# Real-time AI for Science

Mia Liu

COFI Advanced Instrumentation and Data analysis School

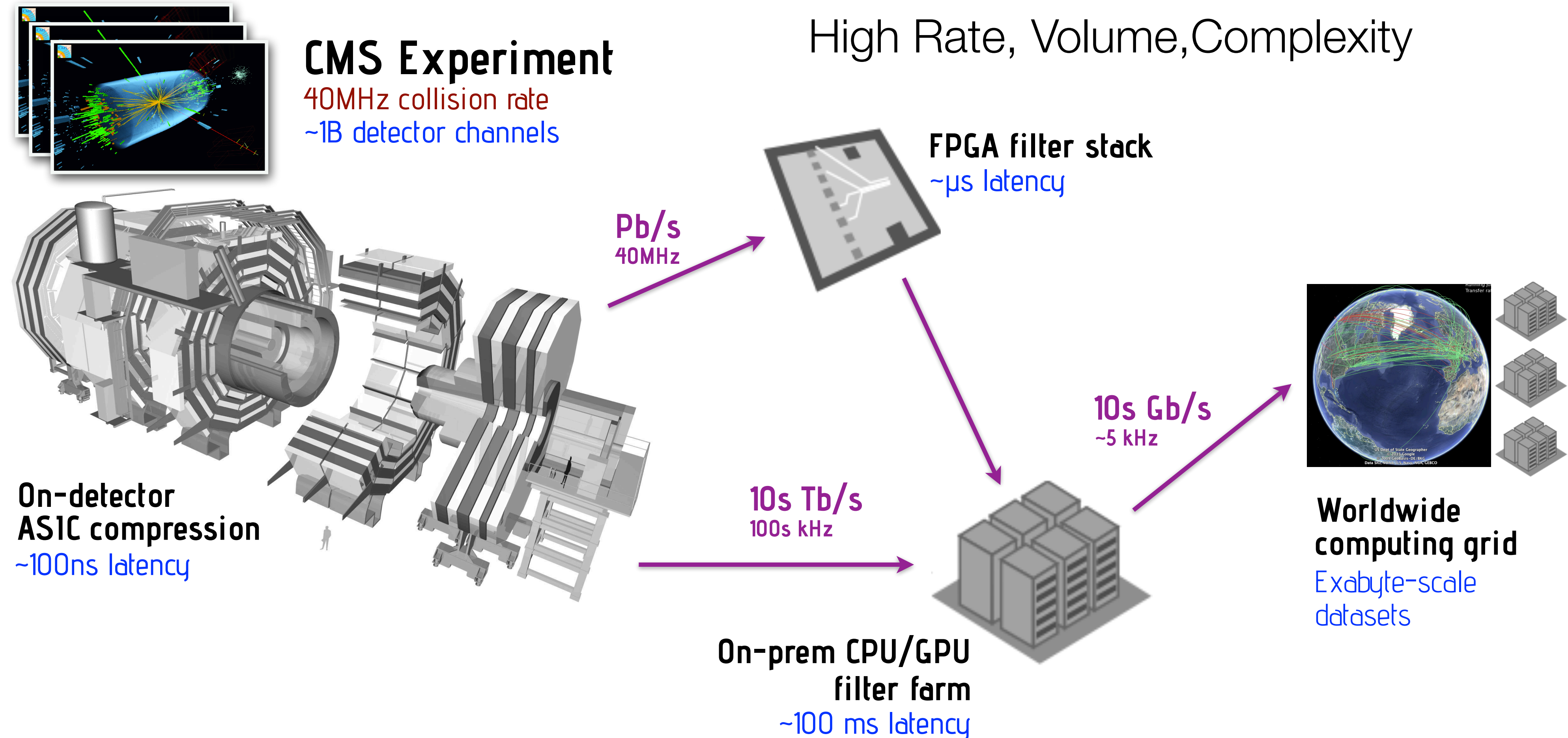Dec.2023

# The Large Hadron Collider

# The CMS detector
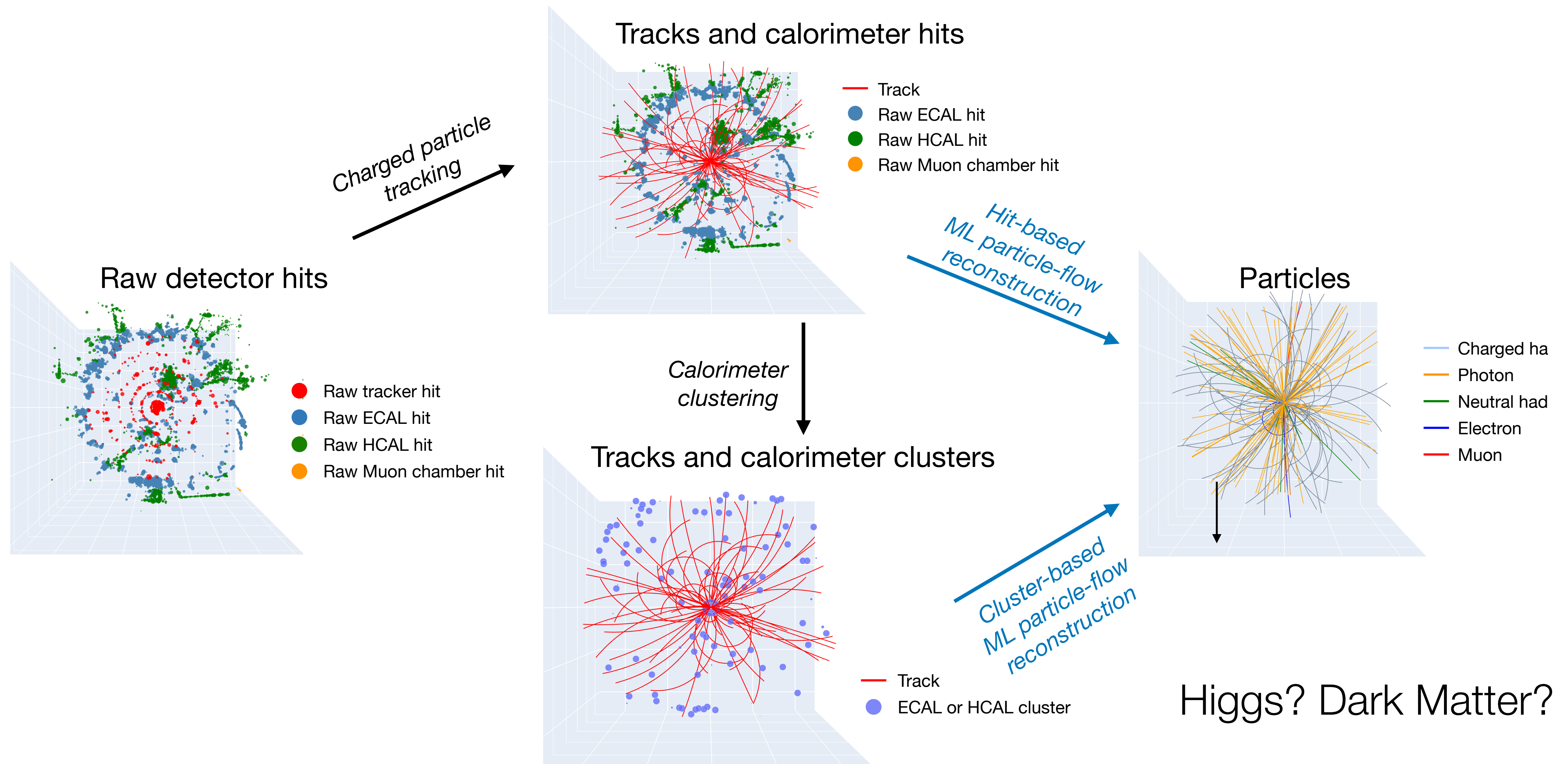
# Today's Lecture

- Real-time system constraints and needs

  - LHC as an example

  - Specialized hardwares: FPGA/ASIC

- Challenges in AI/ML on specialized hardwares

  - Design efficient networks for real-time systems

  - Co-design tools and needs

- Real-time AI for other science domains: quantum system control etc
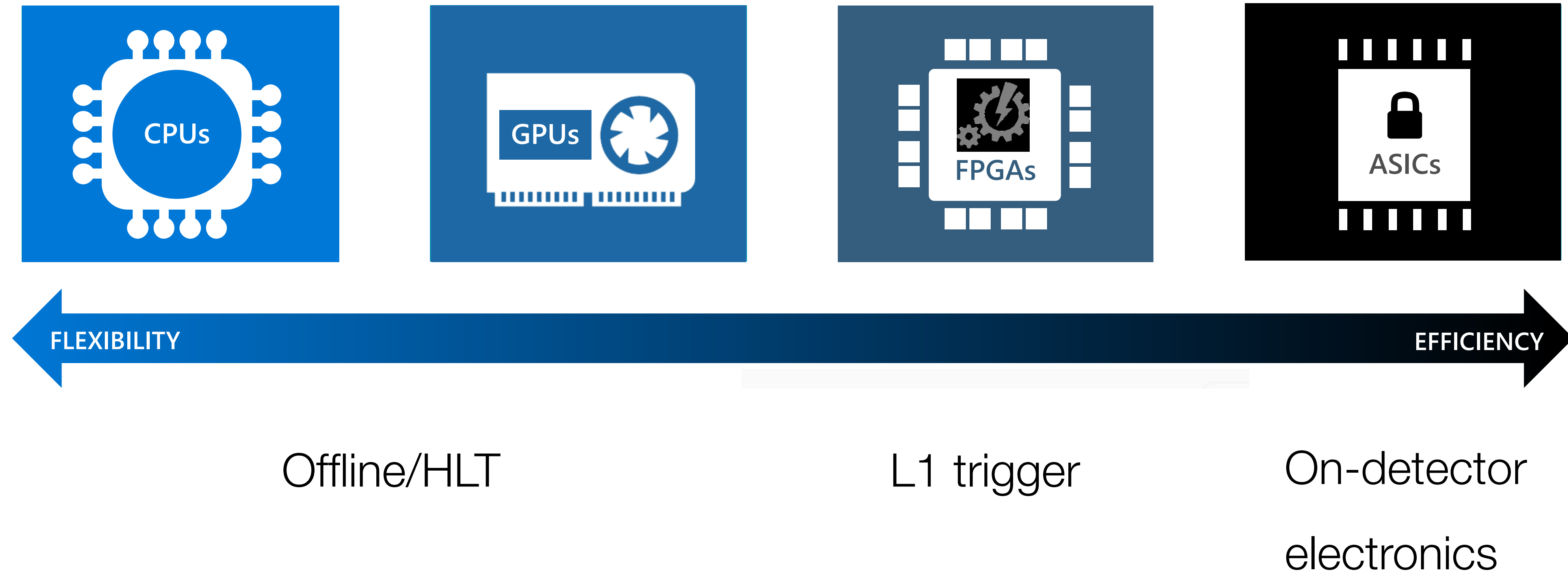
# From Collisions to Discoveries

**Fermilab**

**CMS Experiment**
40MHz collision rate
~1B detector channels

High Rate, Volume, Complexity

**FPGA filter stack**
~µs latency

**Pb/s**
40MHz

**On-detector
ASIC compression**
~100ns latency

**10s Tb/s**
100s kHz

**10s Gb/s**
~5 kHz

**Worldwide
computing grid**
Exabyte-scale
datasets

**On-prem CPU/GPU
filter farm**
~100 ms latency

Science with Big data: Multi-tier Data Processing

# AI/ML: new opportunities for real-time reconstruction



Tracks and calorimeter hits

— Track
● Raw ECAL hit
● Raw HCAL hit
● Raw Muon chamber hit

Charged particle tracking

Raw detector hits

● Raw tracker hit
● Raw ECAL hit
● Raw HCAL hit
● Raw Muon chamber hit

Hit-based ML particle-flow reconstruction

Calorimeter clustering

Tracks and calorimeter clusters

— Track
● ECAL or HCAL cluster

Cluster-based ML particle-flow reconstruction

Particles

— Charged ha
— Photon
— Neutral had
— Electron
— Muon

Higgs? Dark Matter?

Learning grouping of detector elements

6

# Hardware Landscape



CPUs    GPUs    FPGAs    ASICs

**FLEXIBILITY** ← → **EFFICIENCY**

Offline      L1 trigger      On-detector electronics

# What is an FPGA?

# What is an FPGA?

## Field Programmable Gate Arrays



Logic cell:
Flip-flops (FF) and
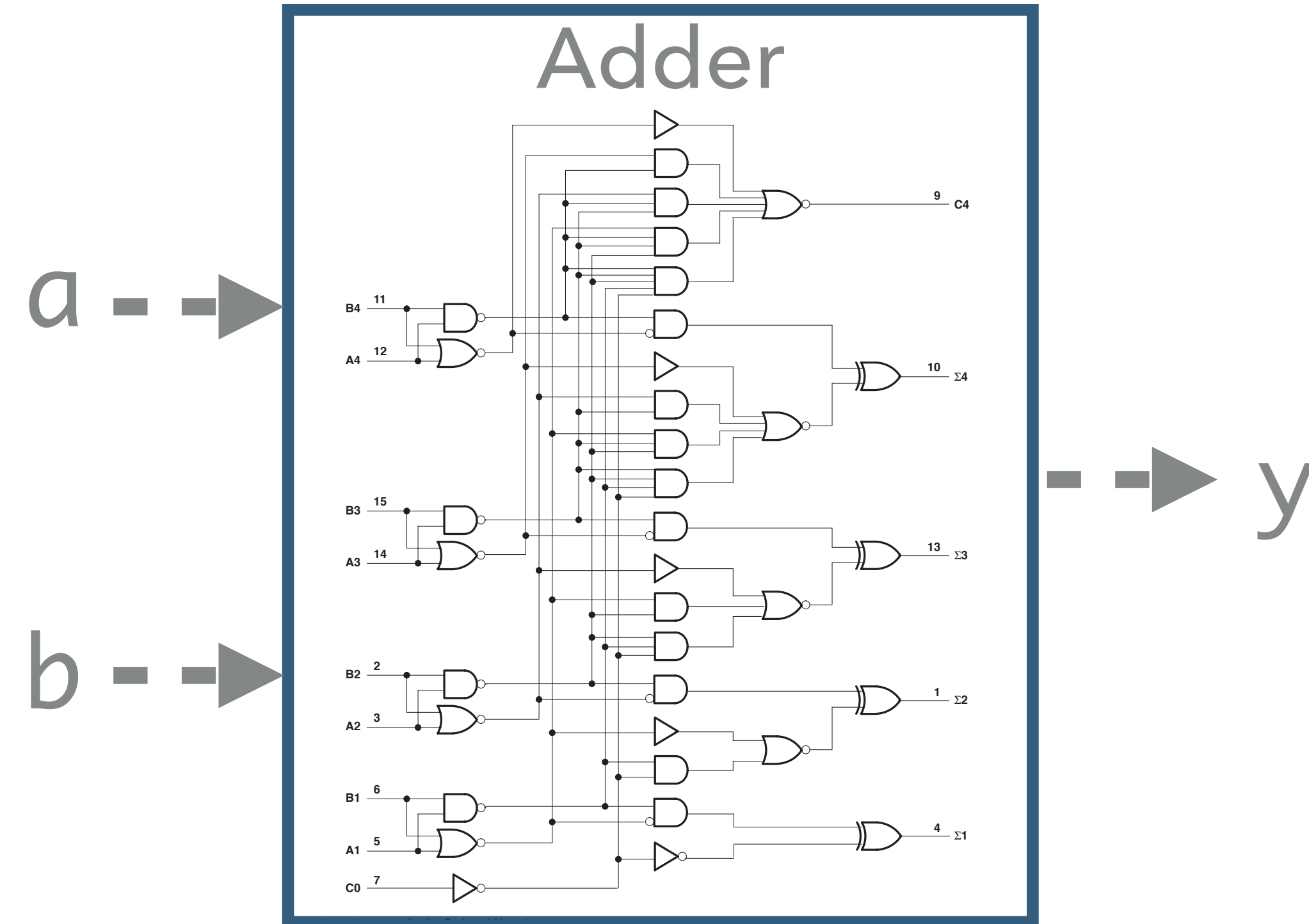look up tables (LUTs)

