**4th SmartData@PoliTO Workshop**

**WIDENING OUR HORIZONS**

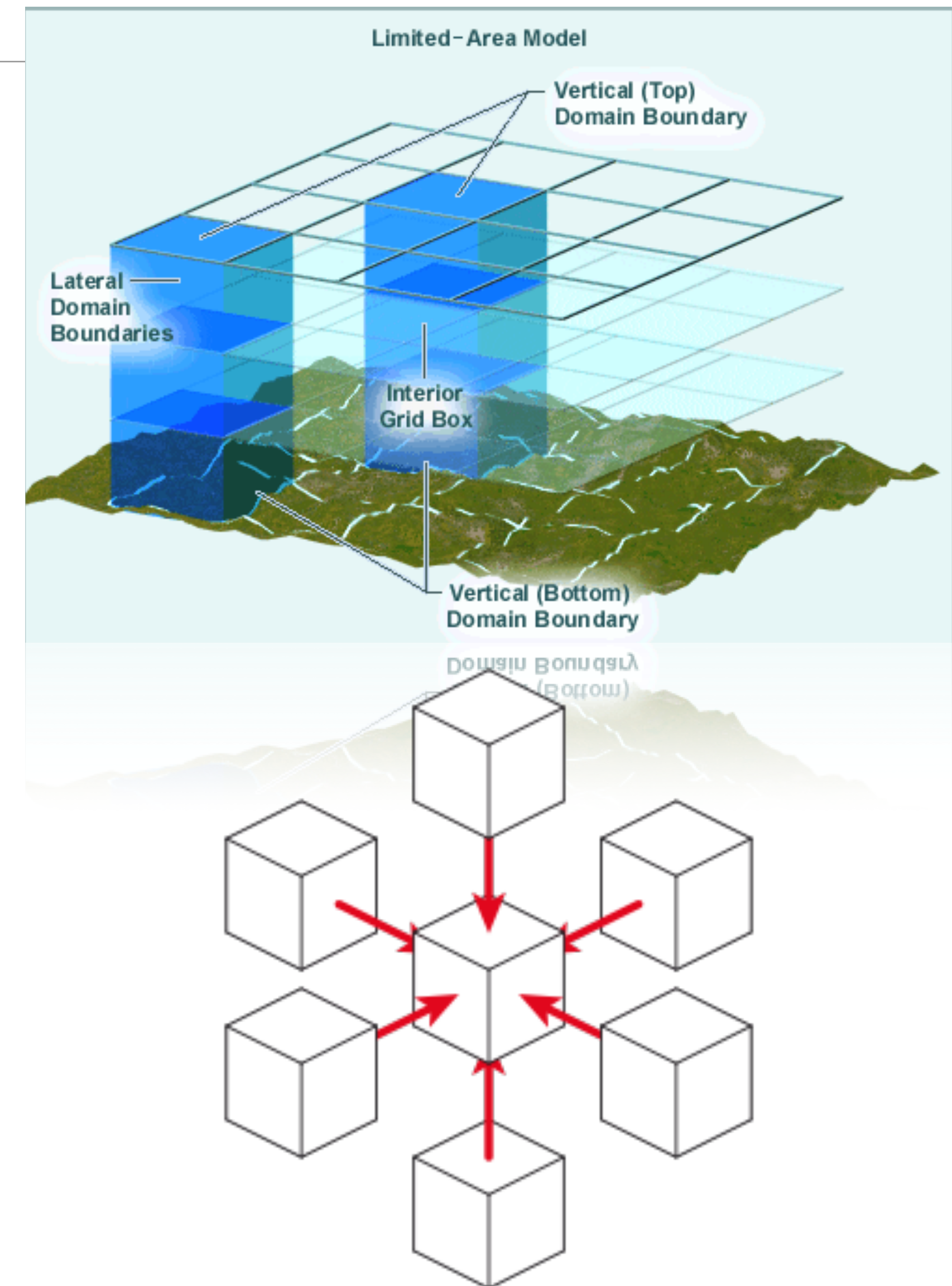# THE EVOLUTION OF HIGH-PERFORMANCE SYSTEMS: FROM HPC TO BIG DATA TO DEEP LEARNING

Marco Aldinucci
Computer Science Department, University of Torino
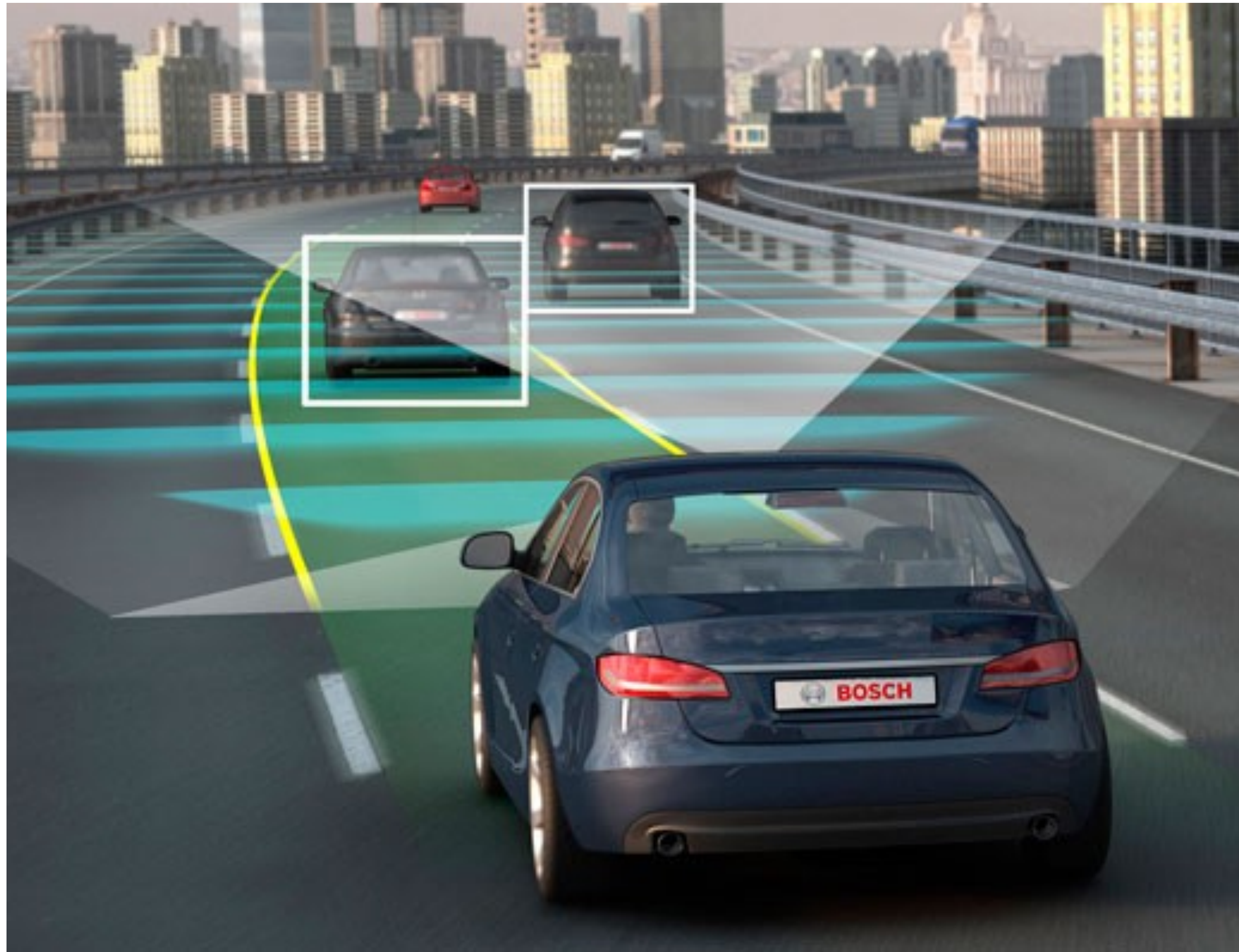
Head of **Parallel Computing group**, coordinator of HPC4AI,
national delegate at EuroHPC JU

http://EuroHPC.eu

# DATA LOSES ITS OPERATIONAL VALUE IN A SHORT TIME

- Partition the atmosfere 1x1x1 Km up to 10 Km altitude (10 cells). Earth: ~ $5 \times 10^9$ cells

- Let us supposed 1 cell needs 200 Floating Point Ops. I.e. $10^{12}$ Ops for each step

- A 7-days forecast with 1 minute time step on a 10GFLOPS CPU needs $10^7$ secs, i.e 10 days
  - Es. Intel Core i7-7500U

- To compute it in 5 mins you need a 35 TFLOPS CPU
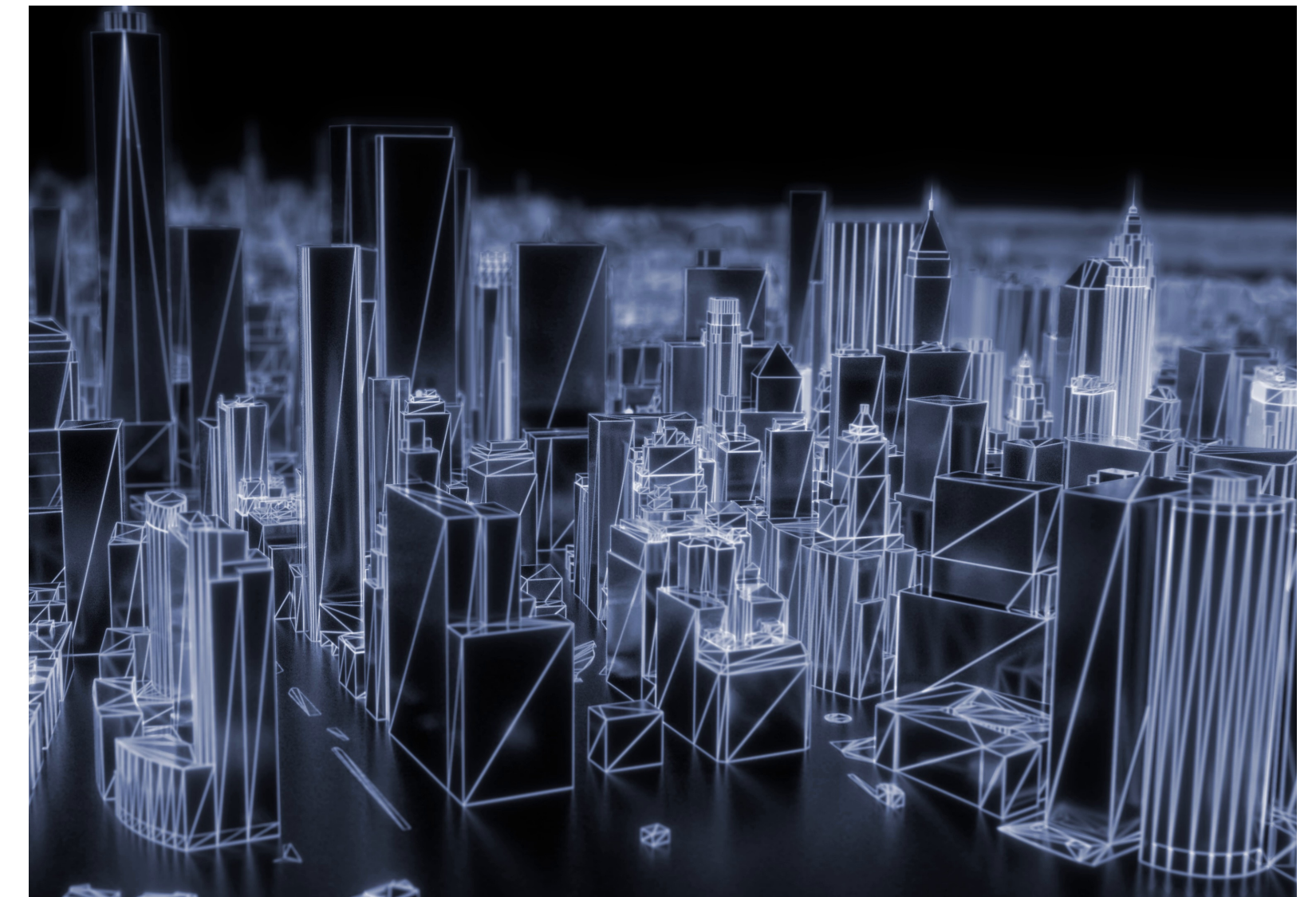  - 3500x (fastest CPU today, ~200GFLOPS)

# Self-Driving Cars Can Handle Neither Rain nor Sleet nor Snow

● To help autonomous vehicles solve inclement conditions, WaveSense will sell a sensor that can see below the ground.
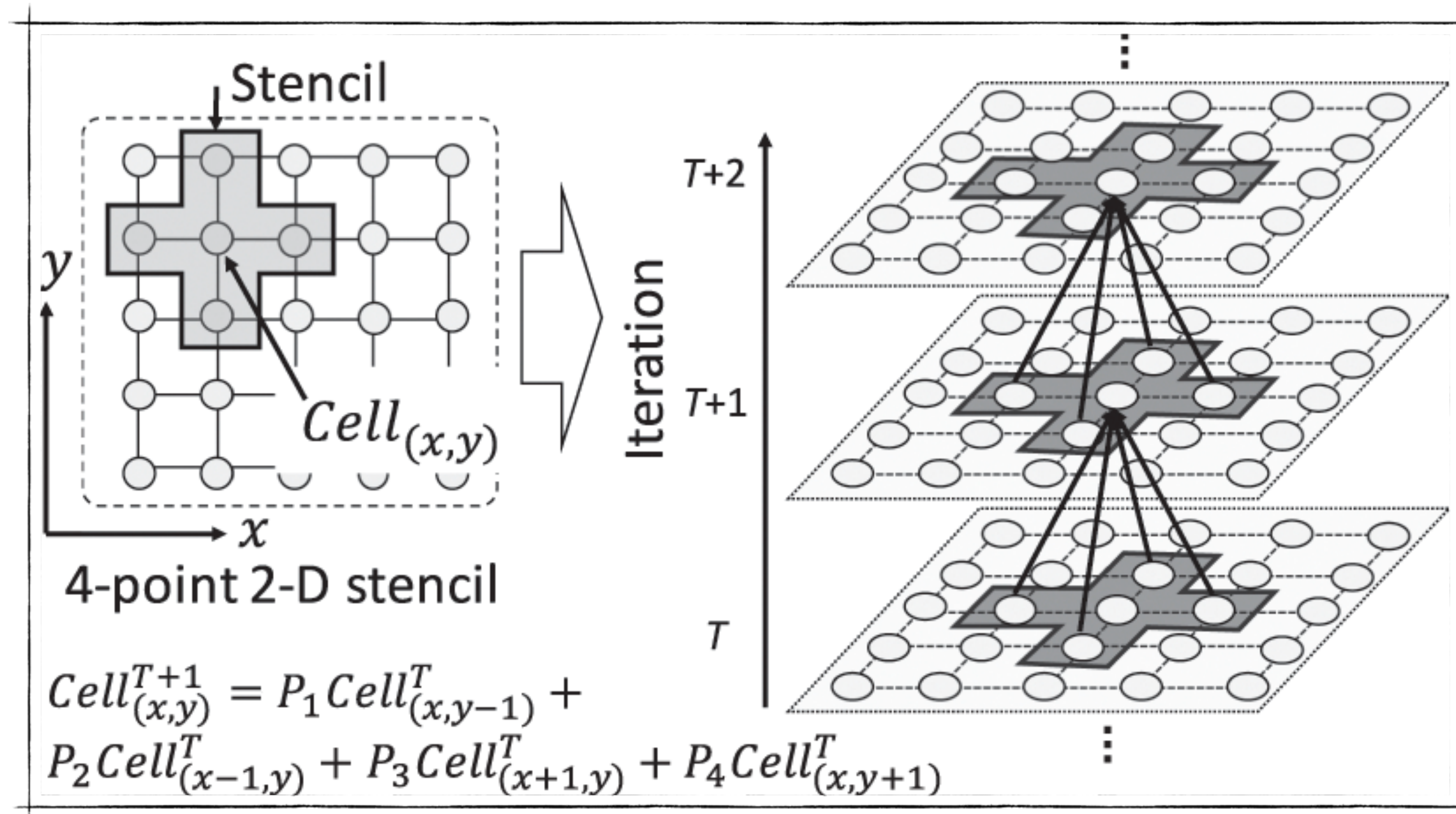
By Kyle Stock

