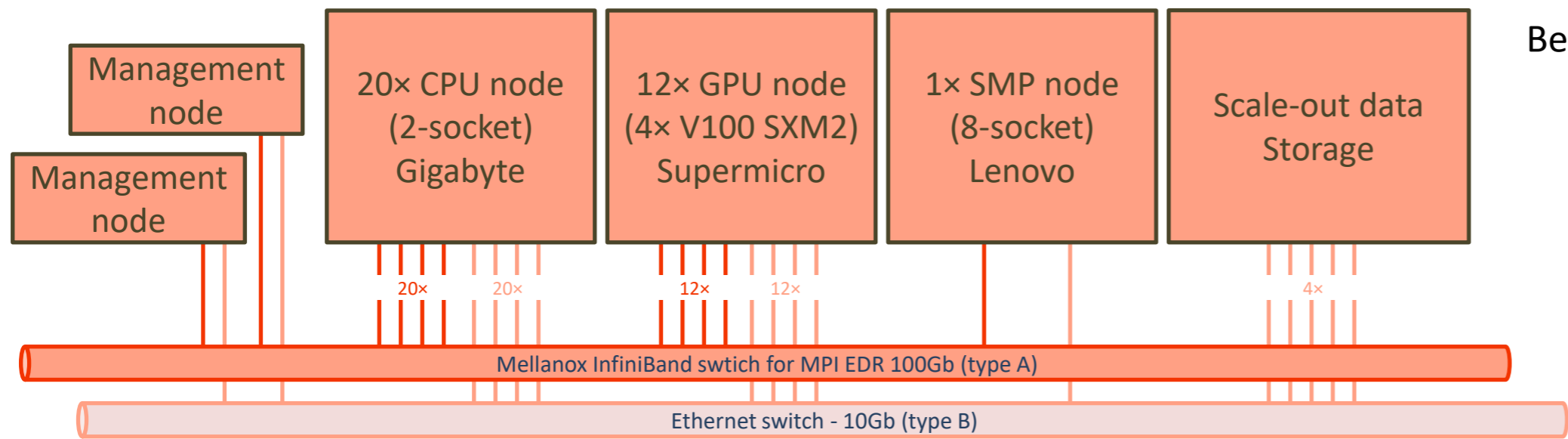
Combinations of different:

- CPU and GP GPU workloads
- different server architectures (cluster, SMP servers), form factors
- fabrics (Infiniband, Ethernet), bandwidth (1Gbps to 200Gbps)
- storage solutions (scratch, scale-out, object, SdS)
- management stack, workload scheduler, administration of container
- cooling (air, hot or cold liquid cooling)



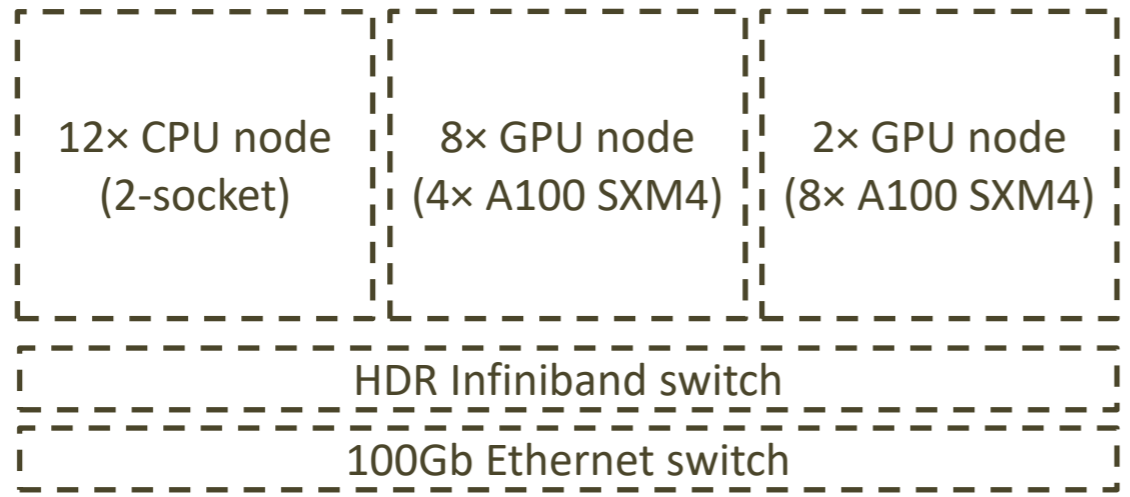
### 3. Tailored solutions with NVIDIA accelerators