Digital Signal
Processors
(DSPs)

**Virtex Ultrascale+ VU9P**
6800 DSPs
1M LUTs
2M FFs
75 Mb BRAM

# Programming FPGAs

▸ Say you want to program an "adder" function on an FPGA

```
module adder(
    input  wire [4:0] a,
    input  wire [4:0] b,
    output wire [4:0] y
);
    assign y = a + b;

endmodule
```
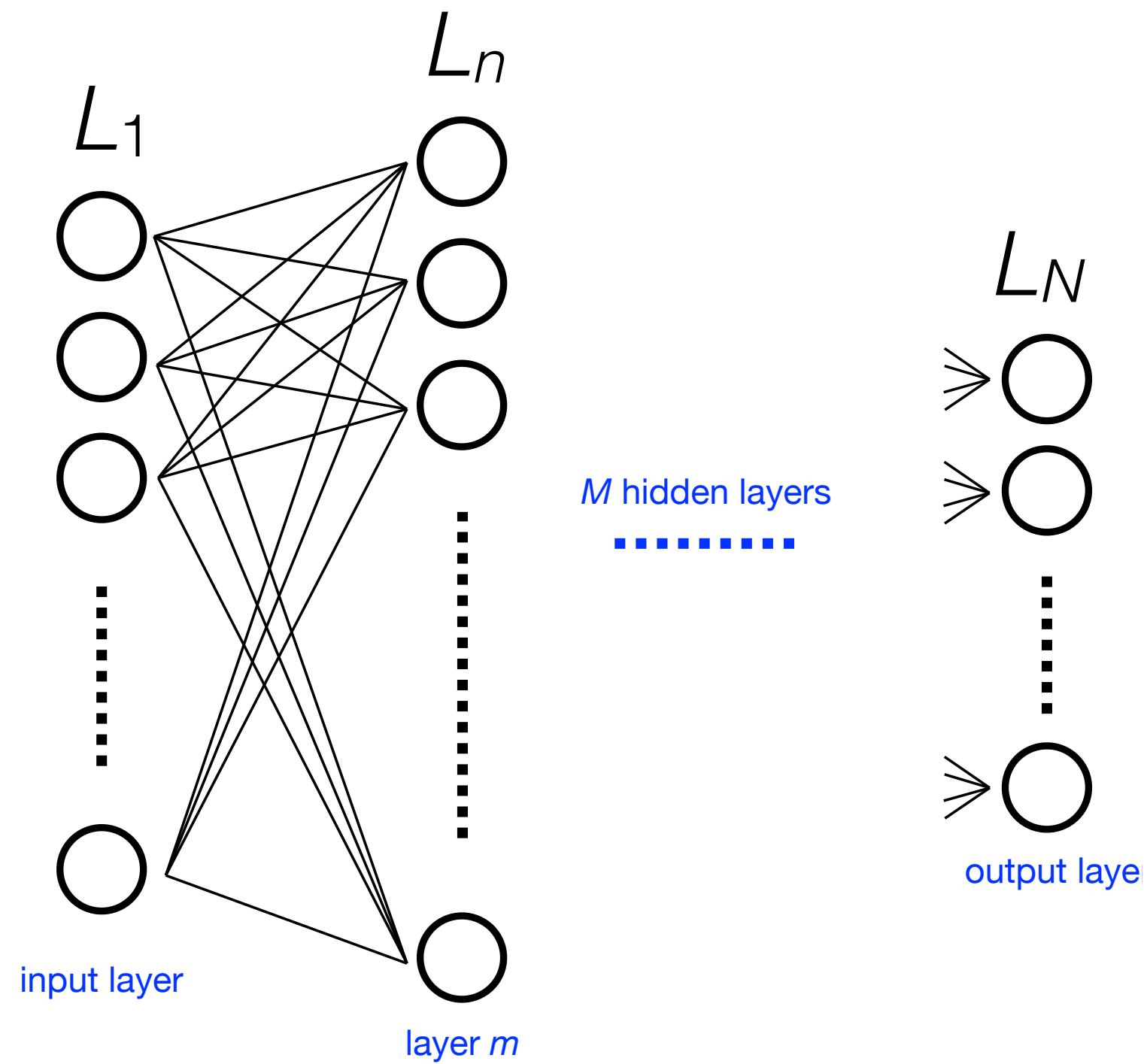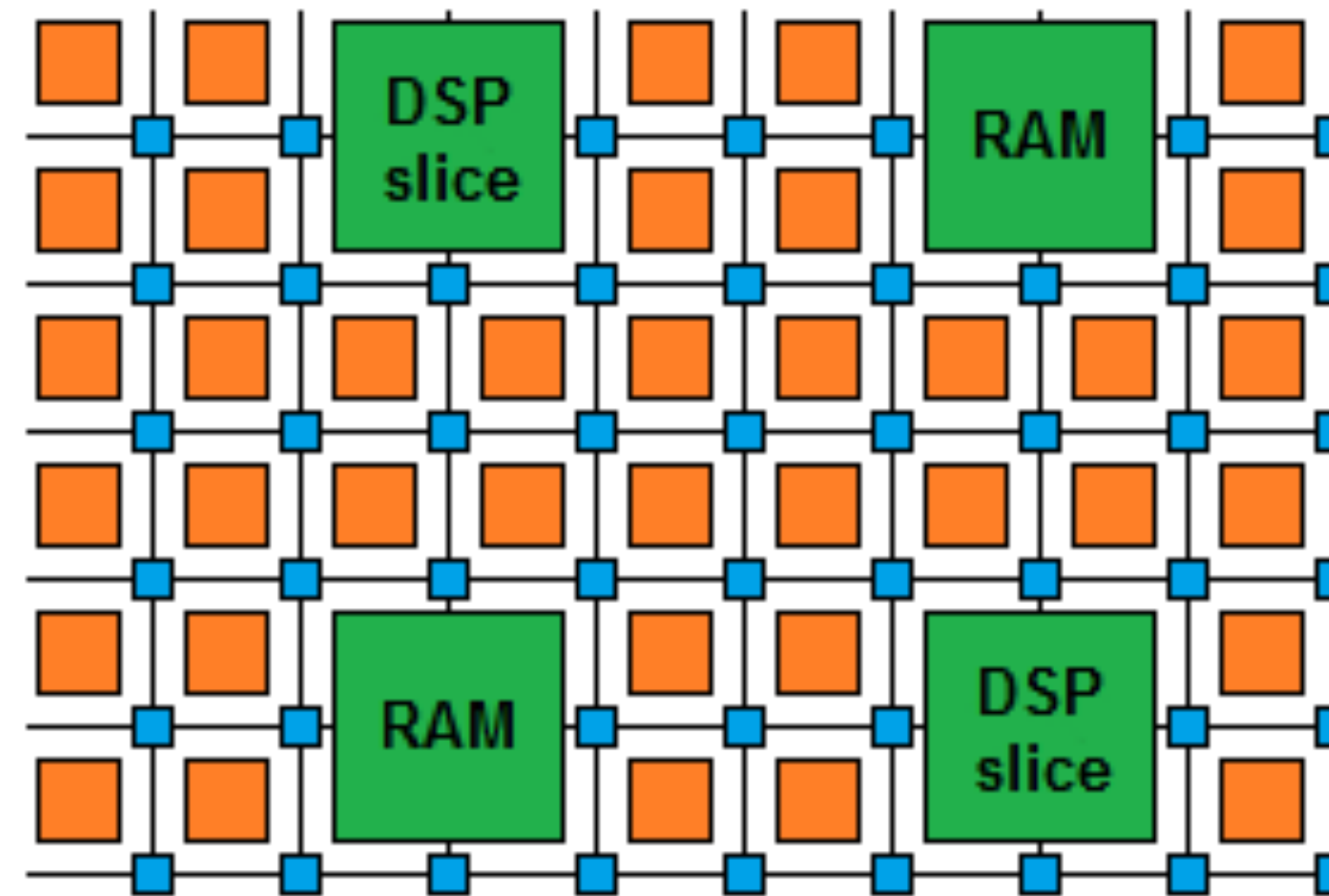
▸ Register transfer-level (RTL) code is "synthesized" into gates

Adder

a - - ▶

b - - ▶

- - ▶ y

Synthesis

**Logic gates**

AND    NAND    OR    NOR

NOT    XOR    XNOR

# Mapping NN onto FPGAs

## FPGA diagram



$$L_n$$

$$L_1$$

*M* hidden layers

$$n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1}$$

Logic cell:

Flip-flops (FF) and
look up tables (LUTs)

Digital Signal
Processors
(DSPs)