**Hyperdrive**



# 100X100M CONTINUOUS WHETHER FORECAST, WHAT FOR?

Stencil

$Cell_{(x,y)}$

4-point 2-D stencil

$$Cell_{(x,y)}^{T+1} = P_1 Cell_{(x,y-1)}^{T} + P_2 Cell_{(x-1,y)}^{T} + P_3 Cell_{(x+1,y)}^{T} + P_4 Cell_{(x,y+1)}^{T}$$
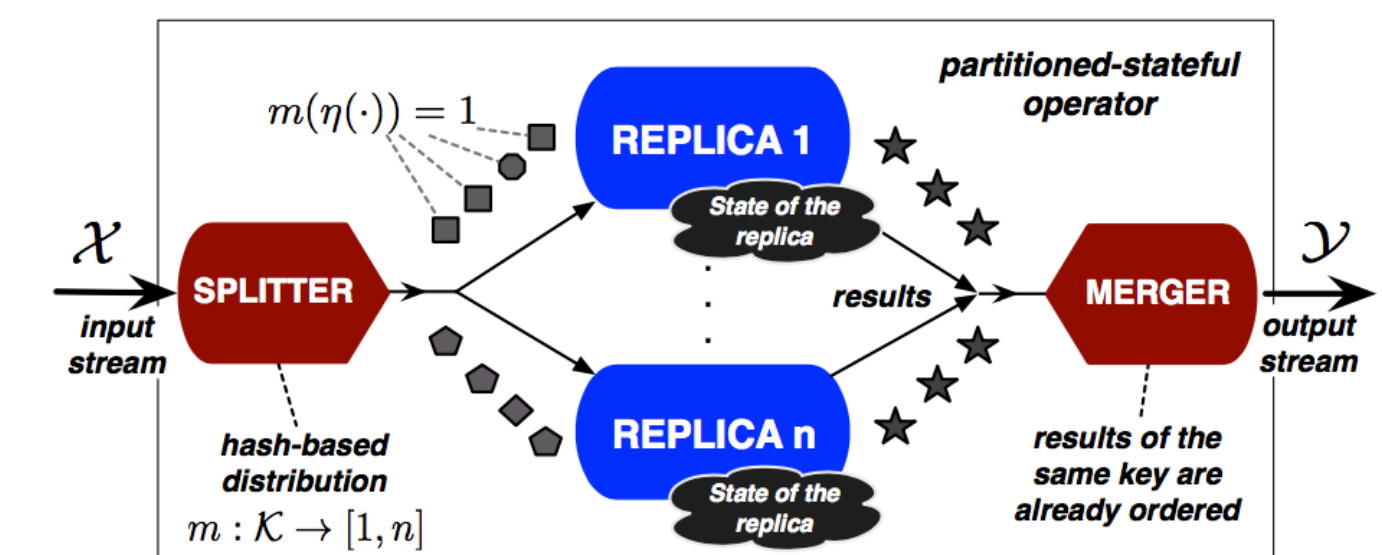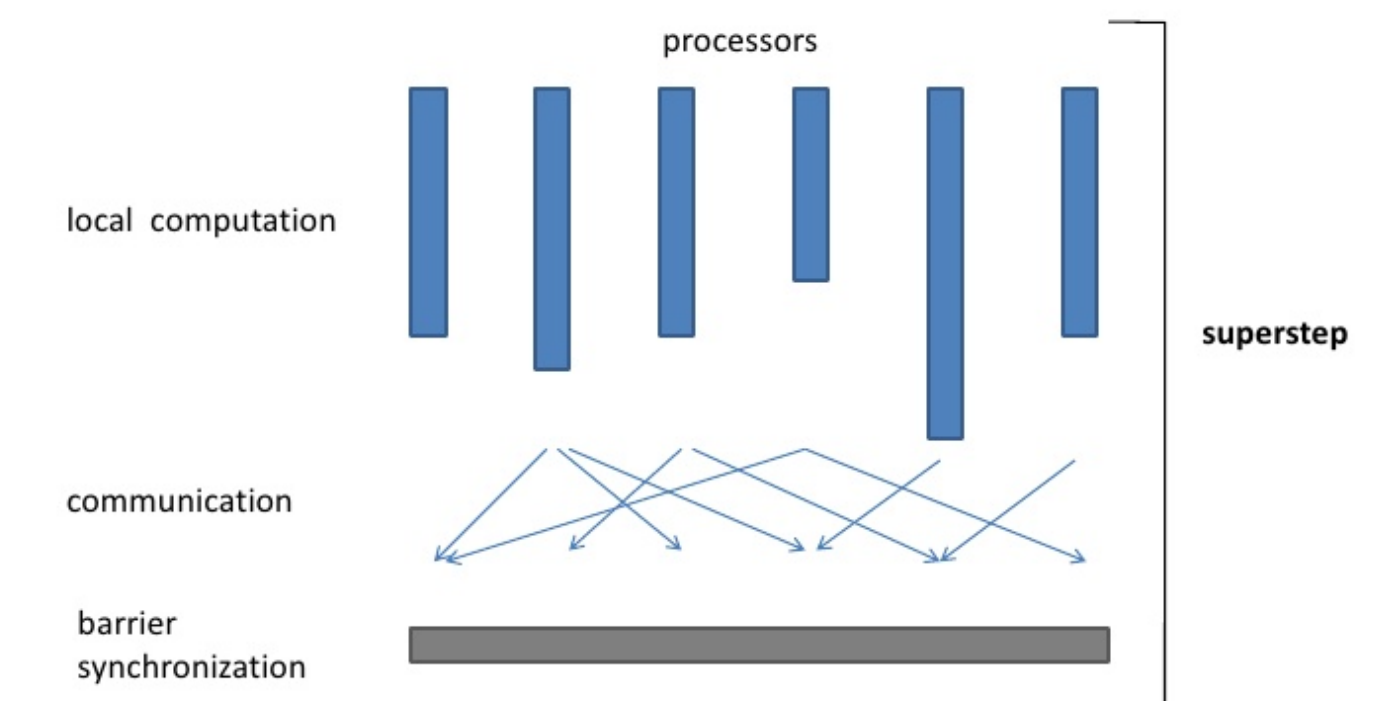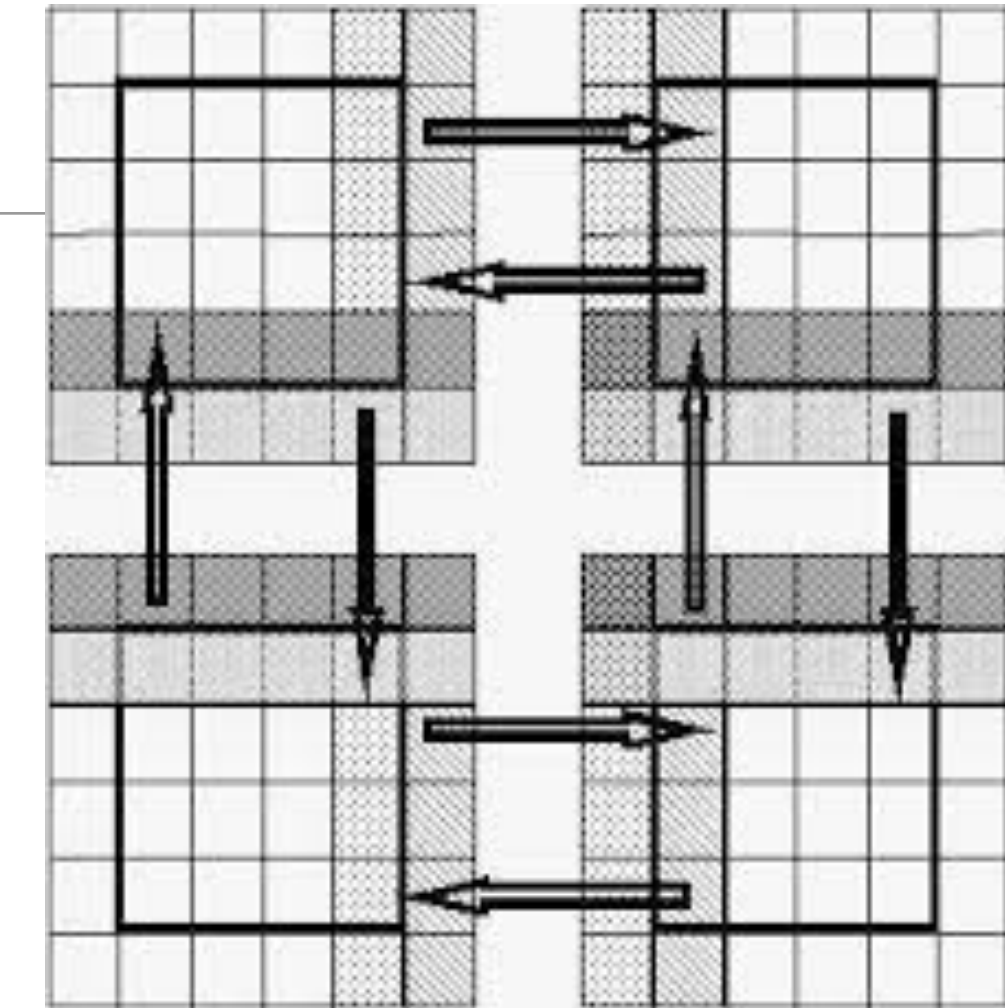
Iteration

$T+2$

$T+1$

$T$

# PROGRAMMING MODELS AND APPLICATIONS
# HPC

- Locally synchronous data parallelism

  - Stencil

  - Scalable and well-understood

  - Simulations, finite elements, AMR, HPLinpack, …

  - Physics, chemistry, engineering, …

- Globally synchronous data parallelism

  - BSP, PGAS, …

  - Barrier = not scalable

- Stream & task parallel

  - Compositional = good, but not scalable





*M. Aldinucci et al. "A Parallel Pattern for Iterative Stencil + Reduce," Journal of Supercomputing, 2018*

6

# HPC APPLICATIONS

## L'OMA SEMPER FAIT PAREI… (SARA PÀ PERICULUS)



Number of applications using a given programming paradigm in the set of 30 candidate exascale applications of the USA DoE Exascale Computing Project (ECP).