#### Architecture of RCI cluster at FEL, Czech Technical University in Prague



BeeGFS is used for scratch storage on top of local NVMe drives

Planned extension:



- 100 Gb Infiniband EDR
- 10 Gb Ethernet



### 3. Tailored solutions with NVIDIA accelerators



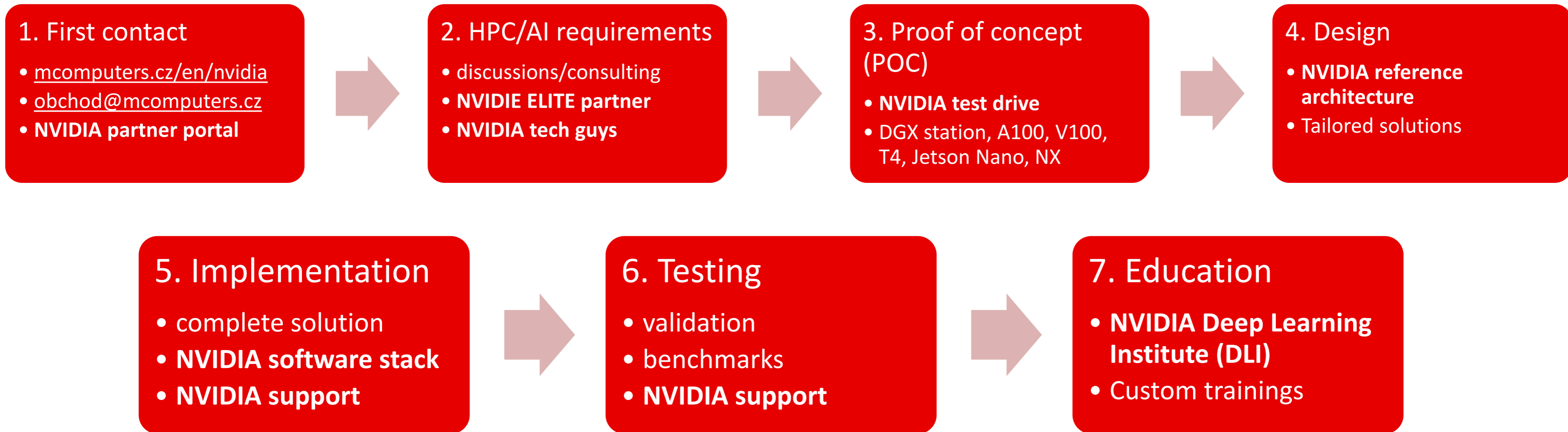
Our references in the CEE region:

- RCI cluster at the **Faculty of Electrical Engineering, Czech Technical University**
  - development & testing of AI applications
- HPC cluster projects at the **Czech Academy of Sciences**
  - several projects for biochemistry, physics, data analysis, ...
- HPC cluster projects at **Masaryk University** in Brno
  - CERIT SC projects for research in life sciences, AI, ...
- AI supercomputer at **Tomas Bata University in Zlín**
- AI supercomputer at Czech startup **Quantasoft**
  - computer vision, biometrics, AI
- NVIDIA A100 accelerators for the **Slovak Academy of Sciences**
  - upgrade of systems for AI research





# Design and implement HPC and AI projects in 7 steps



## References

NVIDIA Enterprise solutions micro-site: <https://www.mcomputers.cz/en/nvidia/>

NVIDIA reference architecture: <https://www.nvidia.com/en-us/data-center/dgx-pod/>

NVIDIA Deep Learning institute: <https://www.nvidia.com/en-us/deep-learning-ai/education/>

NVIDIA NGC repository: <https://ngc.nvidia.com/>

DGX-2 at IT4Innovations: <https://www.it4i.cz/infrastruktura/nvidia-dgx-2>

CERIT Scientific Cloud: <https://www.cerit-sc.cz/>

RCI project at FEL, Czech Technical University in Prague: <http://rci.cvut.cz/>

RICAIP project at CIIRC, Czech Technical University in Prague: <http://ricaip.eu/>

Quantasoft, Czech AI start-up: <https://www.quantasoft.com/>

Tomas Bata University in Zlin:

<https://fai.utb.cz/en/faculty/about-us/structure/departments/the-department-of-informatics-and-artificial-intelligence/>

NetApp ONTAP AI solution: <https://www.netapp.com/us/products/ontap-ai.aspx>



# Questions?



**Petr Plodik, M Computers**  
Prague, Czech Republic  
[petr.plodik@mcomputers.cz](mailto:petr.plodik@mcomputers.cz)  
+420 737 264 480