3 inputs

64 nodes
activation: ReLU

32 nodes
activation: ReLU

layer *m*

$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

Activation function    multiplication    addition

Activation
functions
precomputed and
stored in BRAMs

Precomputed, and
stored in BRAMs

Multiplications    Addition

DSPs    Logic cells

$$N_{\text{multiplications}} = \sum_{n=2}^{N} L_{n-1} \times L_n$$

**Virtex Ultrascale+ VU9P**
6800 DSPs
1M LUTs
2M FFs
75 Mb BRAM

16 inputs

64 nodes
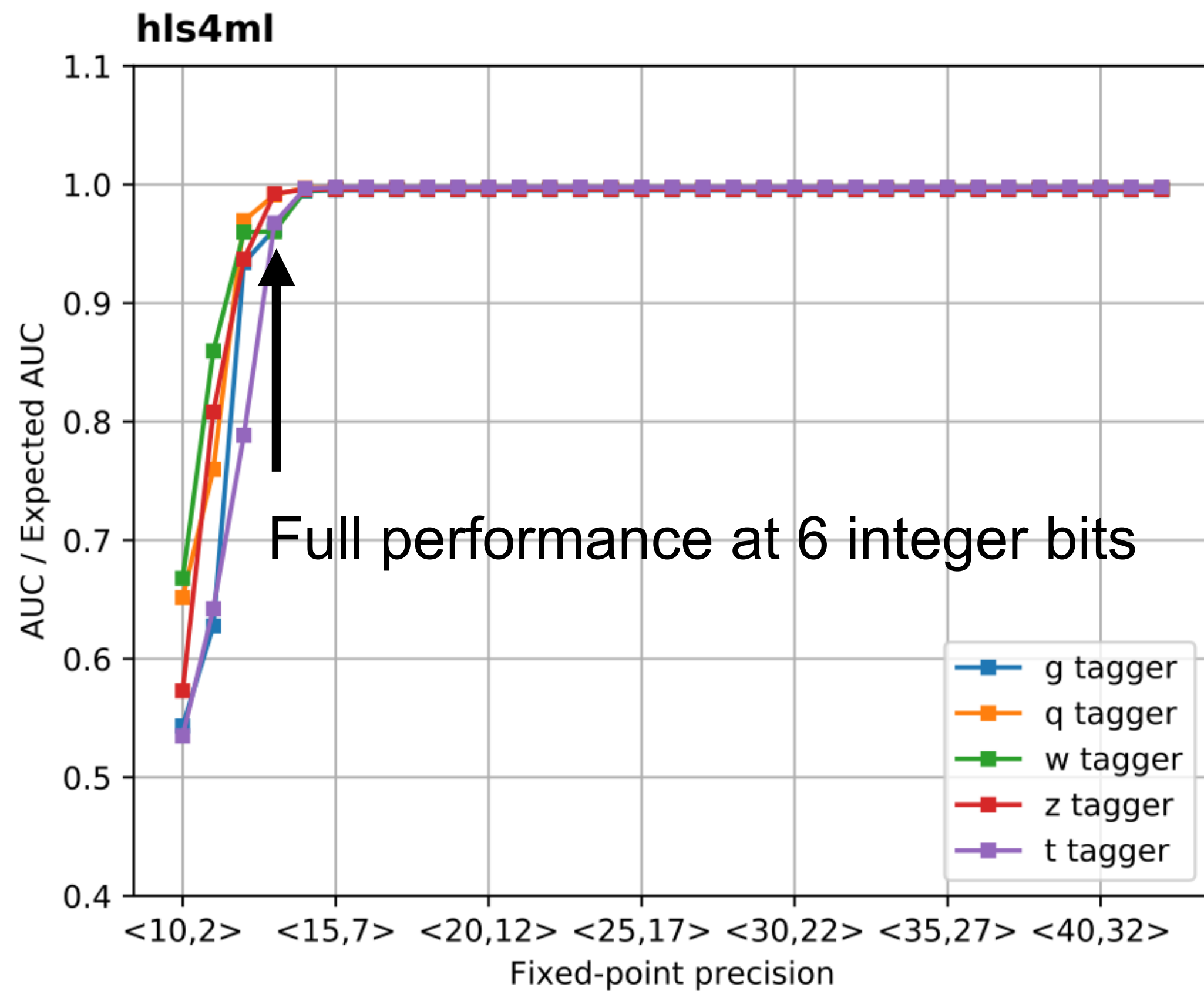
# Today's Lecture

- Real-time system constraints and needs

  - LHC as an example

  - Specialized hardwares: FPGA/ASIC

- **Challenges in AI/ML on specialized hardwares**

  - Design efficient networks for real-time systems

  - Co-design tools and needs

- Real-time AI for other science domains: quantum system control etc
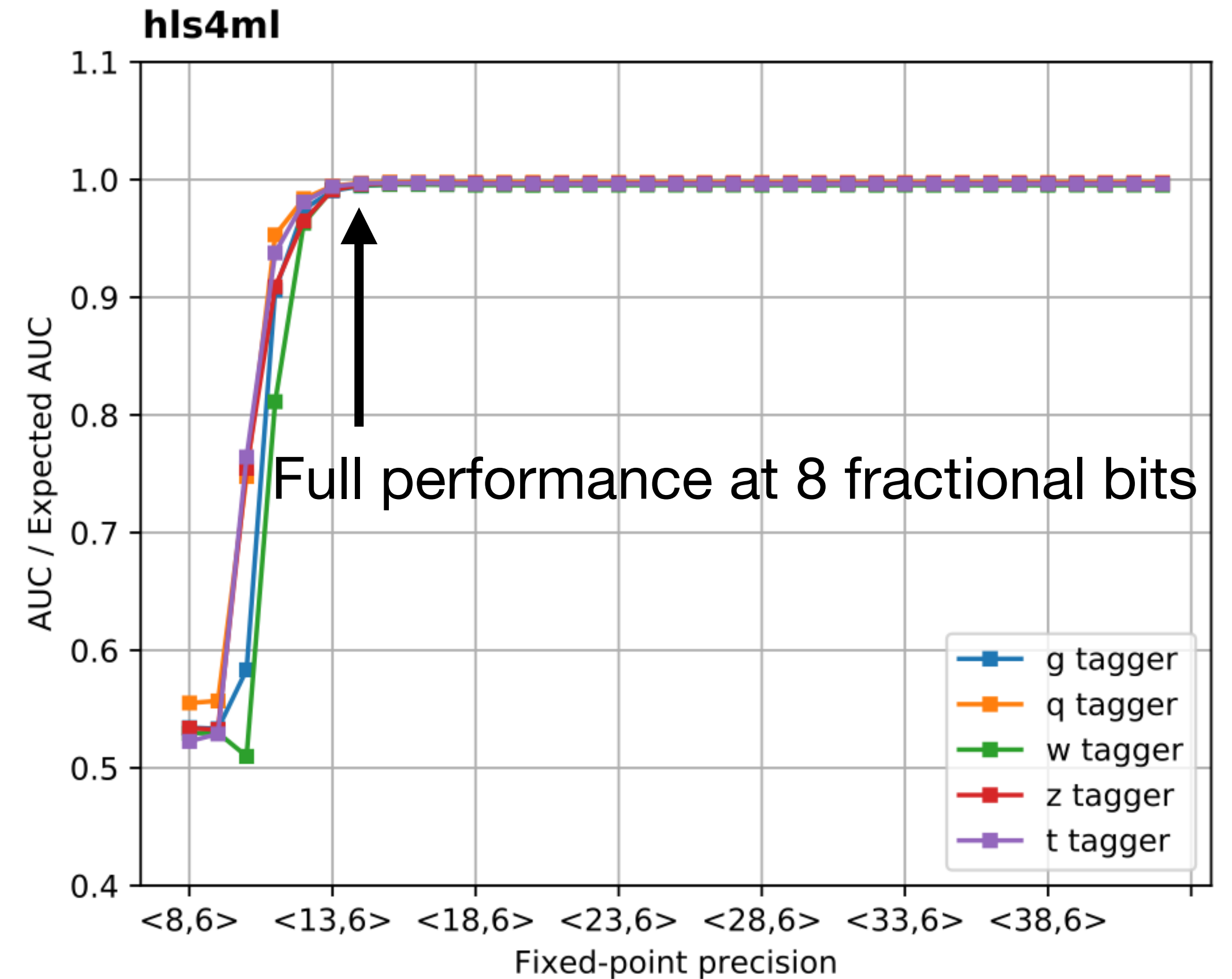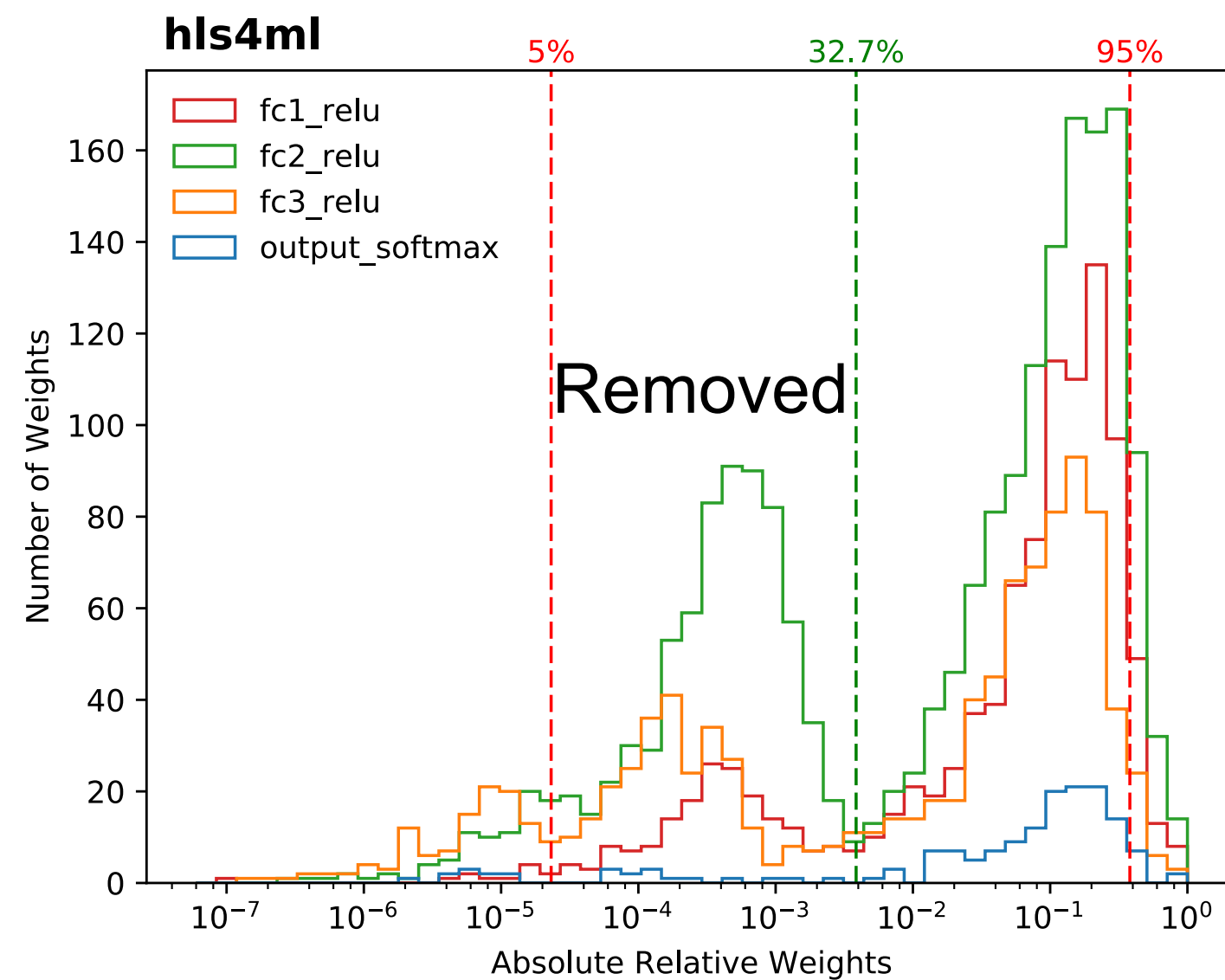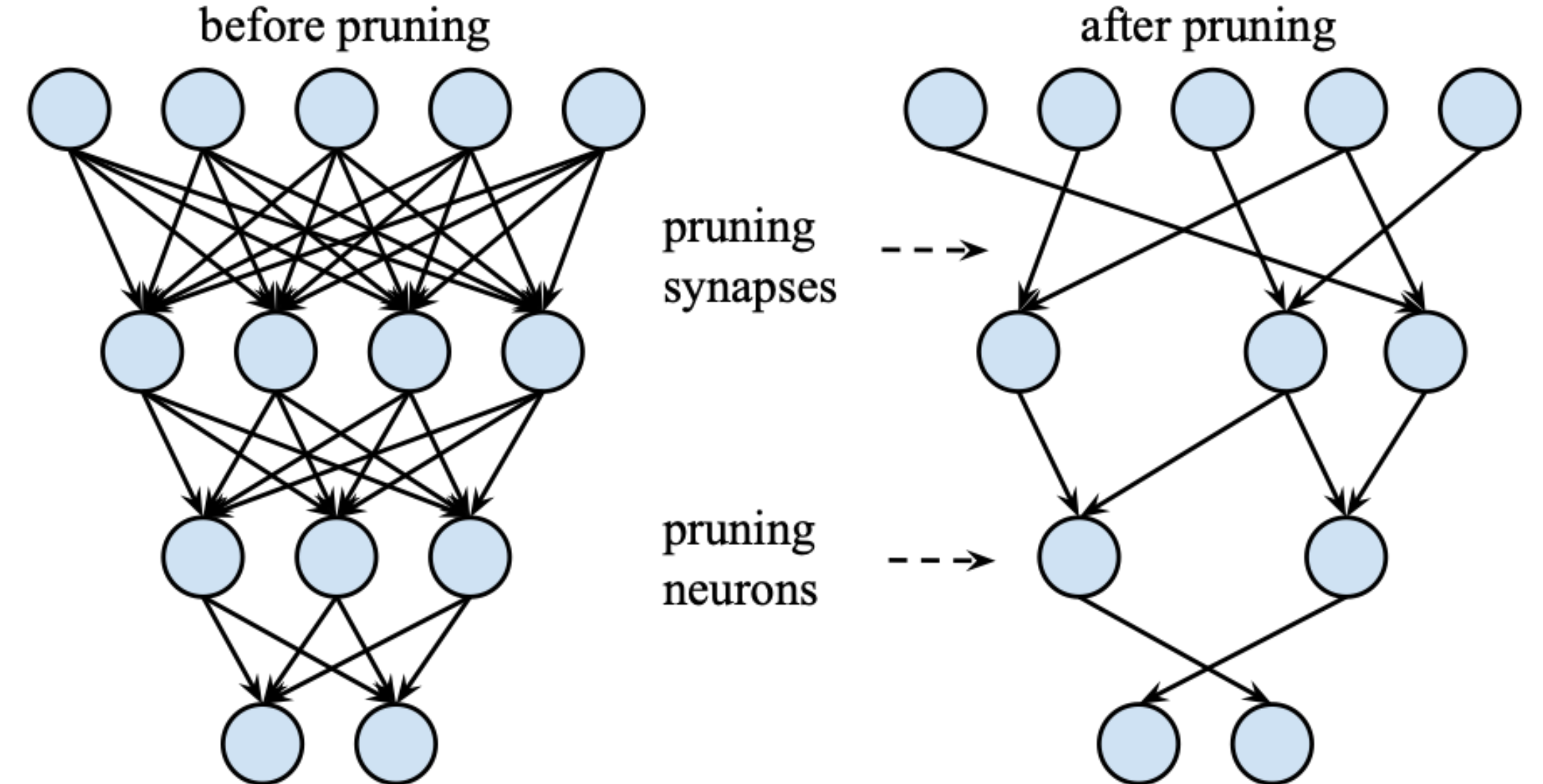
# Quantization

ap_fixed<width bits, integer bits>

0101.1011101010

Integer    Fractional

Width

<Total, Fractional>



Full performance at 6 integer bits

Scan integer bits, fractional bits fixed to 8



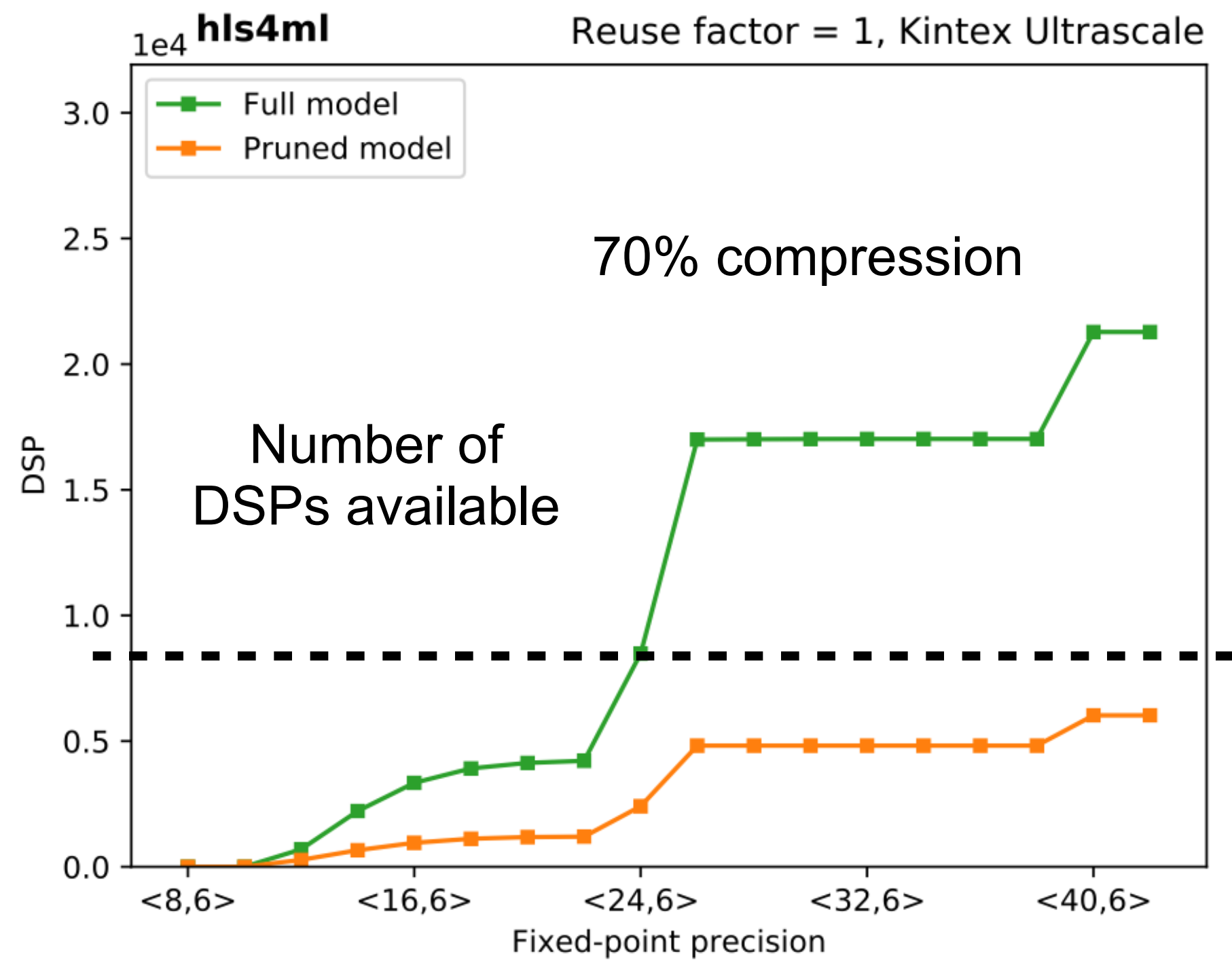Full performance at 8 fractional bits

Scan fractional bits, Integer bits fixed to 6

# Pruning







- DSPs (used for multiplication) are often limiting resource
  - maximum use when fully parallelized
  - Number of DSPs per multiplication changes with precision
- Iterative pruning with L1 norm penalty term: penalizes small weights