# PLATFORMS



- ## Multi-level clusters

  - Multicore + GPUs
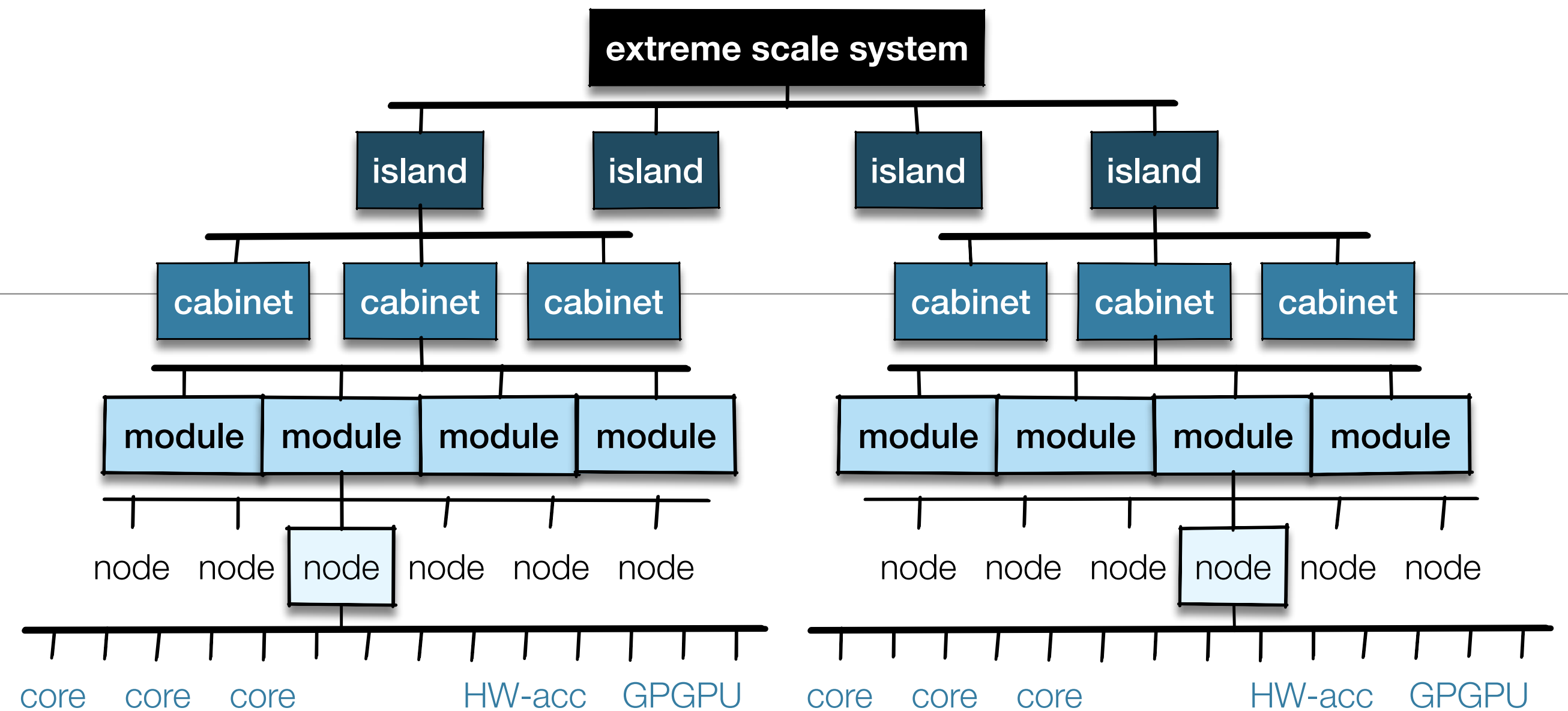
- ## Data movement

  - Optimised for CPU-bound applications, I/O-bound applications problematic

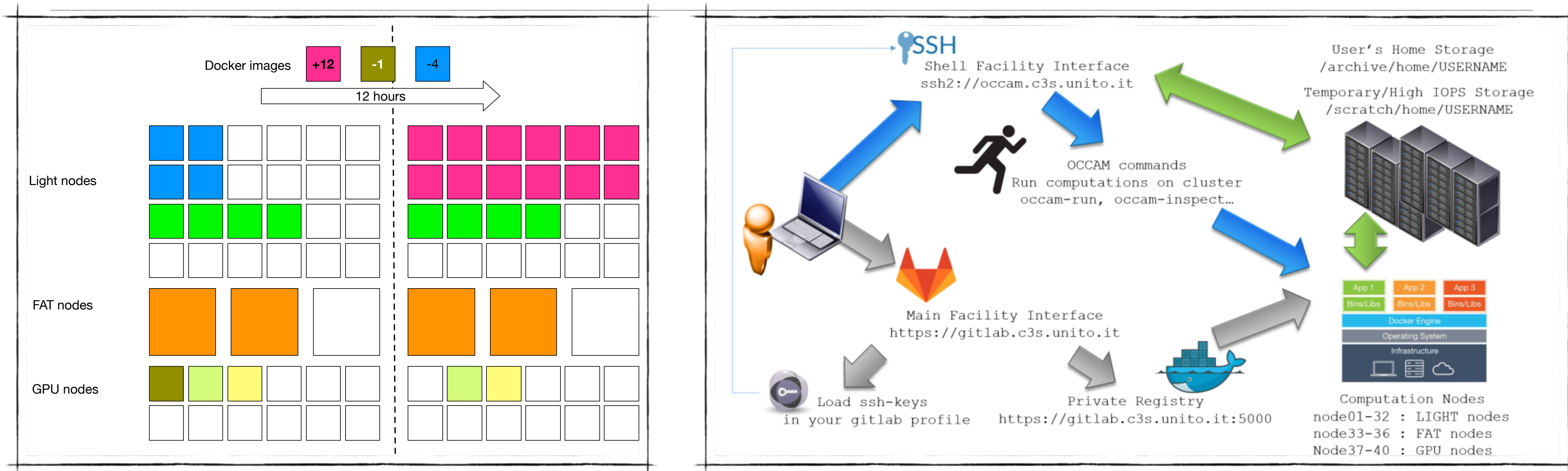- ## High-speed network (e.g. fat-tree IB 100-400 Gb/s)

  - Efficient synchronisation is key for scalability

- ## Managed with job queue (PBS, SLURM)

  - Sometime bare metal virtualisation: singularity, kubernetes, …Occam@UNITO

# Occam@UNITO: calendar + docker



- Advantages: Interactive, virtual farms + queue, user-defined configuration

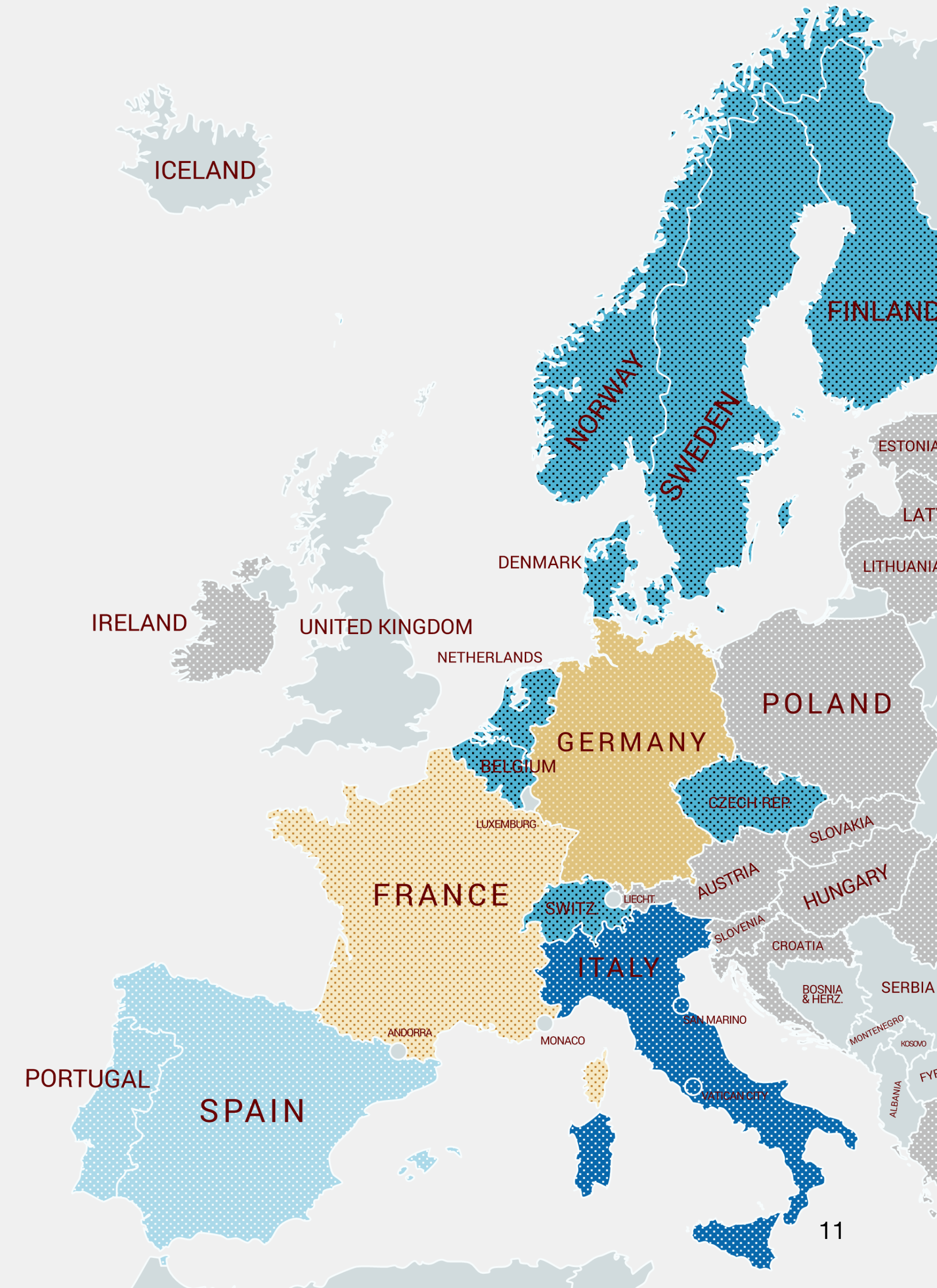- More importantly… **we designed it!**

# EUROHPC
# EU PROCESSOR + EU EXA-MACHINE

Pre-exascale general system specifications

- Applicants must describe how the following general system specifications will be met, for both the EuroHPC supercomputer and the site.

- The hosting entity will host a supercomputer with the following requirements:
  - A capability computing system, with an aggregated performance level capable of executing at least **200 Petaflops** (**sustained performance** measured using linpack benchmark)
  - Covering the needs (including substantial performance increase) of a wide range of key/grand challenge applications that demonstrably require large systems that are precursors of exascale capability computing.
  - A total power consumption of no more than **15 MW** for the hosting of the EuroHPC supercomputer

**(pre)-exascale candidates for EuroHPC**

- Pre-exascale – Finland led consortium
- Pre-exascale – Italy
- Pre-exascale – Spain & Portugal
- Exascale – Germany
- Exascale – France
- Other EuroHPC countries

# "THE MOST DAMAGING PHRASE IN THE LANGUAGE IS L'OMA SEMPER FAIT PAREI" (WE'VE ALWAYS DONE IT THIS WAY!)

## The grandma of Cobol

## Grace Hopper

From Wikipedia, the free encyclopedia

**Grace Brewster Murray Hopper** (née **Murray**; December 9, 1906 – January 1, 1992) was an American computer scientist and United States Navy rear admiral.[1] One of the first programmers of the Harvard Mark I computer, she was a pioneer of computer programming who invented one of the first compiler related tools. She popularized the idea of machine-independent programming languages, which led to the development of COBOL, an early high-level programming language still in use today.

attempted to enlist in the Navy during World War II but was rejected because she was 34 years old. She joined the Navy Reserves. Hopper began her computing career in 1944 when she worked on the Mark I team led by Howard H. Aiken. In 1949, she joined the Eckert–Mauchly Computer Corporation part of the team that developed the UNIVAC I computer. At Eckert–Mauchly she began developing the r. She believed that a programming language based on English was possible. Her compiler converted terms into machine code understood by computers. By 1952, Hopper had finished her program linker lly called a compiler), which was written for the A-0 System.[2][3][4][5]

Eckert–Mauchly chose Hopper to lead their department for automatic programming, and she led the of some of the first compiled languages like FLOW-MATIC. In 1959, she participated in the CODASYL ium, which consulted Hopper to guide them in creating a machine-independent programming language. to the COBOL language, which was inspired by her idea of a language being based on English words. , she retired from the Naval Reserve, but in 1967, the Navy recalled her to active duty. She retired from y in 1986 and found work as a consultant for the Digital Equipment Corporation, sharing her computing nces.

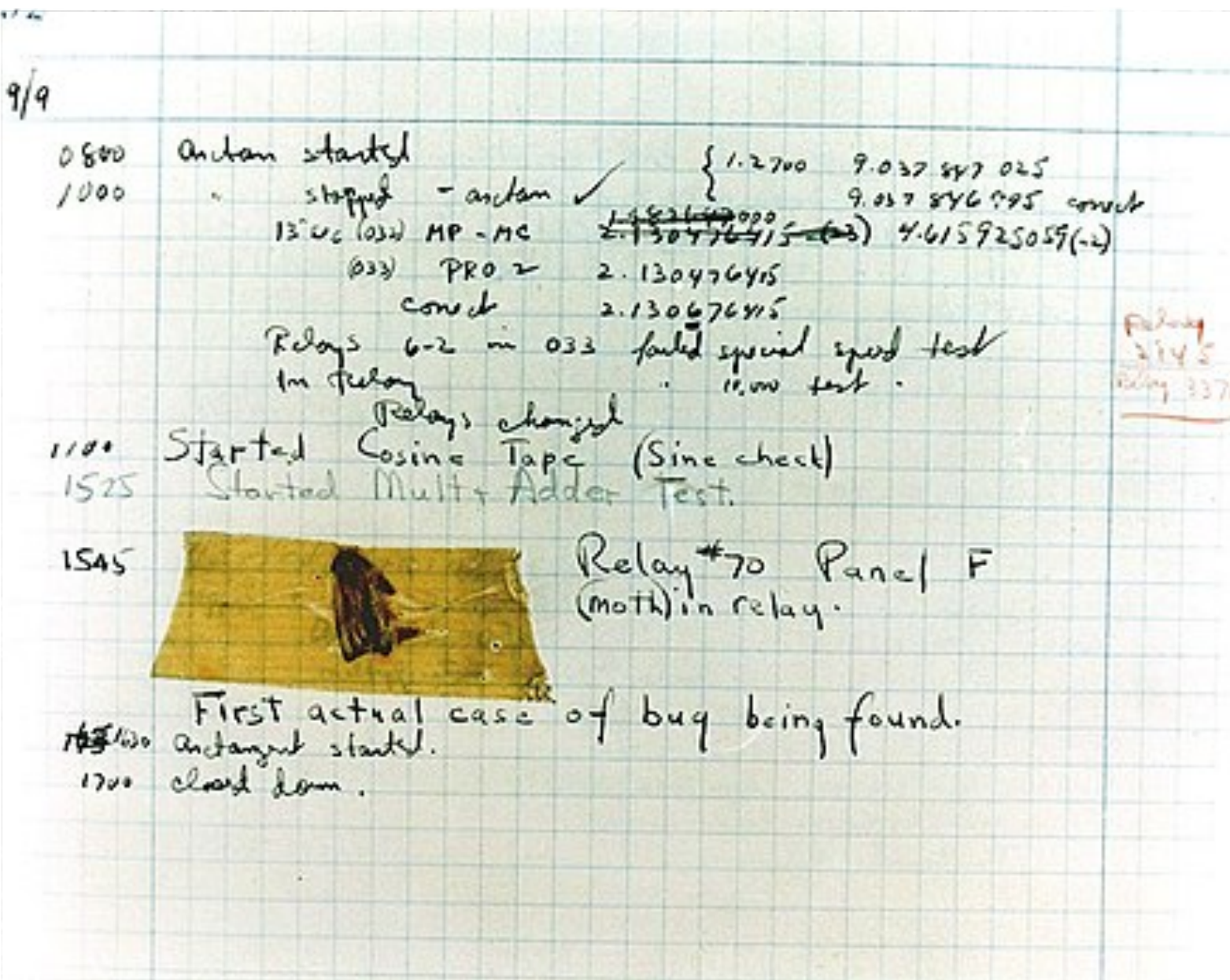to her accomplishments and her naval rank, she was sometimes referred to as "Amazing Grace".[6][7]
S. Navy Arleigh Burke-class guided-missile destroyer USS Hopper was named for her, as was the Cray XE6 "Hopper" supercomputer at NERSC.[8] During her lifetime, Hopper was awarded 40 honorary degrees from

**Grace Murray Hopper**



Rear Admiral Grace M. Hopper, 1984

| | |
|---|---|
| **Born** | December 9, 1906 New York City, New York, U.S. |
| **Died** | January 1, 1992 (aged 85) Arlington, Virginia, U.S. |
| **Other names** | "Amazing Grace", "Grandma COBOL" |
| **Alma mater** | Yale University |
| | **Military career** |

# Parallel Computing group funding perspective (Alpha@UNITO)

- **ParaPhrase** (EC-STREP, 7th FP): Parallel Patterns for Adaptive Heterogeneous Multicore Systems (2011, 42 months, total cost 4.2M €).

- **BETTY** (EC-COST Action1201): Behavioural Types for Reliable Large-Scale Software Systems. (2012, 48 months).

- **CINA** (MIUR PRIN): Compositionality, Interaction, Negotiation, Autonomicity for the future ICT society (2013, 36 months).

- **REPARA** (EC-STREP, 7th FP): Reengineering and Enabling Performance And poweR of Applications (2013, 36 months, total cost 3.5M €).

- **NESUS** (EC-COST Action IC1305): Network for Sustainable Ultrascale Computing (2014, 48 months, total cost 400K).

- **C3S**: Competence Center on Scientific Computing (2014, Compagnia di San Paolo, founding 900K €).

- **Rephrase** (EC-RIA, H2020, ICT-2014-1): Refactoring Parallel Heterogeneous Resource-Aware Applications – a Software Engineering Approach (2015, 36 months, total cost 3.5M €).

- **cHiPSet** (EC-COST Action IC1406): High-Performance Modelling and Simulation for Big Data Applications (2015, 48 months, total cost 500K).

- **OptiBike** (EU I4MS): Robust Lightweight Composite Bicycle design and optimization, an experiment of EU i4MS Fortissimo2 project (2017, 24 months, total cost 230K €).

- **Toreador** (EC-RIA, H2020, ICT-2015-16): TrustwOrthy model-awaRE Analytics Data platfORm (2015, 36 months, total cost 6.2M €).

- **HPC4AI** (Regione Piemonte, INFRA_P): Turin's centre in High-Performance Computing for Artificial Intelligence (2018, 24 months, total cost 4.5M €).

- **DeepHealth** (EC-IA, H2020, ICT-2018-11): Deep-Learning and HPC to Boost Biomedical Applications for Health (2019, 36 months, total cost 12.8M€).

- **MnemoComputing** (Compagnia di SanPaolo): Components for Processing In Memory (2019, 24 months, total cost 70K€)

- **PDEvolve** Deep Learning + HPC (ICT-11-b - under evaluation)

- **ExaTrain** Deep Learning + HPC (Marie Curie - under evaluation)
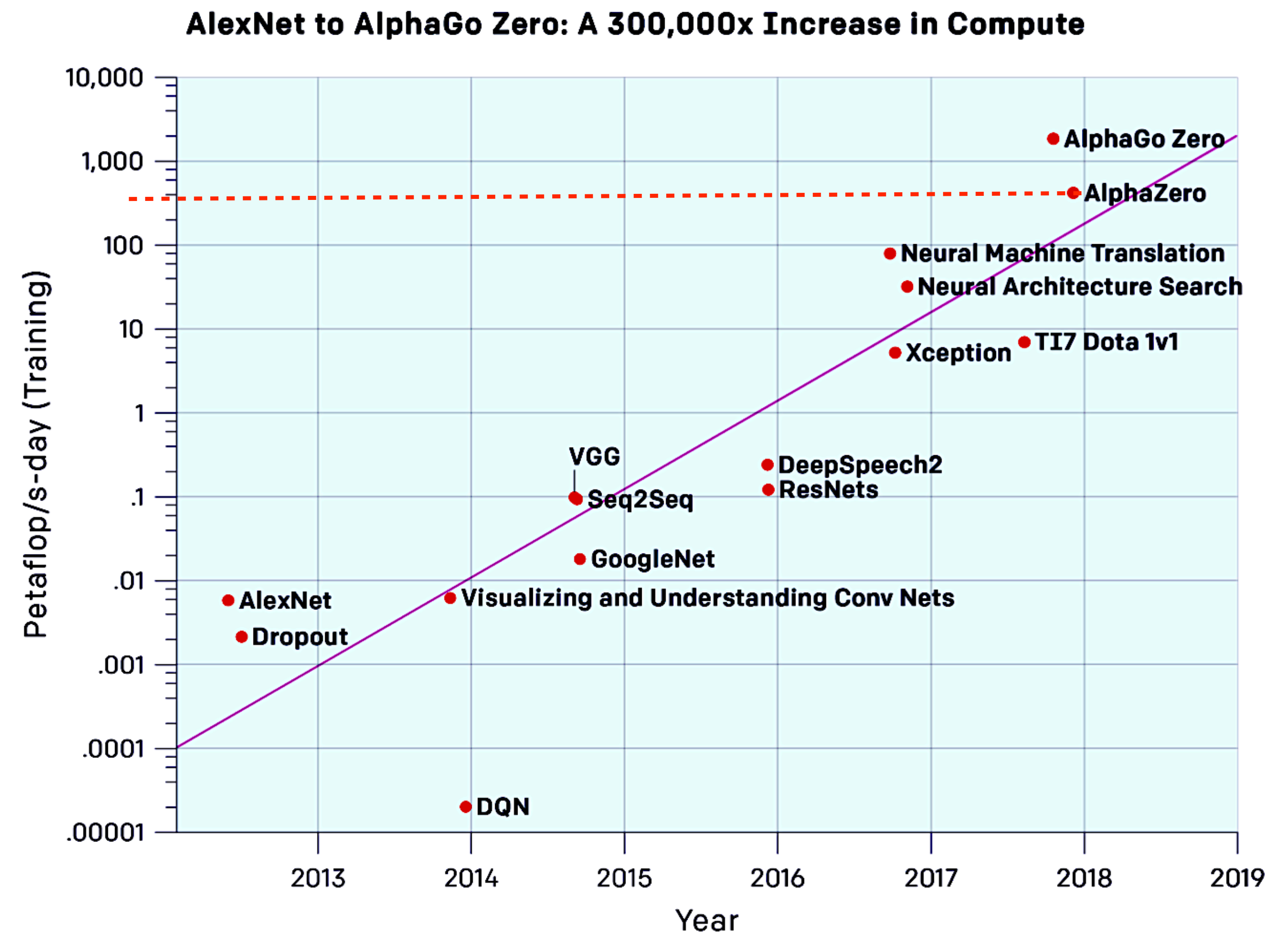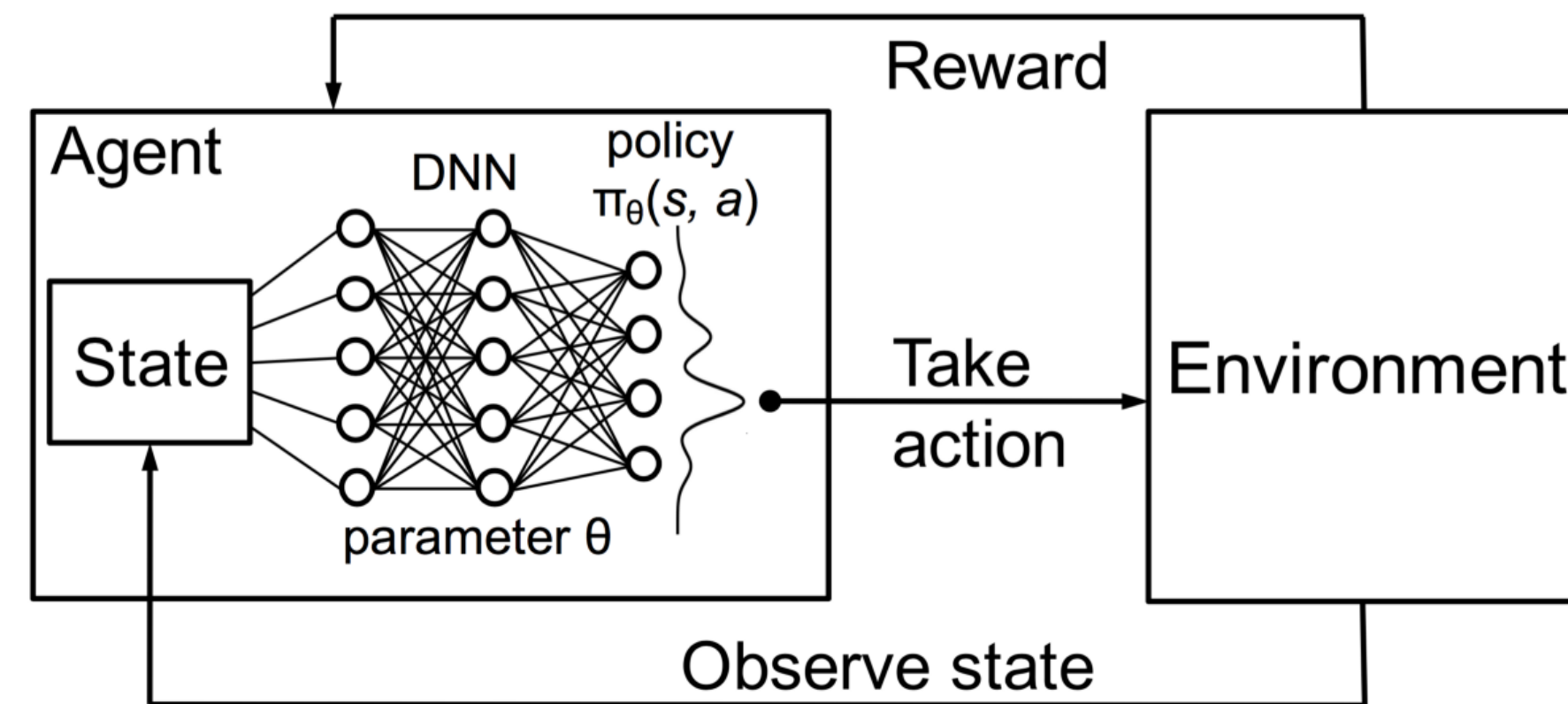
Parallel computing
HPC

2011-

BigData
Cloud

2015-

Machine Learning
Cloud

2018-

# REINFORCED DEEP LEARNING





It might look goofy …

AlexNet to AlphaGo Zero: A 300,000x Increase in Compute

# FLOPS

| Rank | Site | System | Cores | Rmax (TFlop/s) | Rpeak (TFlop/s) | Power (kW) |
|------|------|--------|-------|----------------|-----------------|------------|
| 1 | DOE/SC/Oak Ridge National Laboratory United States | **Summit** - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM | 2,397,824 | 143,500.0 | 200,794.9 | 9,783 |
| 2 | DOE/NNSA/LLNL United States | **Sierra** - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband IBM / NVIDIA / Mellanox | 1,572,480 | 94,640.0 | 125,712.0 | 7,438 |
| 19 | CINECA Italy | **Marconi Intel Xeon Phi** - CINECA Cluster, Lenovo SD530/S720AP, Intel Xeon Phi 7250 68C 1.4GHz/Platinum 8160, Intel Omni-Path Lenovo | 348,000 | 10,384.9 | 18,816.0 | |

- NVIDIA V100
  - 32 bit ~15TFLOPS
  - 5120 cores + 640 tensor cores
  - 300W, cost: ~7K€

- Intel  E5-2699v4
  - 64 bit ~200GFLOPS
  - 22 cores HT + AVX2
  - 150W, cost: ~4K€

- Marconi
  - 64 bit ~10PFLOPS
  - 350,000 KNL 68 cores
  - 3MW, cost: 30M€



Min training time AlphaZero ~36 years



Min training time AlphaZero ~684 years



Min training time

AlphaZero ~1 hour

# THE NEW TURING GPU
# LESS THAN 2500€

## NVIDIA TESLA T4 SPECIFICATIONS

| Performance | |
|---|---|
| | TURING TENSOR CORES |
| | **320** |
| | NVIDIA CUDA® CORES |
| | **2,560** |
| | SINGLE PRECISION PERFORMANCE (FP32) |
| | **8.1** TFLOPS |
| | MIXED PRECISION (FP16/FP32) |
| | **65** FP16 TFLOPS |
| | INT8 PRECISION |
| | **130** INT8 TOPS |
| | INT4 PRECISION |
| | **260** INT4 TOPS |
| Interconnect | GEN3 |
| | **x16** PCIe |
| Memory | CAPACITY |
| | **16** GB GDDR6 |
| | BANDWIDTH |
| | **320+** GB/s |
| Power | **70** watts |



Figure 4. Turing TU102/TU104/TU106 Streaming Multiprocessor (SM)