# Network compression/Efficient Machine Learning Computing

- Many approaches have been studied:

  - Parameter pruning: selective removal of weights based on a particular ranking [arxiv.1510.00149, arxiv.1712.01312]

  - **Neural Network Architecture Search (NAS)** [https://arxiv.org/pdf/2301.08727.pdf]

  - **Knowledge distillation**: training a compact network with distilled knowledge of a large network [https://arxiv.org/abs/1503.02531]

  - Low-rank factorization: using matrix/tensor decomposition to estimate informative parameters [arxiv.1405.3866]

  - Transferred/compact convolutional filters: special structural convolutional filters to save parameters [arxiv.1602.07576]

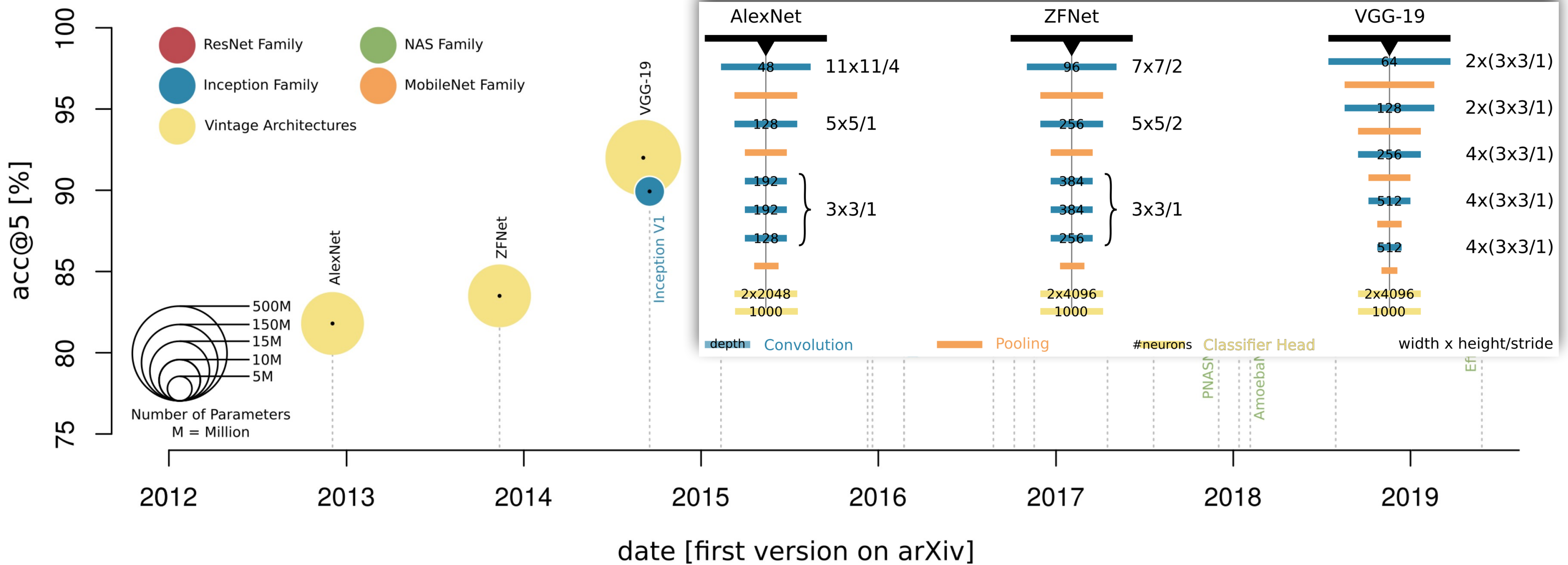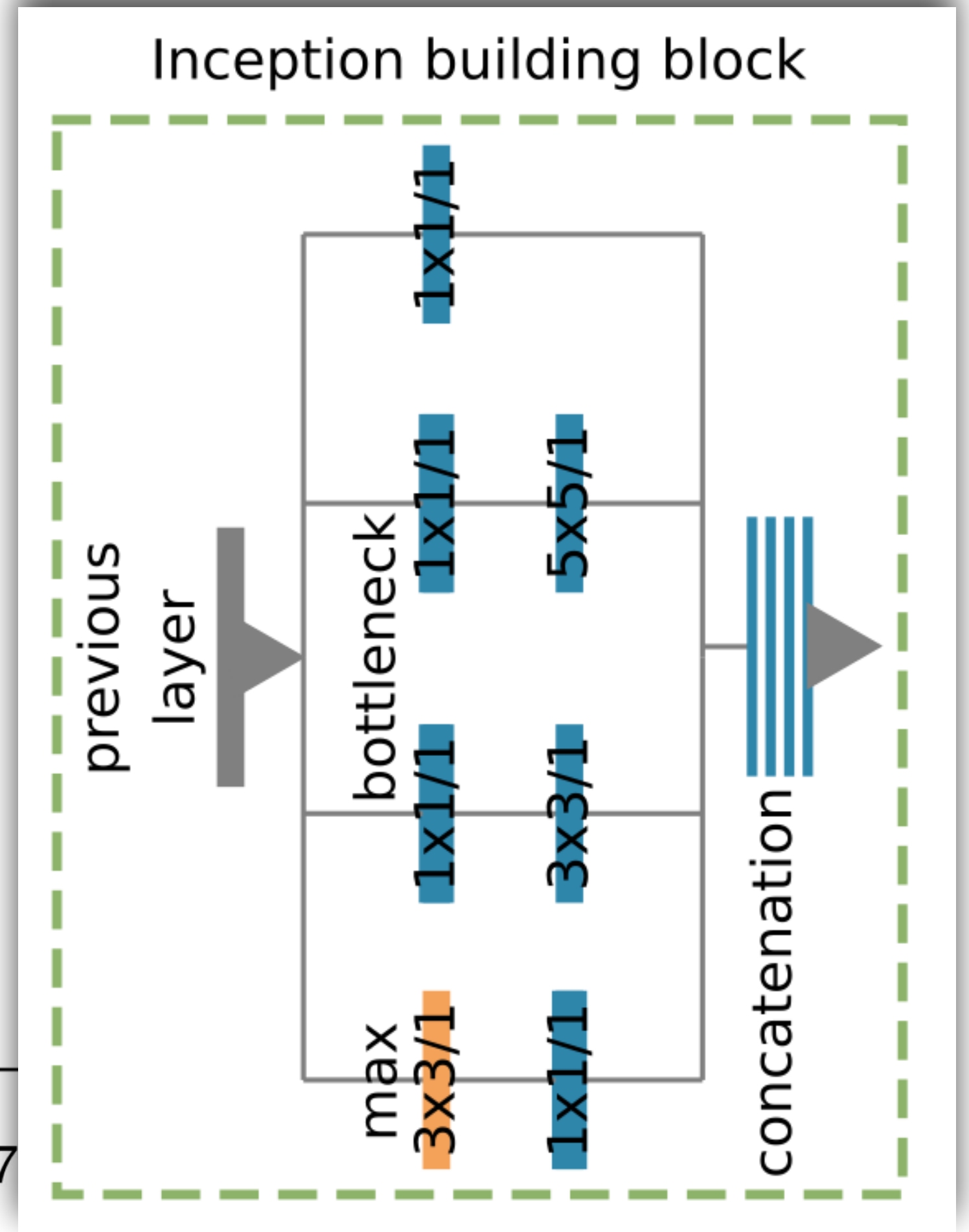  - Tensorflow model sparsity toolkit: https://blog.tensorflow.org/2019/05/tf-model-optimization-toolkit-pruning-API.html

# Image detection network evolution
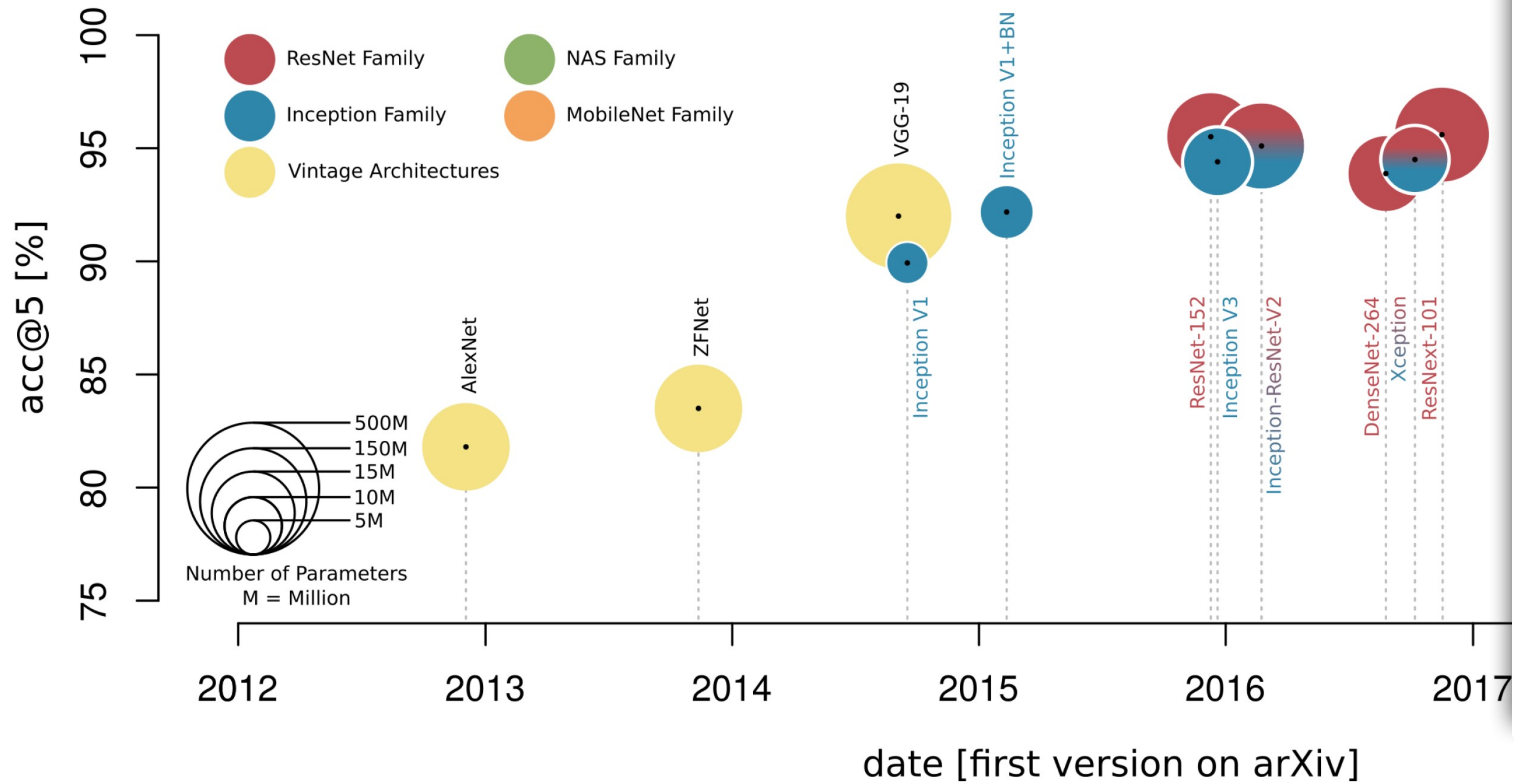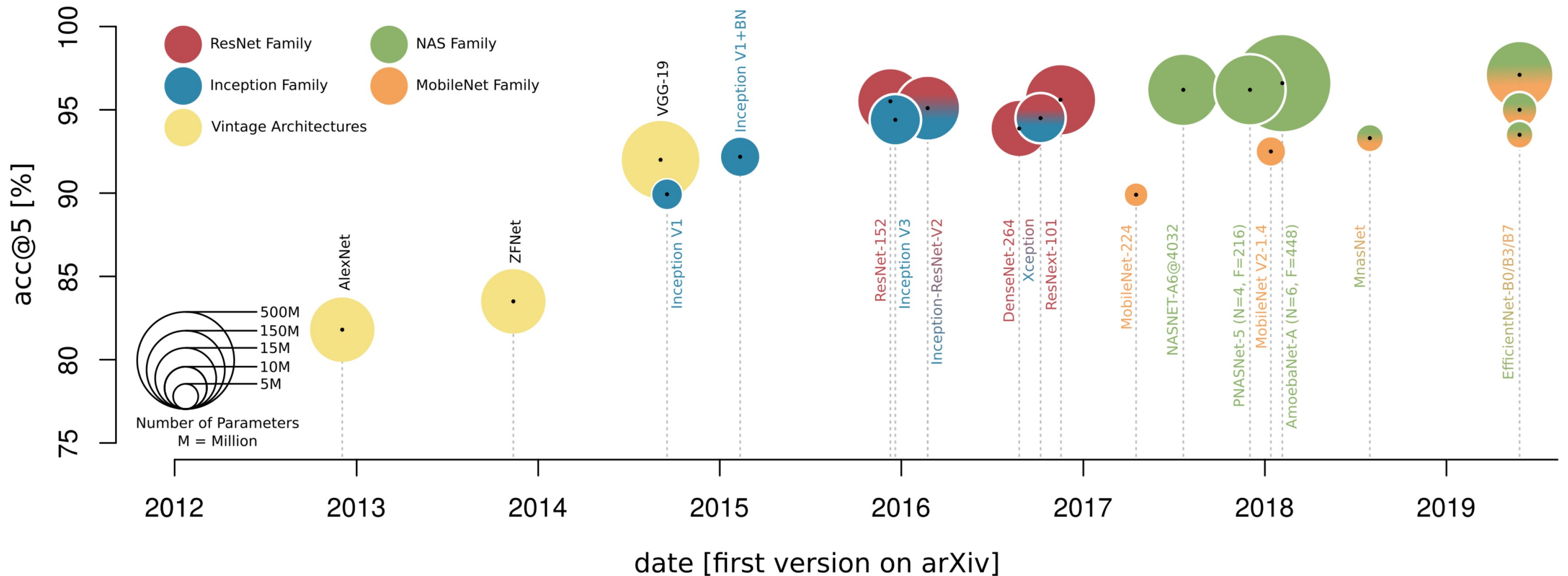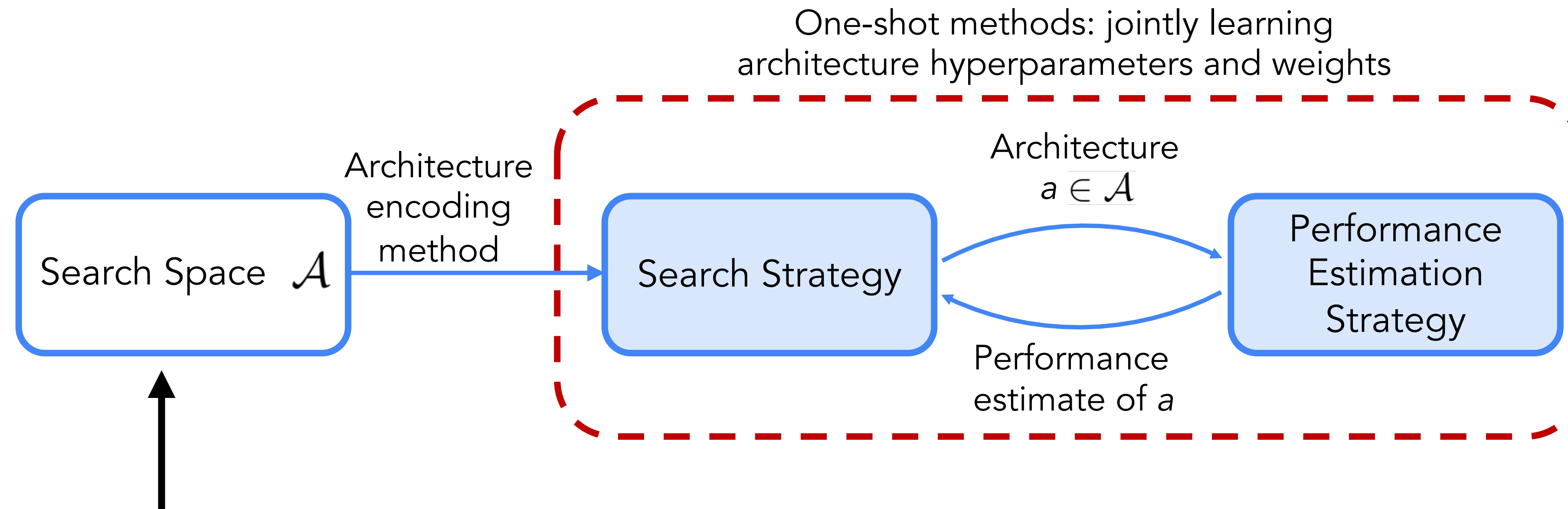
# Image detection network evolution
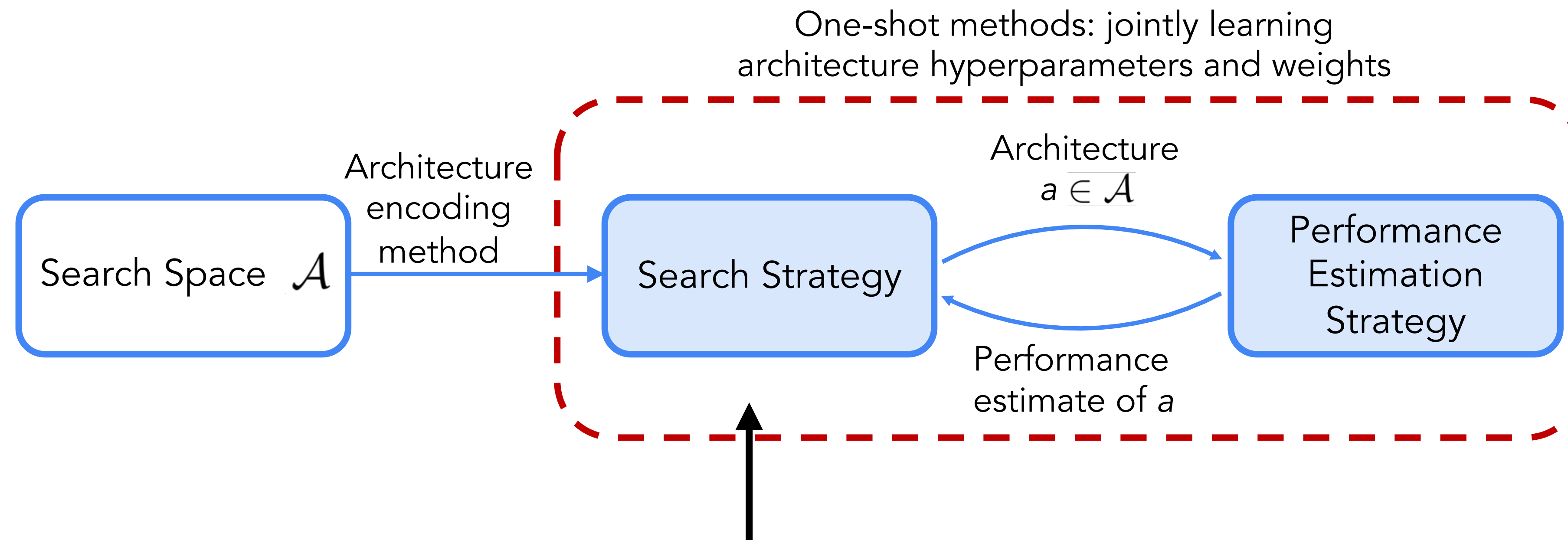
# Image detection network evolution