16

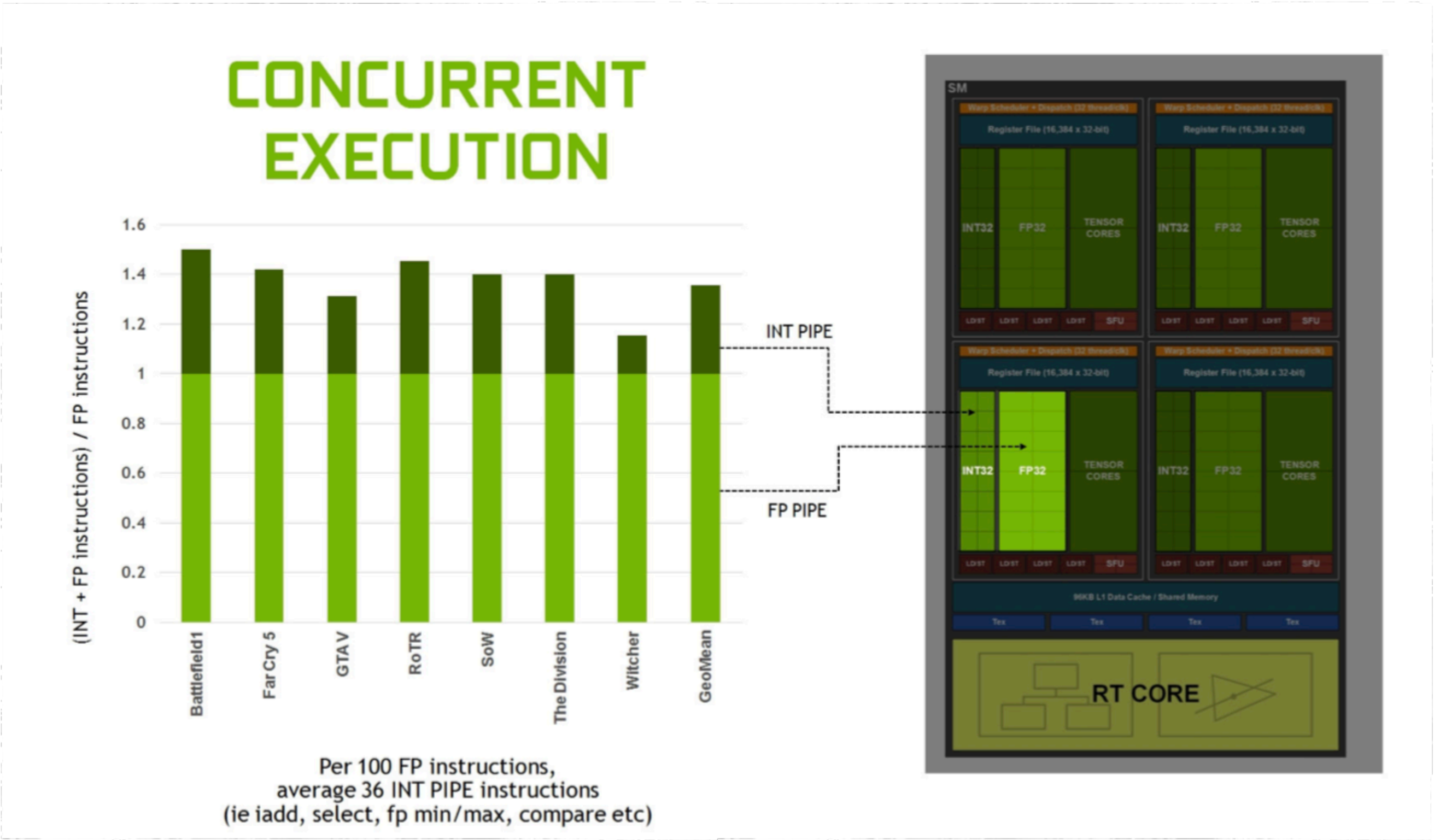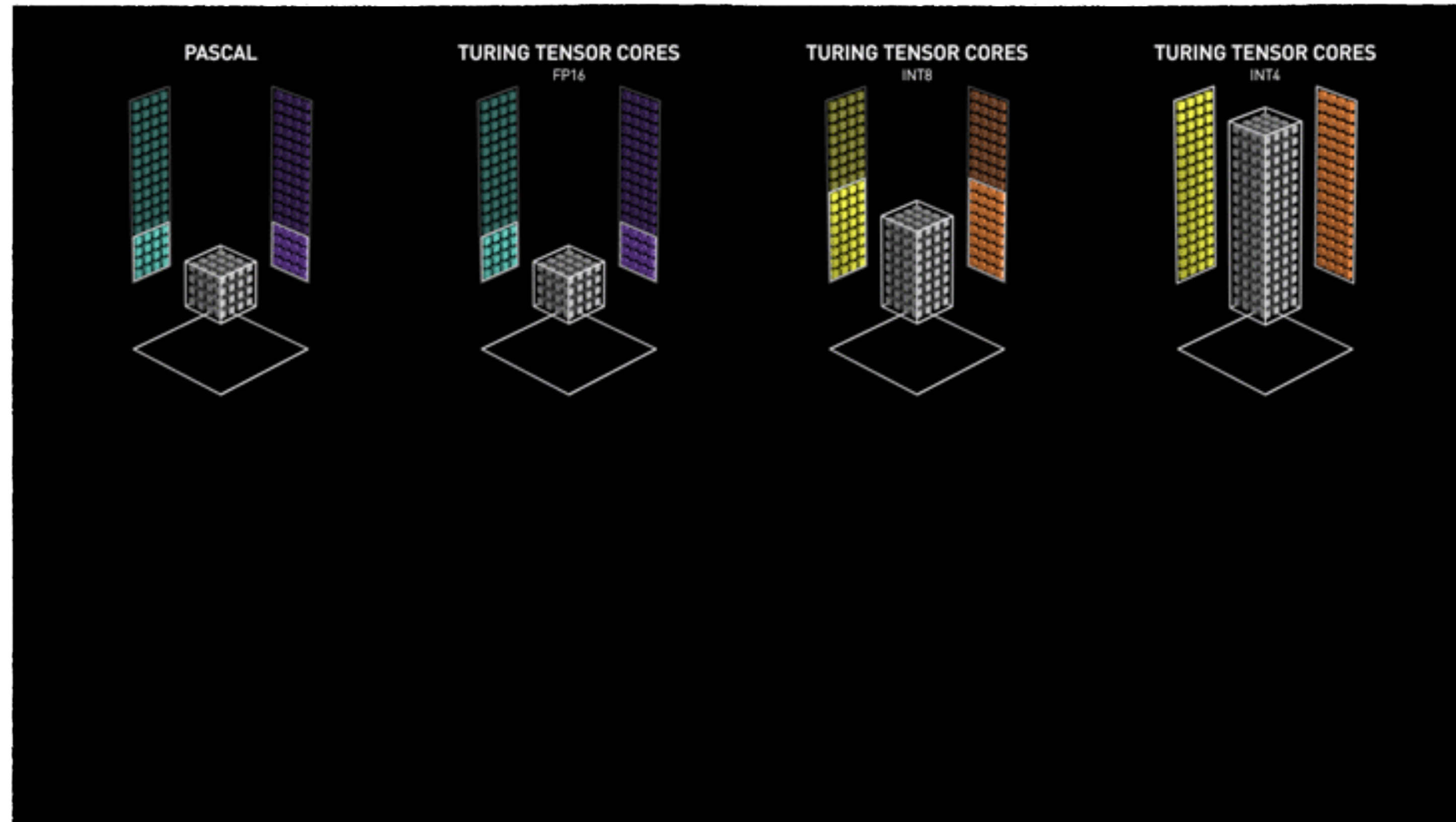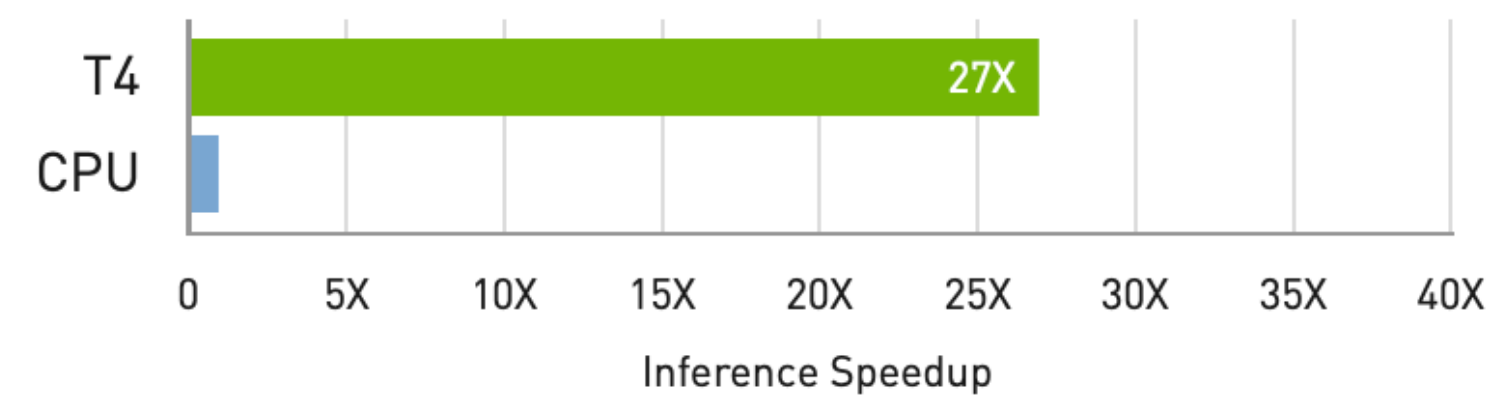# THE NEW TURING GPU
# LESS THAN 2500€



Figure 4.    Turing TU102/TU104/TU106 Streaming Multiprocessor (SM)

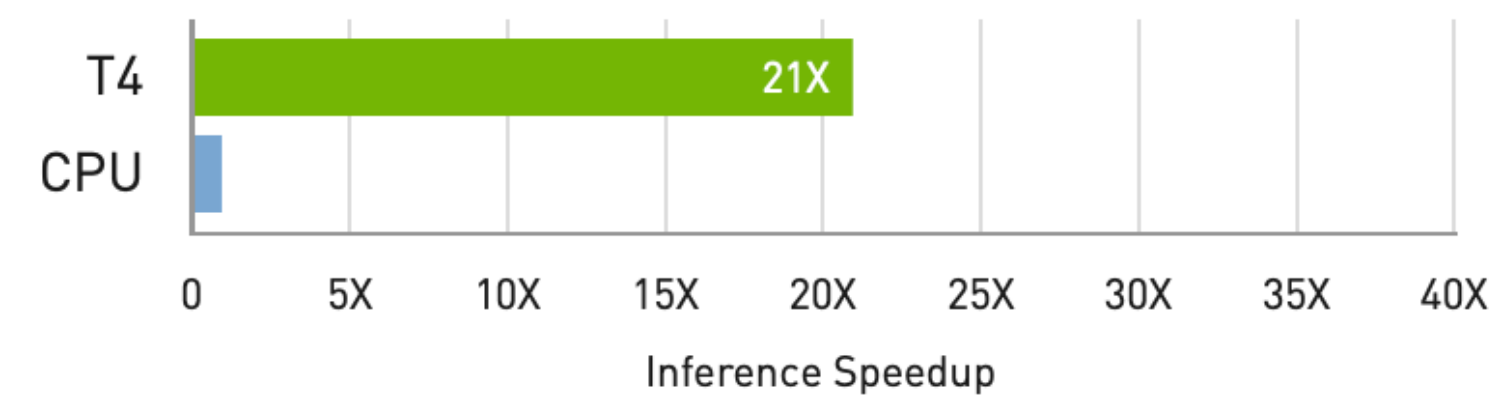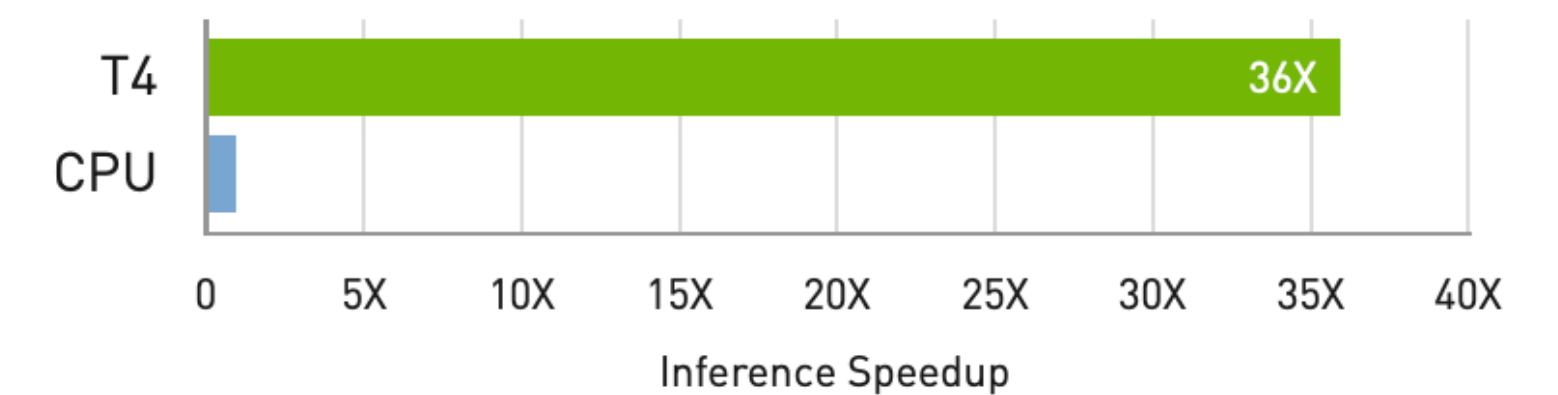# FROM FP32 TO FP16 TO INT8, AS WELL AS INT4



**Resnet50**

| | | |
|---|---|---|
| T4 | | 27X |
| CPU | | |

0  5X  10X  15X  20X  25X  30X  35X  40X
Inference Speedup

**DeepSpeech2**

| | | |
|---|---|---|
| T4 | | 21X |
| CPU | | |

0  5X  10X  15X  20X  25X  30X  35X  40X
Inference Speedup

**GNMT**

| | | |
|---|---|---|
| T4 | | 36X |
| CPU | | |

0  5X  10X  15X  20X  25X  30X  35X  40X
Inference Speedup

# Ternary Neural Networks for Resource-Efficient AI Applications

Hande Alemdar*, Vincent Leroy*, Adrien Prost-Boucle†, Frédéric Pétrot†

*Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG, F-38000 Grenoble, France
†Univ. Grenoble Alpes, CNRS, Grenoble INP, TIMA, F-38000 Grenoble, France
Email: name.surname@univ-grenoble-alpes.fr

*Abstract*—**The computation and storage requirements for Deep Neural Networks (DNNs) are usually high. This issue limits their deployability on ubiquitous computing devices such as smart phones, wearables and autonomous drones. In this paper, we propose ternary neural networks (TNNs) in order to make deep learning more resource-efficient. We train these TNNs using a teacher-student approach based on a novel, layer-wise greedy methodology. Thanks to our two-stage training procedure, the teacher network is still able to use state-of-the-art methods such as dropout and batch normalization to increase accuracy and reduce training time. Using only ternary weights and activations, the student ternary network learns to mimic the behavior of its teacher network without using any multiplication. Unlike its {-1,1} binary counterparts, a ternary neural network inherently prunes the smaller weights by setting them to zero during training. This makes them sparser and thus more energy-efficient. We design a purpose-built hardware architecture for TNNs and implement it on FPGA and ASIC. We evaluate TNNs on several benchmark datasets and demonstrate up to 3.1× better energy efficiency with respect to the state of the art while also improving accuracy.**

eliminating the need for multiplications [10], [11], [12]. The main drawback of these approaches is a significant degradation in the classification accuracy in return for a limited gain in resource efficiency.

This paper introduces ternary neural networks (TNNs) to address these issues and makes the following contributions:

- We propose a teacher-student approach for obtaining Ternary NNs with weights and activations constrained to $\{-1, 0, 1\}$. The teacher network is trained with stochastic firing using back-propagation, and can benefit from all techniques that exist in the literature such as dropout [13], batch normalization [14], and convolutions, The student network has the same architecture and, for each neuron, mimics the behavior of the equivalent neuron in the teacher network without using any multiplications,
- We design a specialized hardware that is able to process TNNs at up to 2.7× better throughput, 3.1× better energy

without extensive data augmentation. TNN's error rate on MNIST is 1.67% with a single 3-layer MLP with 750 neurons in each layer. Bitwise NNs [10] with 1024 neurons in 3 layers achieves a slightly better performance. TNN with an architecture that has similar size to Bitwise NN is worse due to over-fitting. Since TNN selects a different sparsity level for each neuron, it can perform better on smaller networks, and larger networks cause over-fitting on MNIST. Bitwise NN's global sparsity parameter has a better regularization effect on MNIST

### TABLE IV
### CLASSIFICATION PERFORMANCE - ERROR RATES (%)

| | MNIST | CIFAR10 | SVHN | GTRSB | CIFAR100 |
|---|---|---|---|---|---|
| **Fully Discretized** | | | | | |
| **TNN** (This Work) | 1.67 | 12.11 | 2.73 | 0.98 | 48.40 |
| TrueNorth [11], [12] | 7.30 | 16.59 | 3.34 | 3.50 | 44.36 |
| Bitwise NN [10] | 1.33 | | | | |
| **Partially Discretized** | | | | | |
| Binarized NN [6] | 0.96 | 10.15 | 2.53 | | |
| BC [15] | 1.29 | 9.90 | 2.30 | | |
| TC [9] | 1.15 | 12.01 | 2.42 | | |
| TWN [8] | 0.65 | 7.44 | | | |
| EBP [24] | 2.20 | | | | |
| XNOR-Net [16] | | 9.88 | | | |
| DoReFa-Net [7] | | | 2.40 | | |

### A. Hardware Architecture

Figure 3 outlines the hardware architecture of a fully-connected layer in a multi-layer NN. The design forms a pipeline that corresponds to the sequence of NN processing steps. For efficiency reasons, the number of layers and the maximum layer dimensions (input size and number of neurons) are decided at synthesis time. For a given NN architecture, the design is still user-programmable: each NN layer contains a memory that can be programmed at run-time with neuron weights or output ternarization thresholds $b^{lo}$ and $b^{hi}$. As seen in the previous experiments of Section IV, a given NN architecture can be reused for different datasets with success.

Ternary values are represented with 2 bits using usual two's complement encoding. That way, the compute part of each neuron is reduced to just one integer adder/subtractor and one
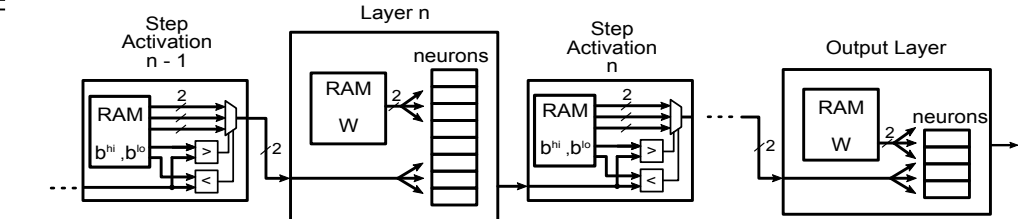


Fig. 3. Hardware implementation scheme of ternary neural network

### TABLE I
### TERNARY NEURAL NETWORK DEFINITIONS FOR A SINGLE NEURON $i$

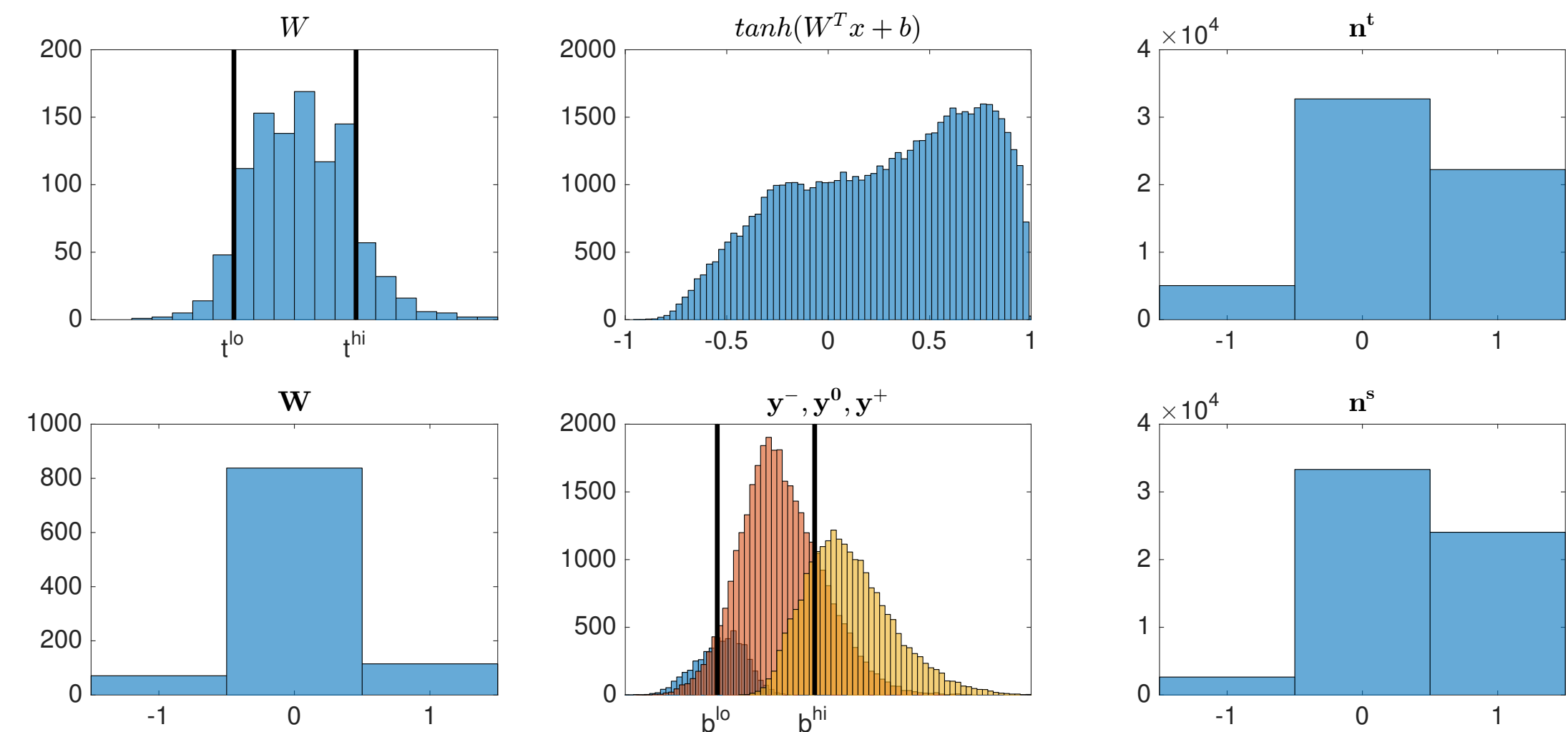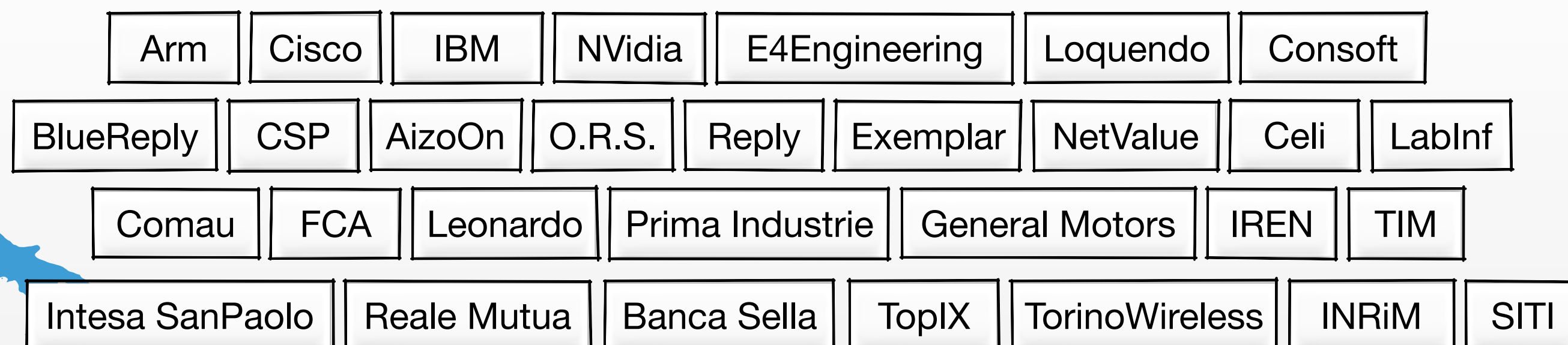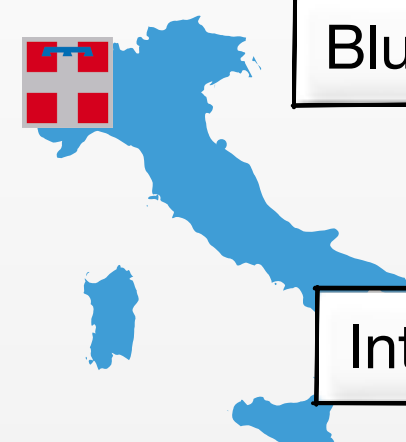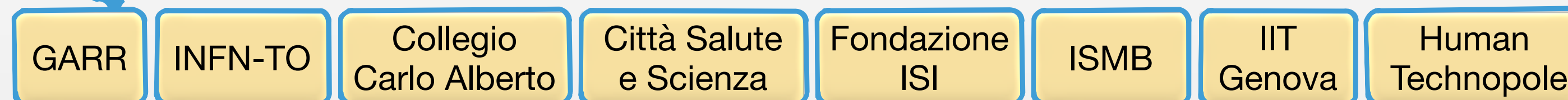| | Teacher network | Student Network |
|---|---|---|
| Weights | $W_i = [w_j], w_j \in \mathbb{R}$ | $\mathbf{W_i} = [\mathbf{w_j}], \mathbf{w_j} \in \{-1, 0, 1\}$ |
| Bias | $b_i \in \mathbb{R}$ | $\mathbf{b_i}^{lo} \in \mathbb{Z}$ $\mathbf{b_i}^{hi} \in \mathbb{Z}$ |
| Transfer Function | $y_i = W_i^\mathsf{T} \mathbf{x} + b_i$ | $\mathbf{y_i} = \mathbf{W_i}^\mathsf{T} \mathbf{x}$ |
| Act. Fun | $\mathbf{n_i^t} = \begin{cases} -1 & \text{with prob. } -\rho \text{ if } \rho < 0 \\ 1 & \text{with prob. } \rho \text{ if } \rho > 0 \\ 0 & \text{otherwise} \end{cases}$ where $\rho = tanh(y_i), \rho \in (-1, 1)$ | $\mathbf{n_i^s} = \begin{cases} -1 & \text{if } \mathbf{y_i} < \mathbf{b_i}^{lo} \\ 1 & \text{if } \mathbf{y_i} > \mathbf{b_i}^{hi} \\ 0 & \text{otherwise} \end{cases}$ |



Fig. 1. Example ternarization for a single neuron

- Training

  - Compute: very demanding

  - Network: asynchronous, associative, commutative → Not demanding

  - Precision requirement: single or even less

- Inference

  - Compute: not demanding quantisable

  - Embarrassingly parallel → Not demanding

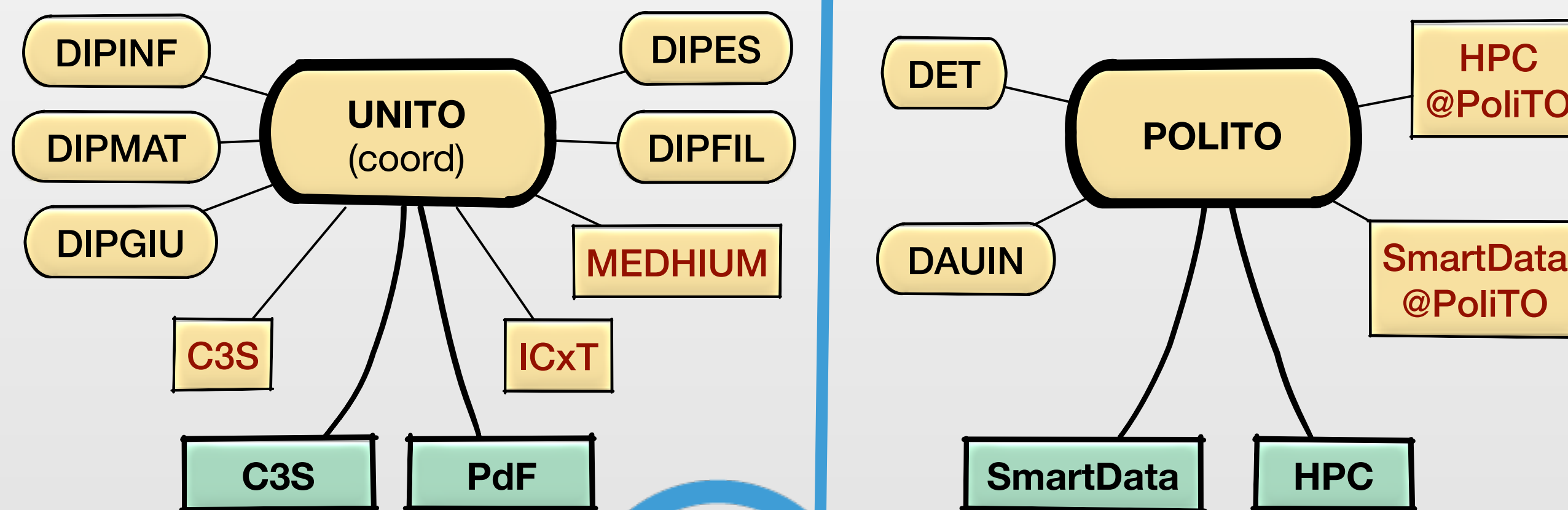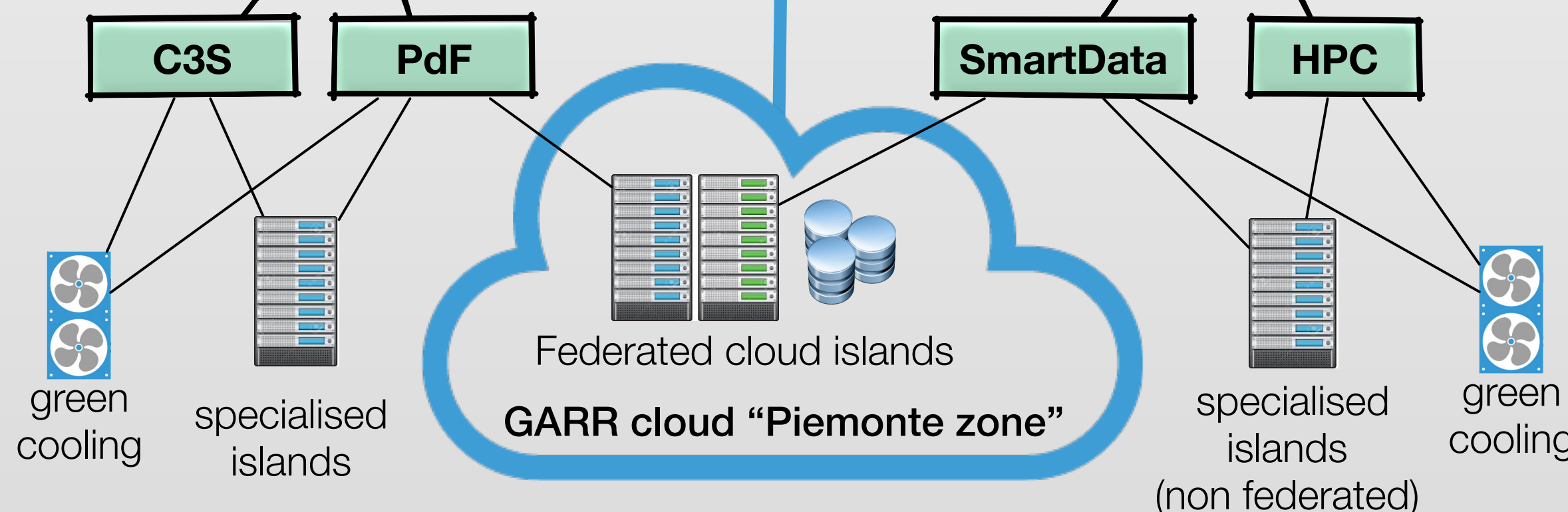  - Precision requirement:  quantisable up to 1 bit with decent performances