The modular concept of CNNs and their building blocks is crucial for the next group of architectures by **Neural Architecture search.**

# Automation: Neural Network Architecture Search



One-shot methods: jointly learning
architecture hyperparameters and weights

Search Space $\mathcal{A}$

Architecture
encoding
method

Architecture
$a \in \mathcal{A}$

Search Strategy

Performance
estimate of $a$
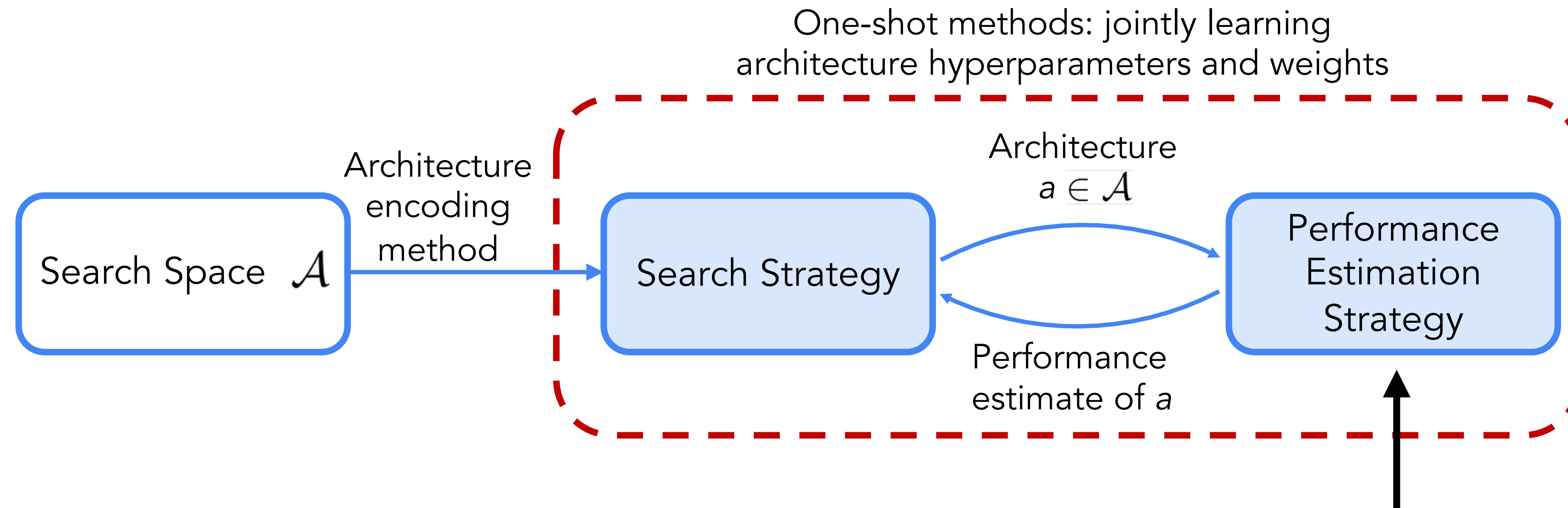
Performance
Estimation
Strategy

- Choose building blocks:

- Operation/primitive: denotes the atomic unit, a popular one is a triplet of a fixed activation, operation, and fixed normalization, such as ReLU-conv 1x1-batchnorm

- Layer, Block, Cell, Motif
.

# Neural Network Architecture Search



One-shot methods: jointly learning
architecture hyperparameters and weights

Search Space $\mathcal{A}$

Architecture
encoding
method

Search Strategy

Architecture
$a \in \mathcal{A}$

Performance
estimate of $a$

Performance
Estimation
Strategy

- Grid search, random search, re-enforcement learning, evolution, Bayesian

# Neural Network Architecture Search

One-shot methods: jointly learning
architecture hyperparameters and weights

Architecture
encoding
method

Architecture
$a \in \mathcal{A}$

Search Space $\mathcal{A}$

Search Strategy

Performance
Estimation
Strategy

Performance
estimate of $a$

- Latency, accuracy, power consumption, hardware types etc

# NAS for MobileNet-v2



**Figure 8.** Conceptual overview of the *Efficient designs*: (**Left**) The MobileNetV2 building block [108], here with an additional Squeeze-and-Excitation (SE) module [47]. In comparison with a ResNet building block, the bottleneck design is inverted, so that first the expansion factor (t) is larger than 1, which leads to intermediate deeper feature maps (td) as the final output depth of the building block (d'). (**Middle**) The Mnas search space with a fixed overall architecture of the network, the skeleton, but fully optional layer designs, based on the MobileNetV2 building block [51]. (**Right**) A recurrent neural network (RNN) [34] controller searches the search space for the best performing combination of layer designs by maximising an optimisation rule [48]. The resulting architecture is scaled in depth, width and resolution to become the EfficientNet-B7 architecture, the sota design in late 2019 [52].

# Distilling the Knowledge in a Neural Network

**Geoffrey Hinton**[*][†]
Google Inc.
Mountain View
geoffhinton@google.com

**Oriol Vinyals**[†]
Google Inc.
Mountain View
vinyals@google.com

**Jeff Dean**
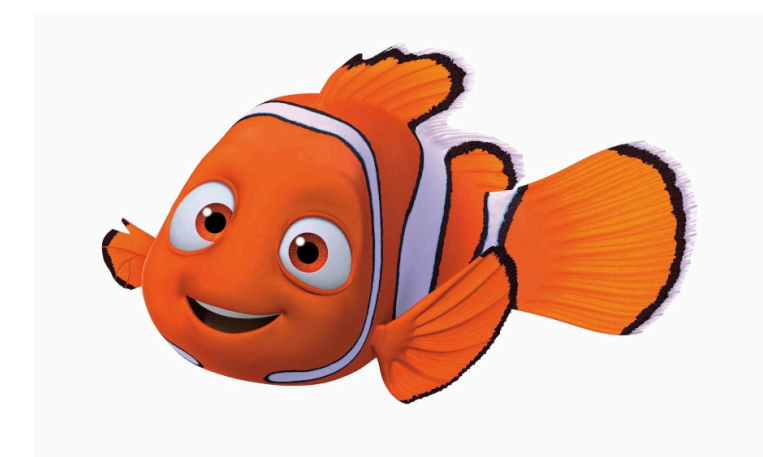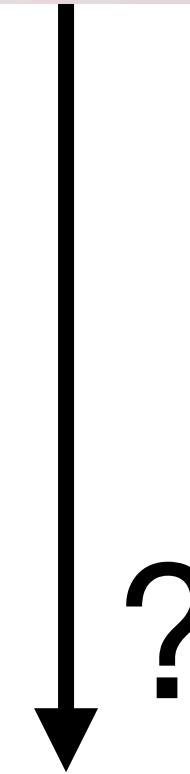Google Inc.
Mountain View
jeff@google.com

## Abstract

A very simple way to improve the performance of almost any machine learning algorithm is to train many different models on the same data and then to average their predictions [3]. Unfortunately, making predictions using a whole ensemble of models is cumbersome and may be too computationally expensive to allow deployment to a large number of users, especially if the individual models are large neural nets. Caruana and his collaborators [1] have shown that it is possible to compress the knowledge in an ensemble into a single model which is much easier to deploy and we develop this approach further using a different compression technique. We achieve some surprising results on MNIST and we show that we can significantly improve the acoustic model of a heavily used commercial system by distilling the knowledge in an ensemble of models into a single model. We also introduce a new type of ensemble composed of one or more full models and many specialist models which learn to distinguish fine-grained classes that the full models confuse. Unlike a mixture of experts, these specialist models can be trained rapidly and in parallel.

# Knowledge distillation

- We have trained a fully supervised model with MLP (fully connected neural networks)

    - Overfit with an overparametrized model

    - Add regularization to improve generalizability

    - Ensemble of models

- Accurate, but big and cumbersome—> not suitable for computing resource constrained use cases

- Small model is not as performant


- Can we transfer the knowledge learned by the large teacher model to a student model?

    - Efficient and performant

Blobfish

?

# Knowledge distillation

- Where is this knowledge stored?

  - Multi-class classification: "Soft labels" that generalizes to unseen datasets

Blobfish

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$ ← Logits

?