# HPC4AI

## AI-on-demand platform
http://www.hpc4ai.it

## Facts

- INFRA-P call Nov. 2017
- Ranked 1st on ~30 submitted projects
- Kick-off mid apr 2018
- 4.5M€ funding
- 2 partners
- 8 associated partners
- Coord. M. Aldinucci
- Many industrial stakeholders

| Users | Kind of service | Services | Artifacts |
|---|---|---|---|
| **Domain experts with no skills** on ML and BDA.<br><br>Training set not required. Off-the-shelf algorithms/networks. | **Service-as-a-Service (SaaS)** | SaaS for ML and BDA designed within HPC4AI partners | Market place for ML and BDA services: Dashboards, trained models in several domains (NLP, Vision, …) |
| **Domain experts skilled** on ML and BDA. **Not expert in parallel computing.**<br><br>New networks or pipelines; training set required. | **Platform-as-a-Service (PaaS)** | PaaS solutions for ML and BDA directly designed within HPC4AI or companion projects | Market place of VMs and Platforms realising software stacks for ML and BDA. Solutions for data ingestion, data lake, etc. |
| Researchers, cloud engineering, ML and BDA framework designers, cloud engineers, stack and automation designers. | **1) Infrastructure-as-a-Service (IaaS)**<br><br>**2) Metal-as-a-Service (IaaS)** | 1) GARR/other cloud able to support federation<br><br><br><br>2) Job scheduler for HPC resources | 1) Openstack, docker, VM, object storage, file storage, kubernetes, etc.<br><br>2) Alternative cloud, job queue, Big Data Stack (Spark, …). |
| Researchers, run-time designers. | **Hardware** | Bare Metal | Multicore, GPU, storage, network, switch, UPS, cooling, etc. |

IaaS

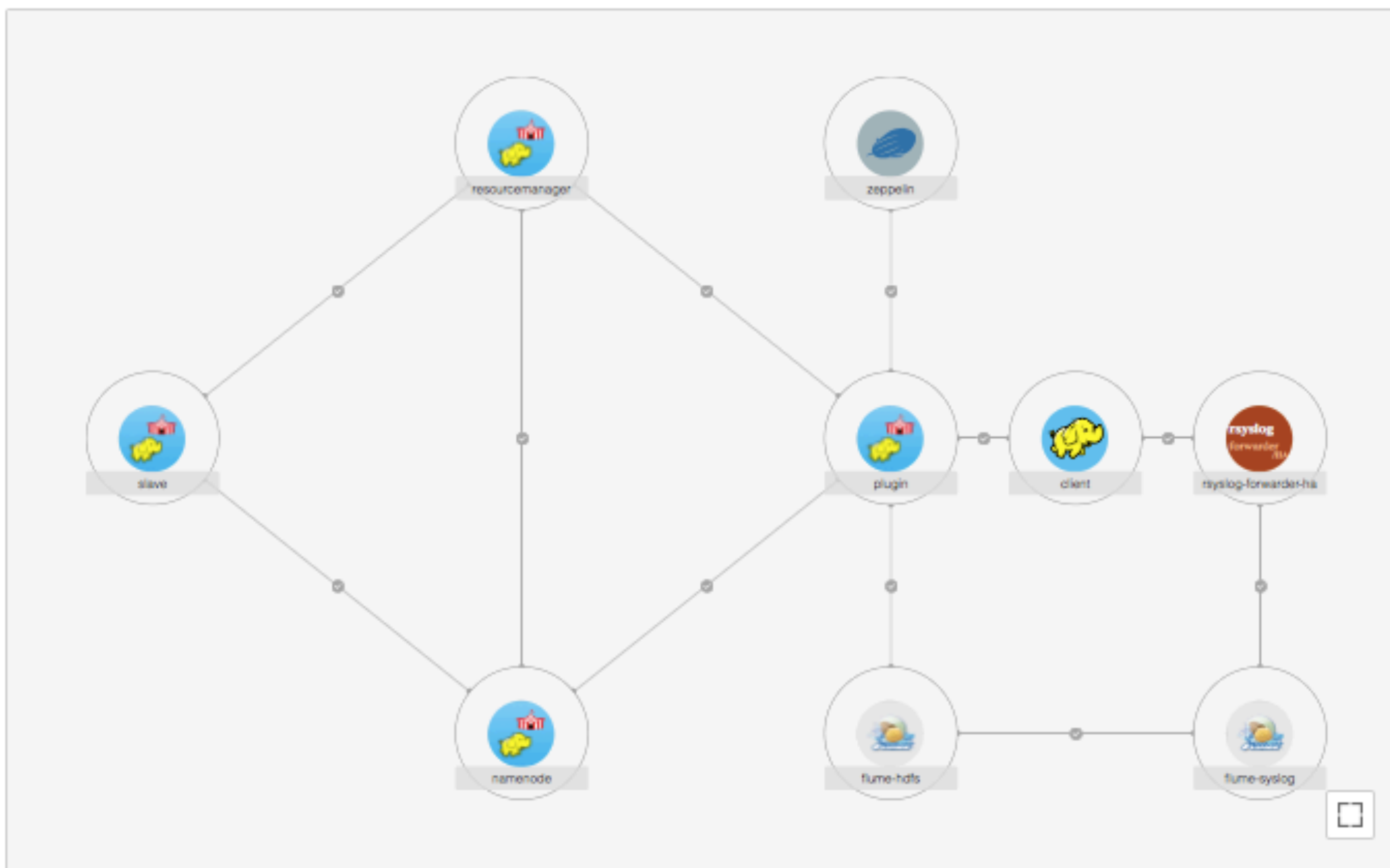PaaS

IaaS

Realtime Syslog Analytics #22

By **bigdata-charmers** • xenial • Stable

This is a six unit big data cluster that includes Hadoop 2.7.3 and other components from Apache Bigtop. By leveraging Rsyslog and Apache Flume, this bundle provides an environment for analysing syslog events in Apache Zeppelin web notebooks.

Add to model

Submit a bug

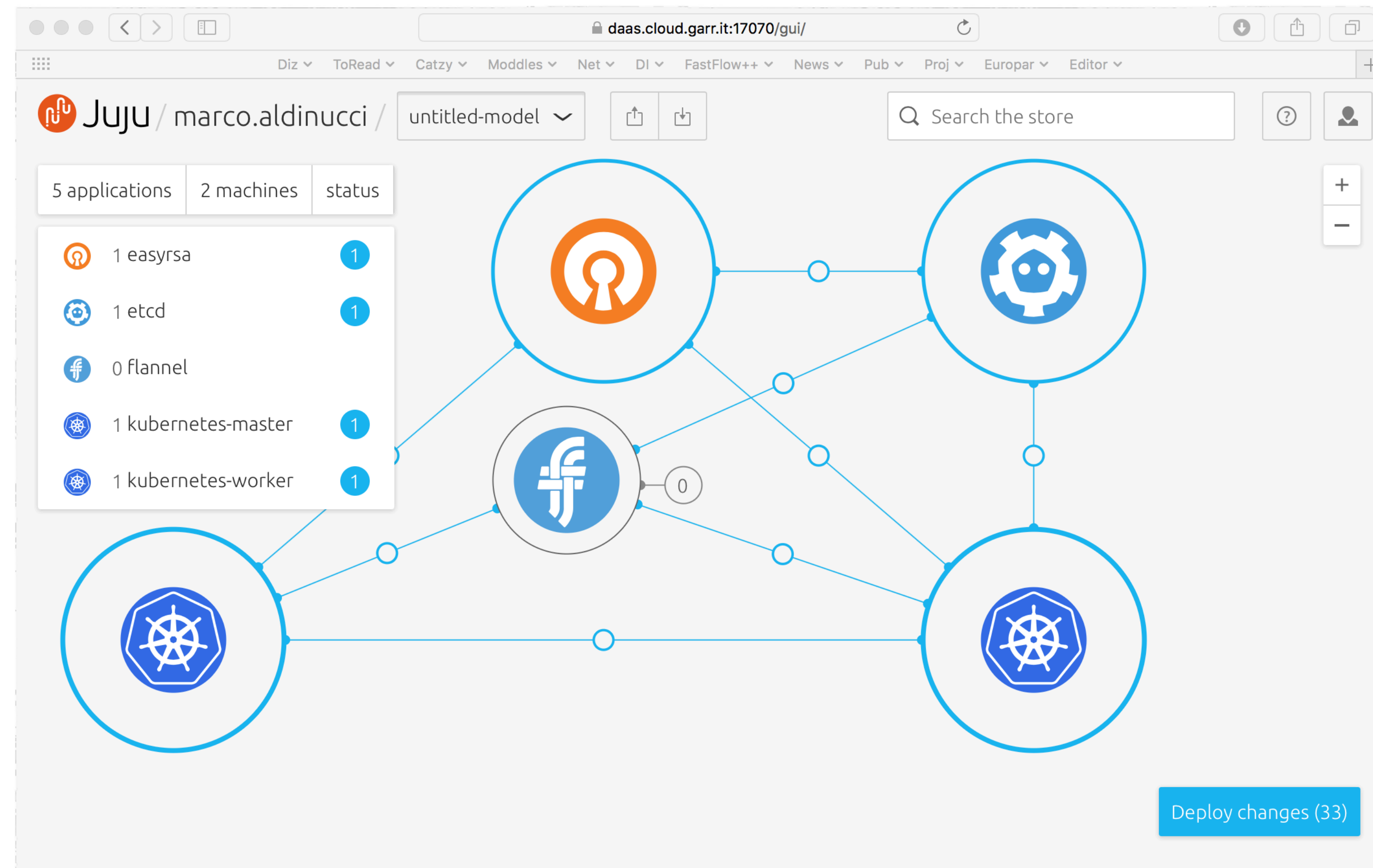Project homepage

Files

README.md

bundle.yaml

copyright

/tests

Download .zip

Embed this charm

Add this card to your website by copying the code below. Learn more.

resourcemanager

zeppelin

slave

plugin

client

rsyslog-forwarder-ha

namenode

flume-hdfs

flume-syslog

Project

Compute

Overview

Instances

Volumes

Images

Access & Security

Network

Object Store

Identity

Overview

# EXAMPLE PAAS ON BARE METAL: KUBERNETES-AS-A-SERVICE



*M. Aldinucci et al. "HPC4AI, an AI-on-demand federated platform endeavour," in ACM Computing Frontiers, Ischia, Italy, 2018.*

# EXAMPLE PAAS ON VMS: BIGDATA-AS-A-SERVICE



M. Aldinucci et al. "HPC4AI, an AI-on-demand federated platform endeavour," in ACM Computing Frontiers, Ischia, Italy, 2018.

# WHAT'S NEXT?



Model Zoo
Discover open source deep learning code and pretrained models.

Browse Frameworks    Browse Categories

**OpenPose** ⭐ 11120

OpenPose represents the first real-time multi-person system to jointly detect human body, hand, and facial keypoints (in total 130 keypoints) on single images.

Caffe
CV

**Mask R-CNN** ⭐ 10157

This is an implementation of Mask R-CNN on Python 3, Keras, and TensorFlow. The model generates bounding boxes and segmentation masks for each instance of an object in the image. It's based on Feature Pyramid Network (FPN) and a ResNet101 backbone.

Keras
CV

**FastPhotoStyle** ⭐ 9228

A Closed-form Solution to Photorealistic Image Stylization

PyTorch
CV

**Image-to-Image Translation with Conditional Adversarial Networks** ⭐ 6843

This is our PyTorch implementation for both unpaired and paired image-to-image translation

**Subscribe**

We send out monthly emails showcasing the best or most notable models released each month. Sign up to receive updates!

✉ youremail@domain.com

Subscribe

**vid2vid** ⭐ 5908

Pytorch implementation of our method for high-resolution (e.g. 2048x1024) photorealistic video-to-video translation.

PyTorch
NLP



*A mostly complete chart of*
# Neural Networks
©2016 Fjodor van Veen – asimovinstitute.org

25

# Follow-up project (2018_ICT-11-a EU IA, 13M€)
## DeepHealth: Deep-Learning and HPC to Boost Biomedical Applications for Health

- 18 Parterns: Everis, Siveco, Wings, Philiphs, SIVECO, IBM, Thales, CEA, Treelogic, EPFL, UPV, UNITO, UNIMORE, ...

- Design and develop:
  - European Library for Distributed Deep Learning for health
  - AI-on-demand cloud platforms (HPC4AI, ...)