$\sigma$ = softmax

$\vec{z}$ = input vector

$e^{z_i}$ = standard exponential function for input vector

$K$ = number of classes in the multi-class classifier

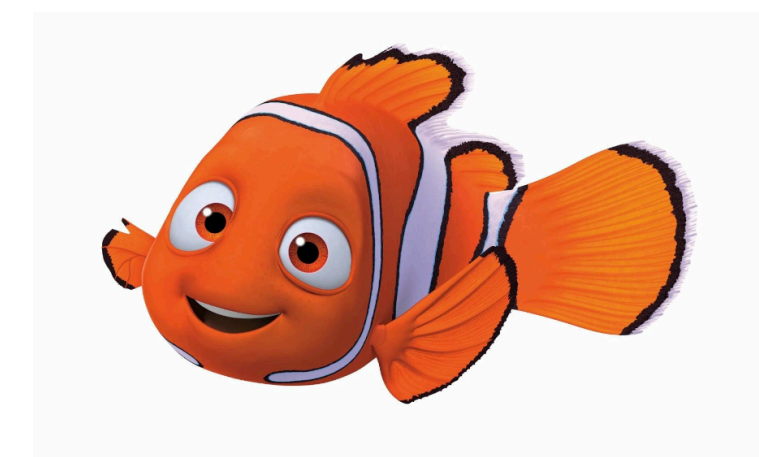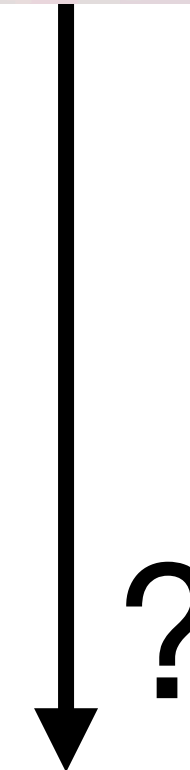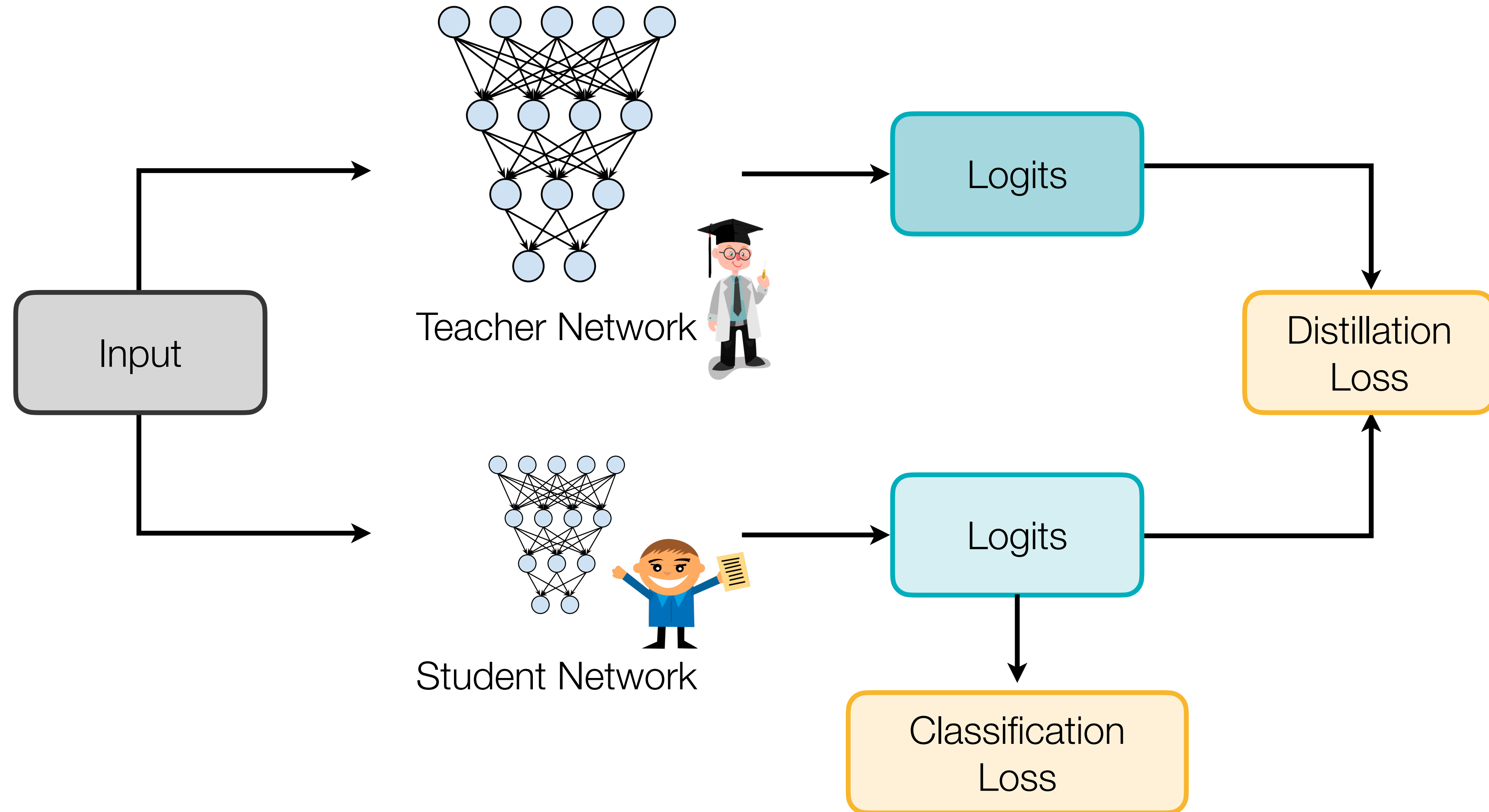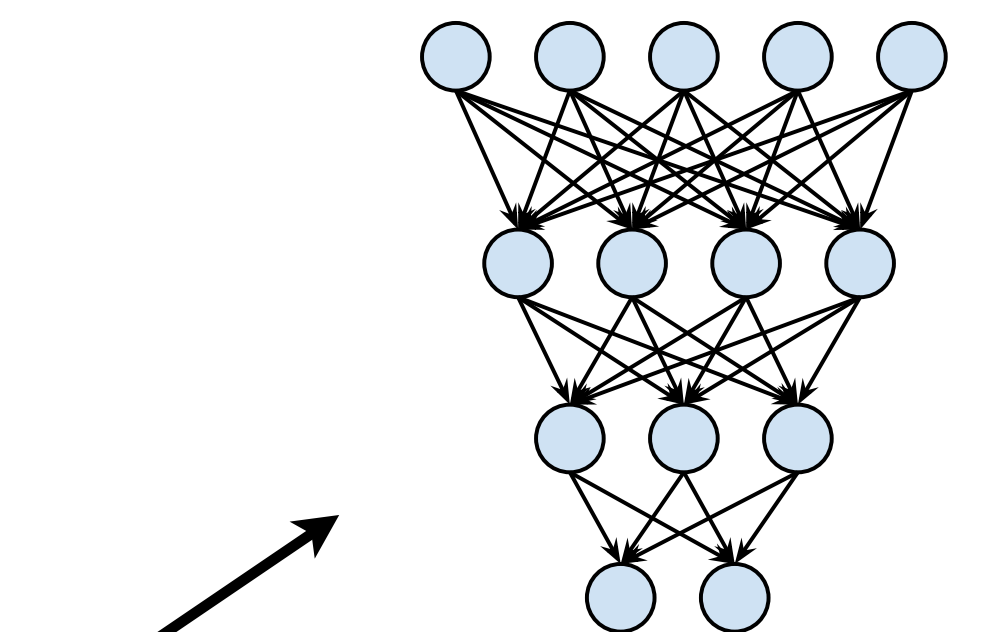$e^{z_j}$ = standard exponential function for output vector

# Illustration of knowledge distillation

# Matching prediction probabilities between teacher and student



| | Logits | Probabilities | |
|---|---|---|---|
| Cat | 5 | 0.982 | $\dfrac{\exp(5)}{\exp(5) + \exp(1)}$ |
| Dog | 1 | 0.017 | |
| Cat | 5 | 0.982 | $\dfrac{\exp(1)}{\exp(5) + \exp(1)}$ |
| Dog | 5 | 0.017 | |
| Cat | 8 | 0.731 | |
| Dog | 2 | 0.269 | The student model is less confident |
| Cat | 8 | 0.731 | |
| Dog | 3 | 0.500 | |
| Cat | 3 | 0.731 | |

Teacher Network

Student Network

# Student model gives less confident predictions



| | Logits | Probabilities |
|---|---|---|
| Cat | 5 | 0.982 |
| Dog | 1 | 0.017 |
| Cat | 5 | 0.982 |
| Dog | 1 | 0.017 |
| Cat | 3 | 0.731 |
| Dog | 2 | 0.267 |
| Cat | 3 | 0.731 |
| Dog | 2 | 0.267 |

# Concept of temperature

before pruning  after pruning

pruning
synapses

pruning
neurons

before pruning  after pruning

pruning
synapses

pruning
neurons

Teacher Network

before pruning  after pruning

pruning
synapses

pruning
neurons

before pruning  after pruning

$$\frac{\exp(5/1)}{\exp(5/1) + \exp(1/1)}$$

| | Logits | Probabilities (T=1) | Probabilities (T=10) |
|---|---|---|---|
| Cat | 5 | 0.982 | |
| Dog | 1 | 0.018 | 0.599 |
| Cat | 5 | 0.982 | 0.401 |
| Dog | 1 | 0.018 | |
| Cat | 3 | 0.731 | |
| Dog | 2 | 0.269 | |

$$\frac{\exp(5/10)}{\exp(5/10) + \exp(1/10)}$$

# Align the class probability distributions from teacher and student networks

Soft labels with increased 'temperature'



$$q_i = \frac{exp(z_i/T)}{\sum_j exp(z_j/T)}$$

Logits

Logits

Layer 1    Layer 2

Teacher Network

Layer N

Input

Input

Layer 1    Layer 2

Student Network

Layer N

Logits

Logits

Distillation Loss

Classification Loss

Classification Loss

Cross entropy loss:
$$\mathbb{E}(-p_t \log p_s);$$

L2 loss:
$$E(\|p_t - p_s\|_2^2)$$

# Match intermediate feature maps



Like What You Like: Knowledge Distill via Neuron Selectivity Transfer [Huang and Wang, arXiv 2017]

# Match intermediate attention maps

**Gradients of feature maps are used to characterize "attention" of DNNs**

- The attention of a CNN feature map $x$ is defined as $\dfrac{\partial L}{\partial x}$, where $L$ is the learning objective.

- Intuition: If $\dfrac{\partial L}{\partial x_{i,j}}$ is large, a small perturbation at $i, j$ will significantly impact the final output. As a result, the network is putting more attention on position $i, j$.
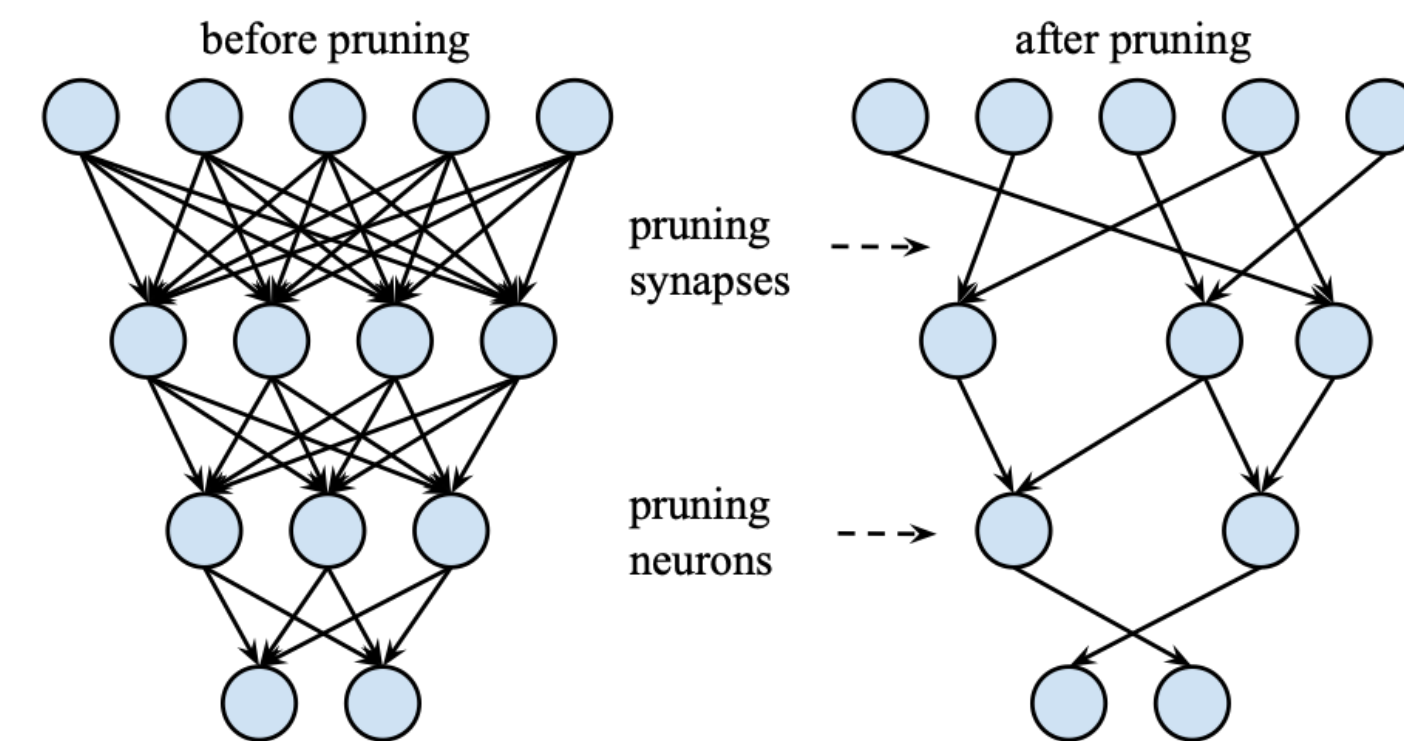
input image                    attention map



Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer [Zagoruyko and Komodakis, ICLR 2017]