# From HPC to BigData to Deep Learning

- ## HPC
  - Send/recv, batch, CPU intensive
  - Programmer should have the direct knowledge of all processes and all communicat
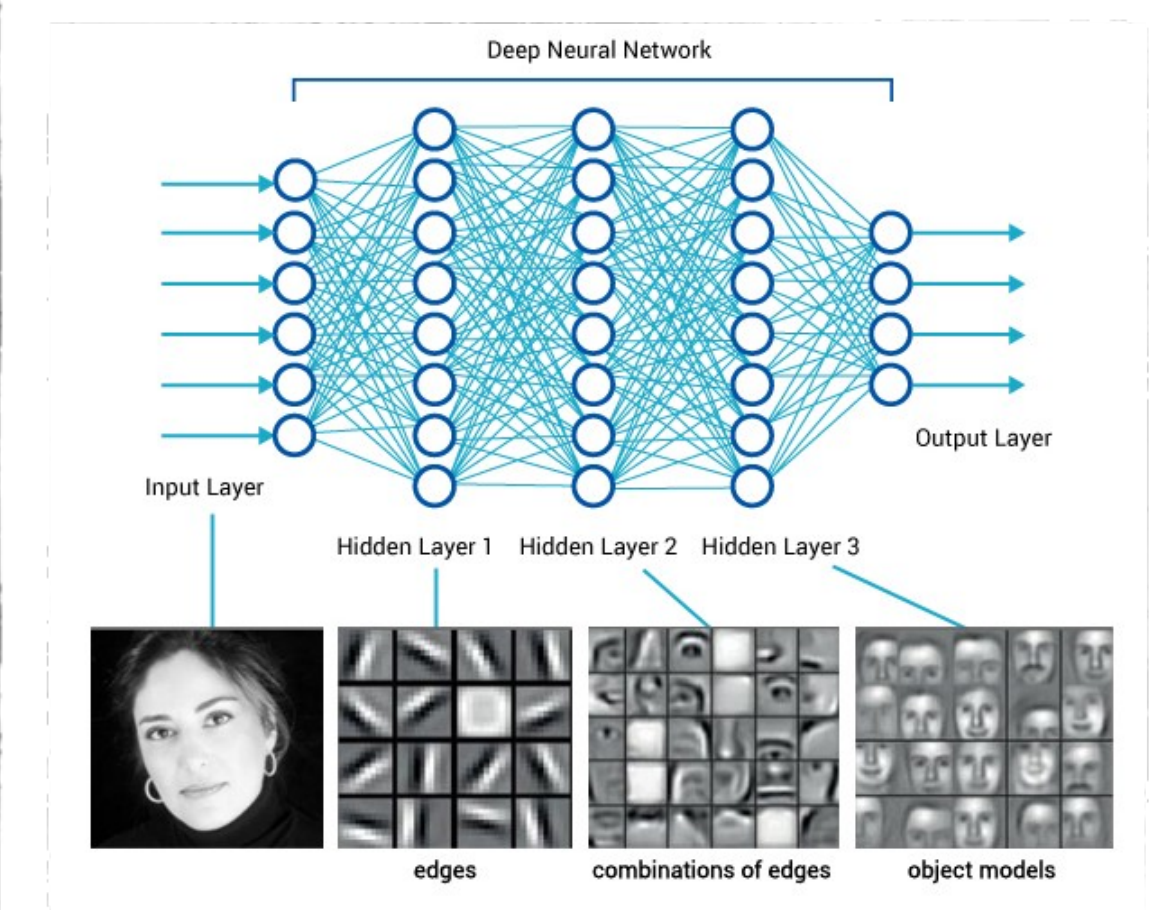- ## Exascale HPC
  - Data movements rather than compute power
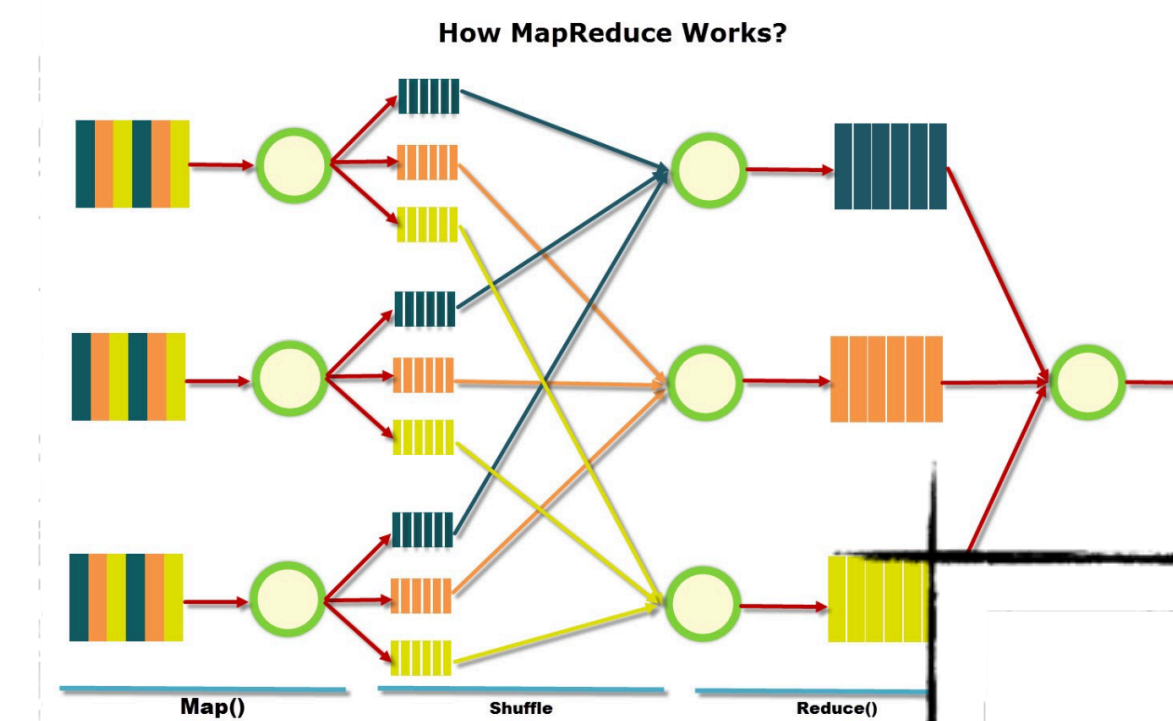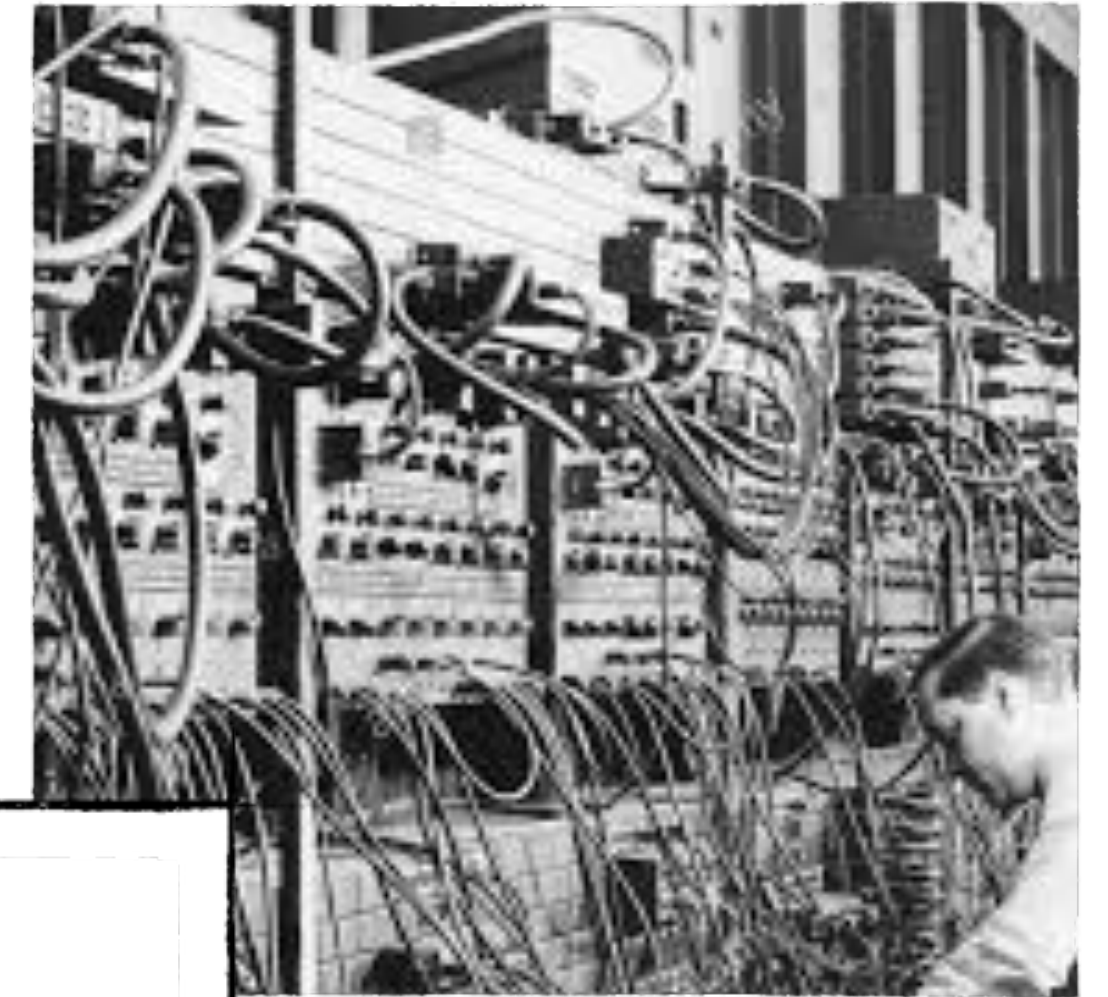  - Beyond locally synchronous data parallelism? Tasks?
- ## BigData (more abstract)
  - Mostly streams & I/O: intensive, interactive
  - Virtualised, cloud stack, Platform/Service-as-a-Service, service composition
  - O(n) algoritms - no more!
- ## Deep Learning (much more abstract)
  - Naturally asynchronous, permissive, low precision
  - Both batch (training) and stream (inference)

# PROBLEMS SOLVED (PROGRAMMING MODEL)

- ## HPC: large-scale locally synchronous (stencil)

  - Low-level: MPI

  - High-level: open problem. Tasks?

- ## Big Data analytics

  - MapReduce programming model (data flow run-time)

- ## Deep learning

  - Training: API and architecture for **scale-up** well understood (GPU/Tensor)

  - Inference: Quantisation + Prcessing-in-Memory new unfolding

*M. Aldinucci, M. Drocco, C. Misale, and G. Tremblay, "Languages for Big Data analysis," in Encyclopedia of Big Data Technologies, Springer, 2019*

- ## DL: Distributed training at scale

  - ### Today: lock-step all-to-all weight exchange, i.e. globally synchronous → not scalable

**Dataset**

Partitions

Gradients

Gradients

Neighbors

time

$\mathbf{w}_0$

Worker A

Worker B

$\delta?$

$\mathbf{w}_1^A = \mathbf{w}_0 + \delta_1^A(\mathbf{w}_0)$

$\mathbf{w}_1^B = \mathbf{w}_0 + \delta_1^B(\mathbf{w}_0)$

Not received

Received

Not received

Received

$\mathbf{w}_1^A + \delta_2^A(\mathbf{w}_1^A)$

$\mathbf{w}_1^A + \delta_1^B(\mathbf{w}_0)$
$+\delta_2^A(\mathbf{w}_1^A + \delta_1^B(\mathbf{w}(0))$

$\mathbf{w}_1^B + \delta_2^B(\mathbf{w}_1^B)$

$\mathbf{w}_1^B + \delta_1^A(\mathbf{w}_0)$
$+\delta_2^B(\mathbf{w}_1^B + \delta_1^A(\mathbf{w}_0))$

*P. Viviani, M. Drocco, D. Baccega, and M. Aldinucci, "Deep Learning at Scale," PDP 2019*
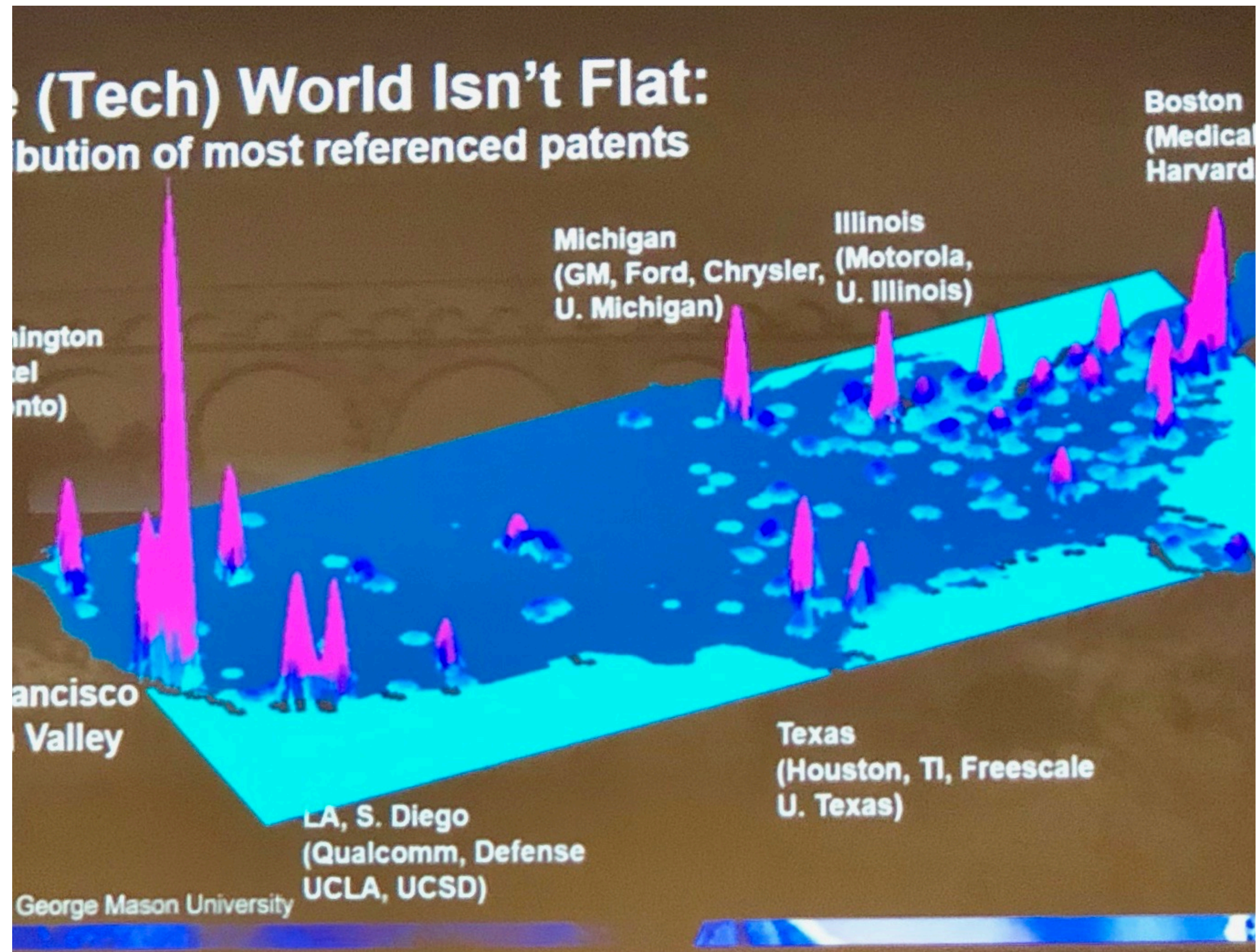
# OPEN PROBLEMS

- ## General methodology

  - Most of the paper "observe" network x on data y with tuning z

- ## Distributed memory hierarchy

  - Data movement prevails computation

- ## Privacy and data marketplace

  - Multi-party computation

  - Federated learning

  - Inference: a complete programming ecosystem

# FUTURE HORIZONS (PROGRAMMING MODEL)

- ## Specialisation: accelerators are the key

  - Deep Learning

  - Quantum computing

  - Neuromorphic

- ## Programming models: methodology is the key

  - Should be simple and compositional as sequential or web services coding

  - Performance is not the only issue: time-to-market, cost, correctness, …

# Future horizons