# Efficient Algorithms

# Efficient Algorithms

# Efficient Algorithms



## Pruning



before pruning        after pruning

pruning
synapses

pruning
neurons

# Efficient Algorithms



Quantization          Compression/Pruning

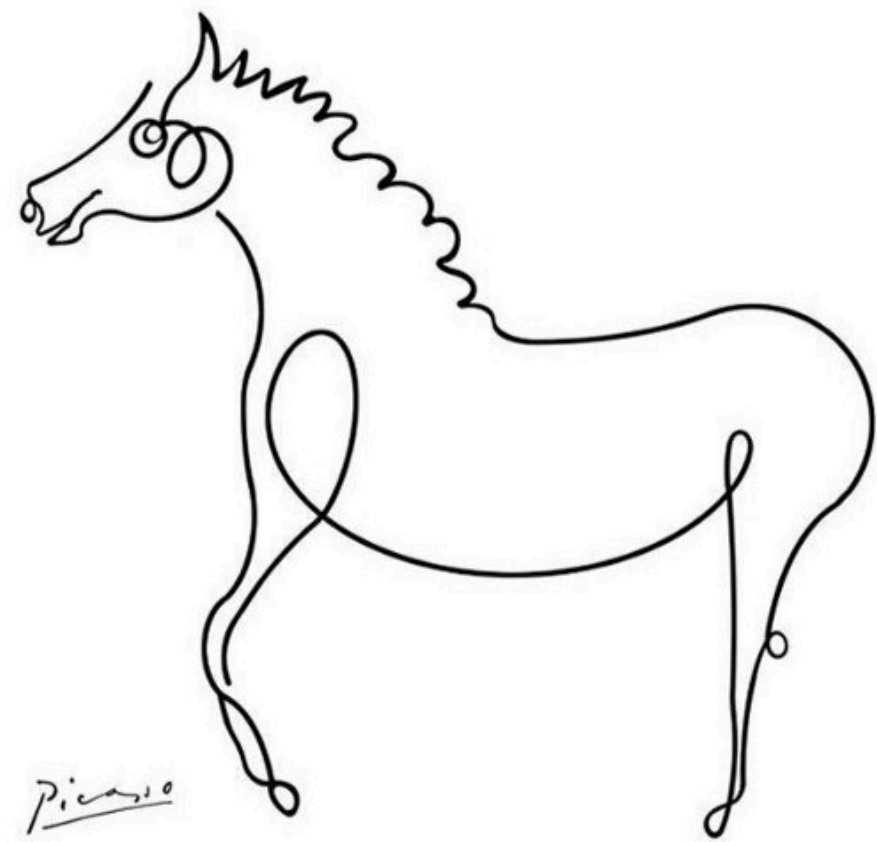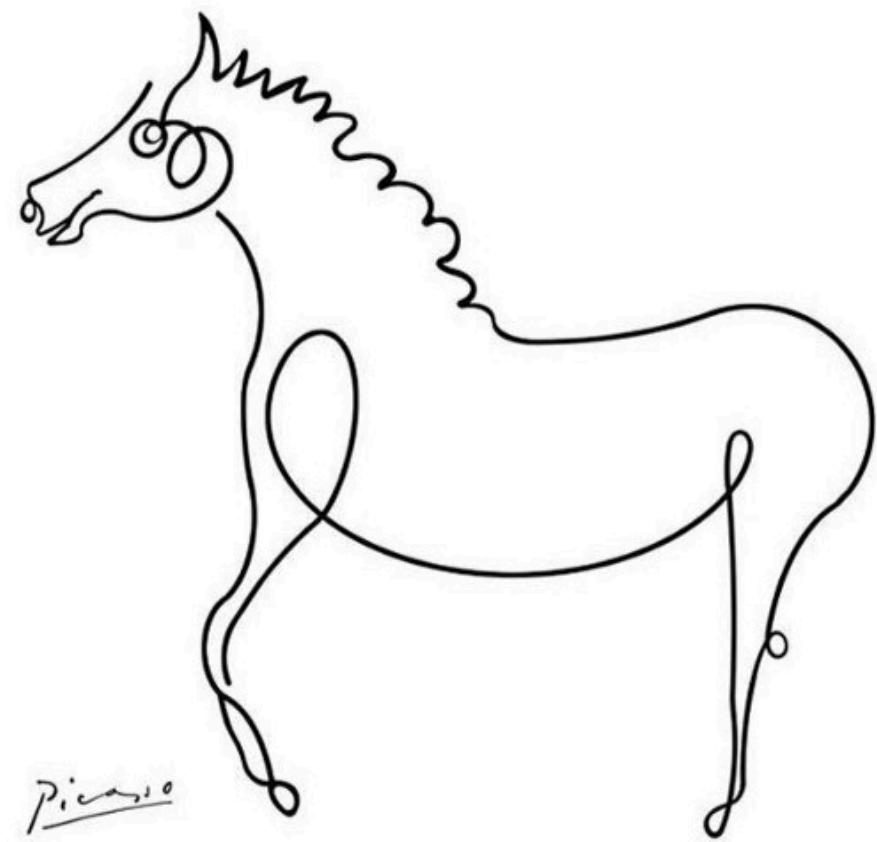'Ultimate optimization' of 'bits of information':
Quantization Aware Pruning

https://arxiv.org/abs/2102.11289
https://arxiv.org/abs/2304.06745

# Efficient Algorithms



Compress it creatively:
knowledge distillation. e.g

# Efficient Algorithms



馬

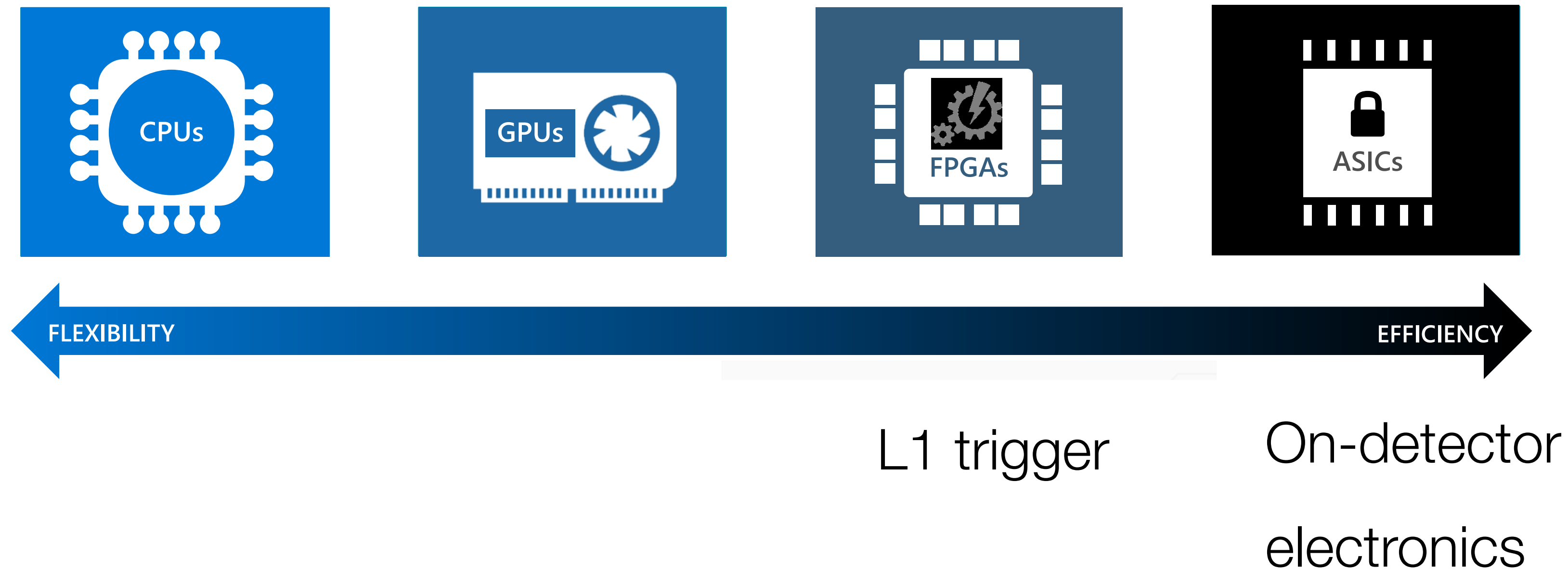Neural Architecture search
e.g. EfficientNet for image
detection

# Efficient Algorithms

# Co-design tool for Specialized Hardware



FLEXIBILITY ← → EFFICIENCY

L1 trigger

On-detector electronics

**Co-design tool: crucial for prototyping AI at edge solutions**

Algorithm hardware co-design with limited computing

Prototype with manageable programming barrier for domain scientists

# Hardware Pros and Cons

More Flexible

Homogeneous

Specialized

More Efficient

|  | Perf/W | |
|---|---|---|
| CPUs | 1X | Today's standard, most programmable, good for services changing rapidly |
| Manycore CPUs | 3X | Many simple cores (10s to 100s per chip), useful if software can be fine-grain parallel, difficult to maintain. |
| GPUs | 5-30X | Good for data parallelism by merged threads (SIMD), High memory bandwidth, power hungry |
| **FPGAs** | 5-30X | Most radical fully programmable option. Good for streaming/irregular parallelism. Power efficient but currently need to program in H/W languages. |
| Structured ASICS | 20-100X | Lower-NRE ASICs with lower performance/efficiency. Includes domain-specific (programmable) accelerators. |
| Custom ASICs | > 100X | Highest efficiency. Highest NRE costs. Requires high volume. Good for functions in very widespread use that are stable for many years. |

Conventional programming

Alternative programming

Can't change functionality

C/C++

CUDA

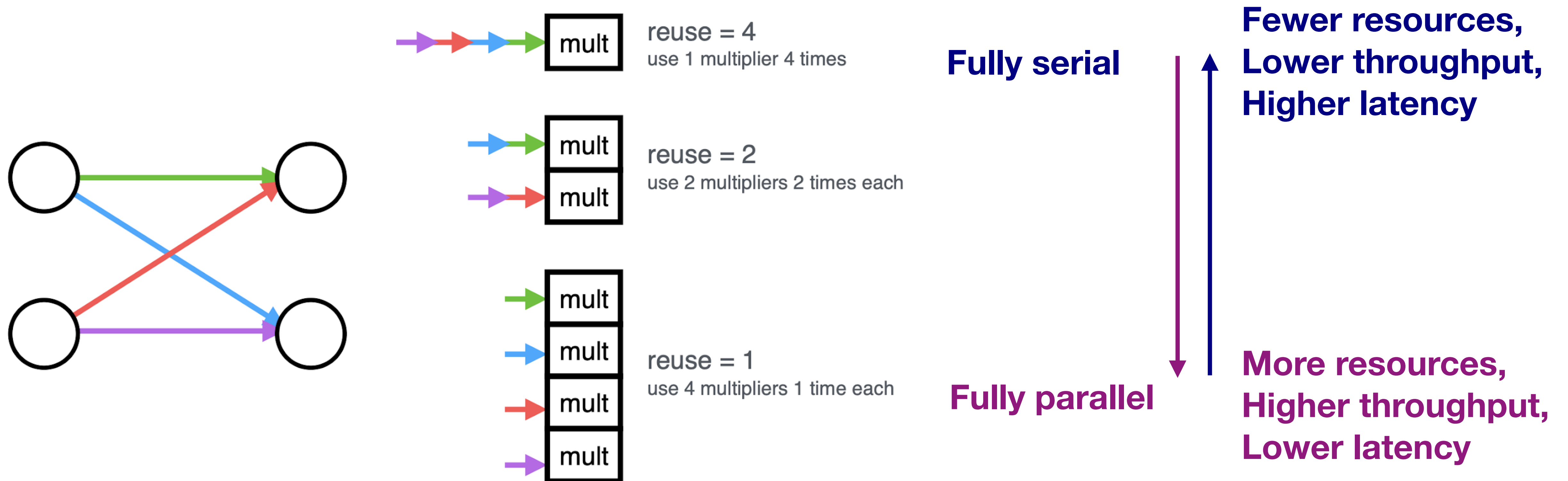Verilog

Verilog

# Co-design
## Connecting domain scientists with prototype solutions
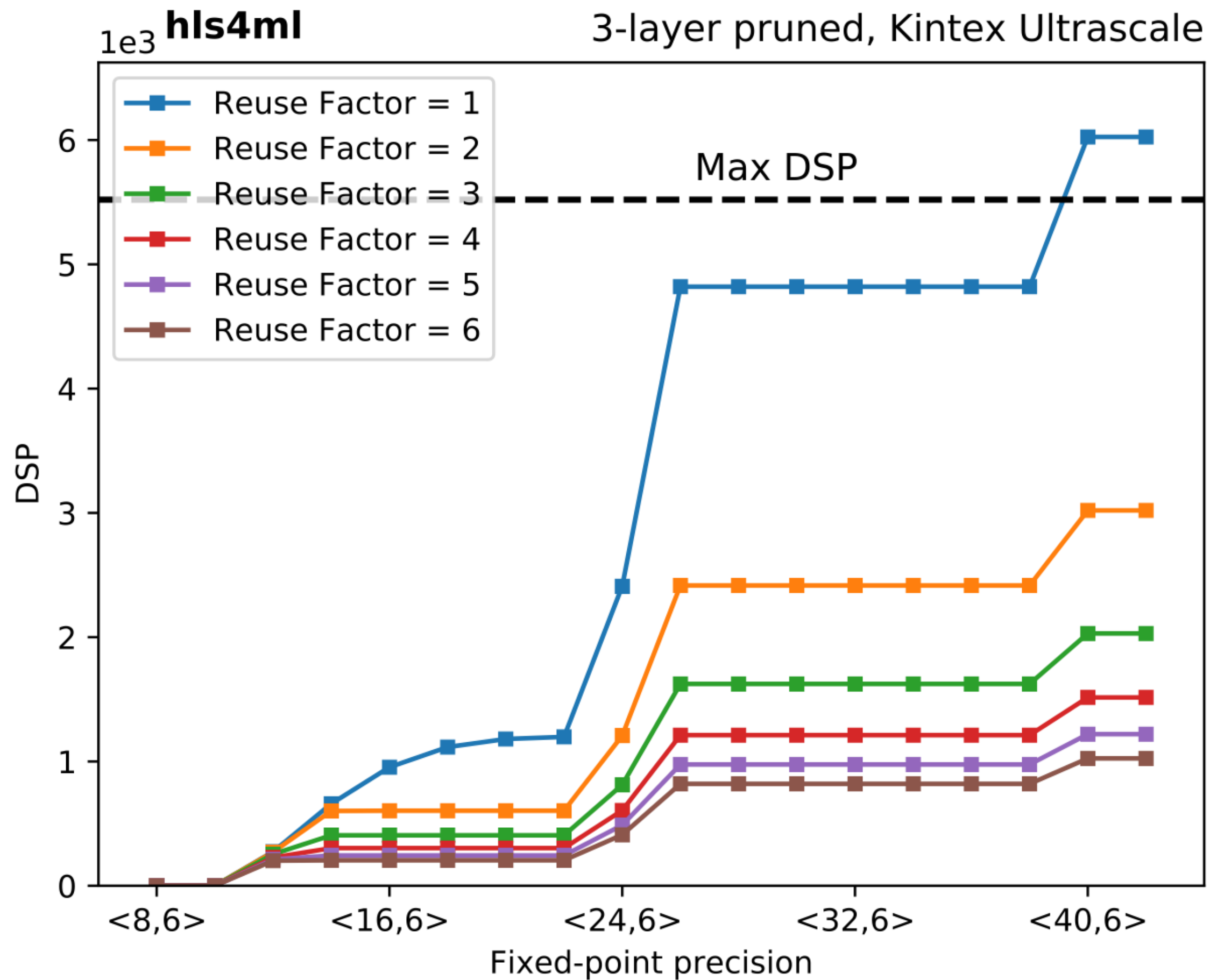


HLS4ML: to aid prototype science application solutions.

# Example Knobs to tune: Reuse factor to balance latency and resources

- Trade-off between latency and FPGA resource usage determined by the parallelization of the calculations in each layer
- Configure the **"reuse factor"** = number of times a multiplier is used to do a computation



reuse = 4
use 1 multiplier 4 times

reuse = 2
use 2 multipliers 2 times each

reuse = 1
use 4 multipliers 1 time each

**Fully serial**

**Fully parallel**

**Fewer resources, Lower throughput, Higher latency**

**More resources, Higher throughput, Lower latency**

# Reuse factor to balance latency and resources



3-layer pruned, Kintex Ultrascale

**More resources**

Fully parallel
Each mult. used 1x

Each mult. used 2x

Each mult. used 3x

**Longer latency**

# Trade-off Between Efficiency and Accuracy



Image source: 1

Image source: 2

# HLS4ML tutorial

## hls4ml-tutorial: Tutorial notebooks for `hls4ml`

[JB jupyter book] [deploy-book passing] [code style black] [pre-commit enabled] [launch binder]
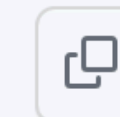
There are several ways to run the tutorial notebooks:

## Online

[launch binder]

## Conda

The Python environment used for the tutorials is specified in the `environment.yml` file. It can be setup like:

```
conda env create -f environment.yml
conda activate hls4ml-tutorial
```

## Docker without Vivado

Pull the prebuilt image from the GitHub Container Registry:

```
docker pull ghcr.io/fastmachinelearning/hls4ml-tutorial/hls4ml-0.8.0:latest
```

- https://github.com/fastmachinelearning/hls4ml-tutorial

46

# ML everywhere in CMS Phase 2 L1 Trigger

**B-tagging**



**CMS** *Phase-2 Simulation Preliminary*    14 TeV, 200 PU

HH → bbbb events
- b-tag NN @ 15 kHz
- QuadJet+HT OR Jets+Muon @ 15 kHz
- QuadJet+HT @ 11 kHz

Efficiency / Ratio to QuadJet+HT vs Offline $H_T$ [GeV]

**Missing transverse energy**



(13 TeV)

$\sigma(u_\perp)$ [GeV] vs $q_T$ [GeV]

- PF MET
- PUPPI MET
- DeepMET



**CMS** *Phase-2 Simulation*    (14 TeV, PU 200)

- NN Puppi Tau
- Calo Tau

*Minbias*

Tau

Rate (kHz) vs Offline $p_T$ (GeV)

---

t = 0

Each small box is one
Xilinx Ultrascale+ FPGA

These have NNs
and/or BDTs inside

hls 4 ml    Conifer

t = 12.5 µs



Detector hits

Clusters & Tracks

Particles

Event Categorisation

1 bit: keep / discard

# HLS4ML: user driven development

- hls4ml 2023 roadmap plans new developments and a regular release schedule
- Q1 release: v0.7.0 & v0.7.1
  - Backend redesign to support multiple compilation targets [395]
  - Documentation updates [710, 744, 774]
  - Efficient network implementations [503, 509, 509]
  - Recurrent neural networks [560, 575]
  - Alveo accelerator FPGA card support [552]
  - Support for Vitis HLS [629]
  - Extension API [528]
- Q2 release: v0.8.0
  - Configuration editor [784]
  - PyTorch parsing improvements [799]
  - Symbolic expressions [660]
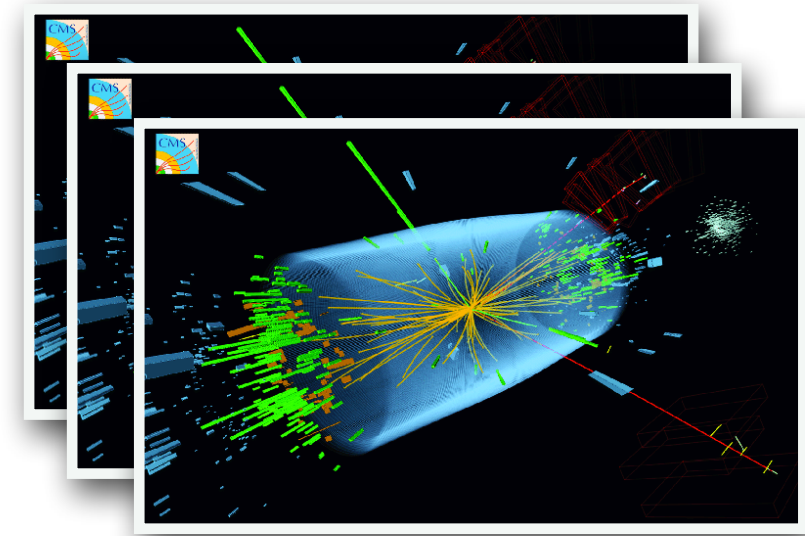  - Optimization API [768, 809]
  - Large streaming CNN [PR Soon]

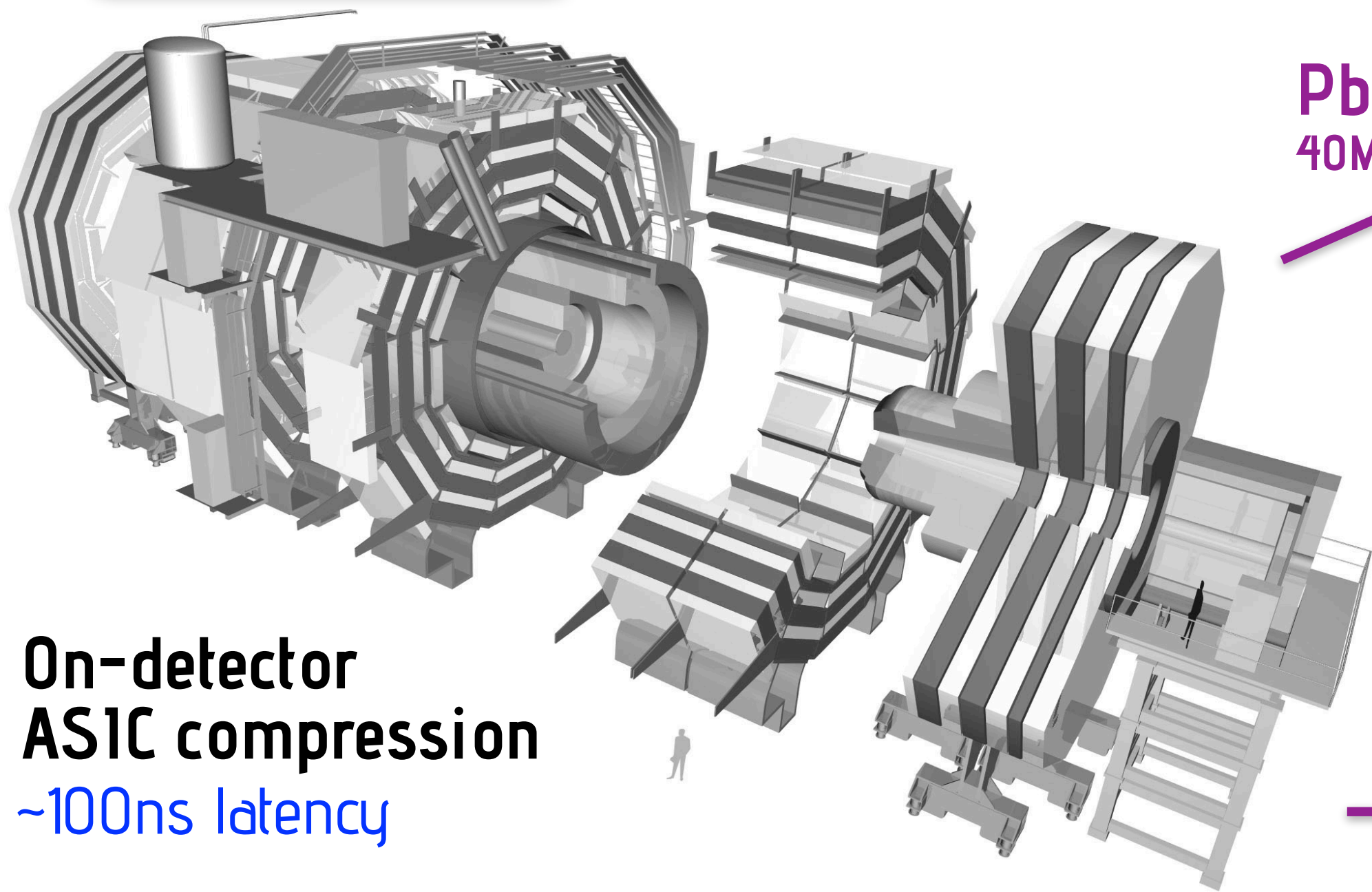  - Upcoming:

    - QONNX ingestion [591]
    - Catapult HLS

# HLS4ML: user driven development

## sPHENIX tracking GNN hls4ml synthesis results

- Network inputs: nodes=80, edges=100    **Extremely preliminary - DO NOT TRUST NUMBERS**
- Input network
    - Can be parallelized to be "nodes" times faster (i.e., 15ns)

| Latency | BRAMs | DSPs | FFs | LUTs |
|---------|-------|------|-----|------|
| 1.2 us | 6.5% | 0.3% | 5% | 7.5% |

- Edge network

| Latency | BRAMs | DSPs | FFs | LUTs |
|---------|-------|------|-----|------|
| 3 us | 15% | 2% | 20% | 65% |

- Node network (results from HLS synthesis, vivado synthesis OOM'd)
    - Neet to optimize the scatter_add function (expecting ~2us for the net)

| Latency | BRAMs | DSPs | FFs | LUTs |
|---------|-------|------|-----|------|
| 12 us | 42% | 7% | - | - |



Input network
Repeat n_graph_iters times
Edge network
Node network
Edge network

- Configuration editor [784]
- PyTorch parsing improvements [799]
- Symbolic expressions [660]
- Optimization API [768, 809]
- Large streaming CNN [PR Soon]

Support GNN/transformers

More work needed for CMS applications (100 ns latency)

# HLS4ML: user driven development

## Symbolic Regression on FPGAs for Fast Machine Learning Inference

*Ho Fung* Tsoi[1]*, *Adrian Alan* Pol[2], *Vladimir* Loncar[3,4], *Ekaterina* Govorkova[3], *Miles* Cranmer[2,5], *Sridhara* Dasu[1], *Peter* Elmer[2], *Philip* Harris[3], *Isobel* Ojalvo[2], and *Maurizio* Pierini[6]

[1]University of Wisconsin-Madison, USA
[2]Princeton University, USA
[3]Massachusetts Institute of Technology, USA
[4]Institute of Physics Belgrade, Serbia
[5]Flatiron Institute, USA
[6]European Organization for Nuclear Research (CERN), Switz

schedule

- PyTorch parsing improvements [799]
- Symbolic expressions [660]
- Optimization API [768, 809]
- Large streaming CNN [PR Soon]



Bit width

10X resource reduction

# HLS4ML: collaboration with XILINX FINN

- hls4ml 2023 roadmap plans new developments and a regular release schedule
- Q1 release: v0.7.0 & v0.7.1
    ◦ Backend redesign to support multiple compilation targets [395]
    ◦ Documentation updates [710, 744, 774]
    ◦ Efficient network implementations [503, 509, 509]
    ◦ Recurrent neural networks [560, 575]
    ◦ Alveo accelerator FPGA card support [552]
    ◦ Support for Vitis HLS [629]
    ◦ Extension API [528]
- Q2 release: v0.8.0
    ◦ Configuration editor [784]
    ◦ PyTorch parsing improvements [799]
    ◦ Symbolic expressions [660]
    ◦ Optimization API [768, 809]
    ◦ Large streaming CNN [PR Soon]

    ◦ Upcoming:

        ◦ QONNX ingestion [591]
        ◦ Catapult HLS



QONNX: Extension to the ONNX intermediate representation
format to represent arbitrary-precision quantized neural networks

# Heterogeneous computing



**CMS Experiment**
40MHz collision rate
~1B detector channels

**FPGA filter stack**
~µs latency

**Pb/s**
40MHz

**On-detector**
**ASIC compression**
~100ns latency

**10s Tb/s**
100s kHz

**10s Gb/s**
~5 kHz

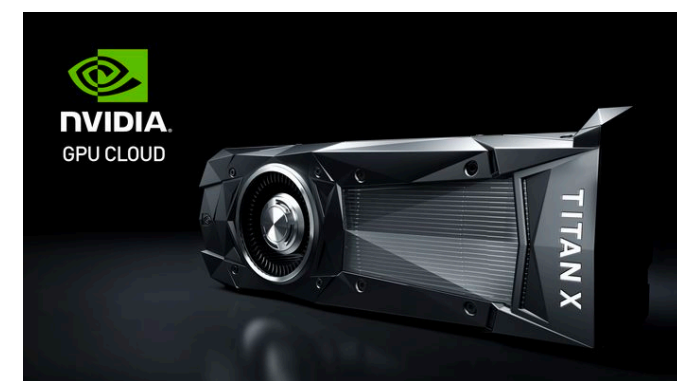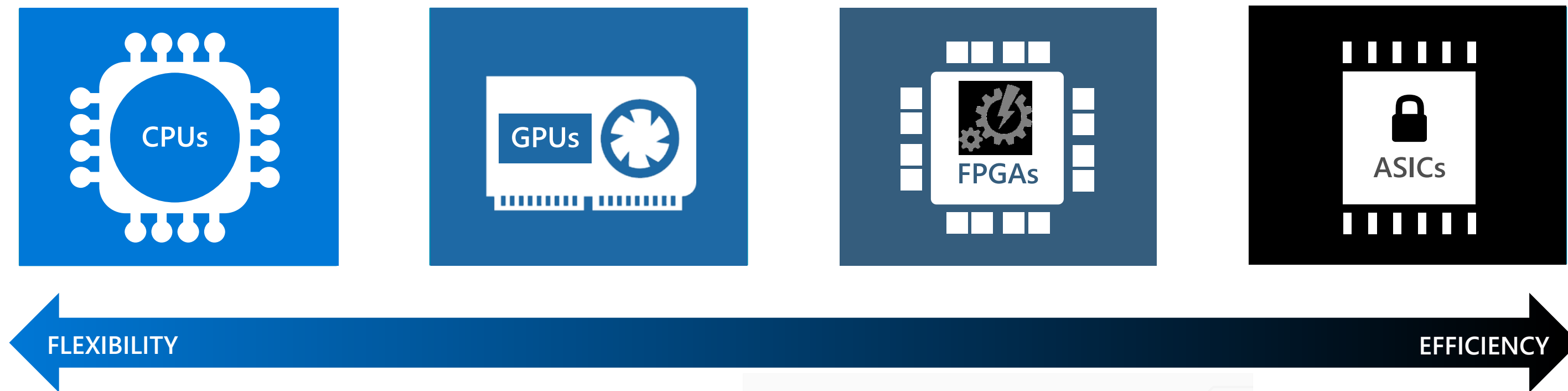**On-prem CPU/GPU**
**filter farm**
~100 ms latency

**Worldwide**
**computing grid**
Exabyte-scale
datasets

# Heterogeneous Computing for ML/AI



FLEXIBILITY ← → EFFICIENCY

CPUs | GPUs | FPGAs | ASICs

Advances driven by
big data explosion
& machine learning

A 5 year old slide, message
remains…

**Discontinued:** October 18, 2022

Groq

Graphcore
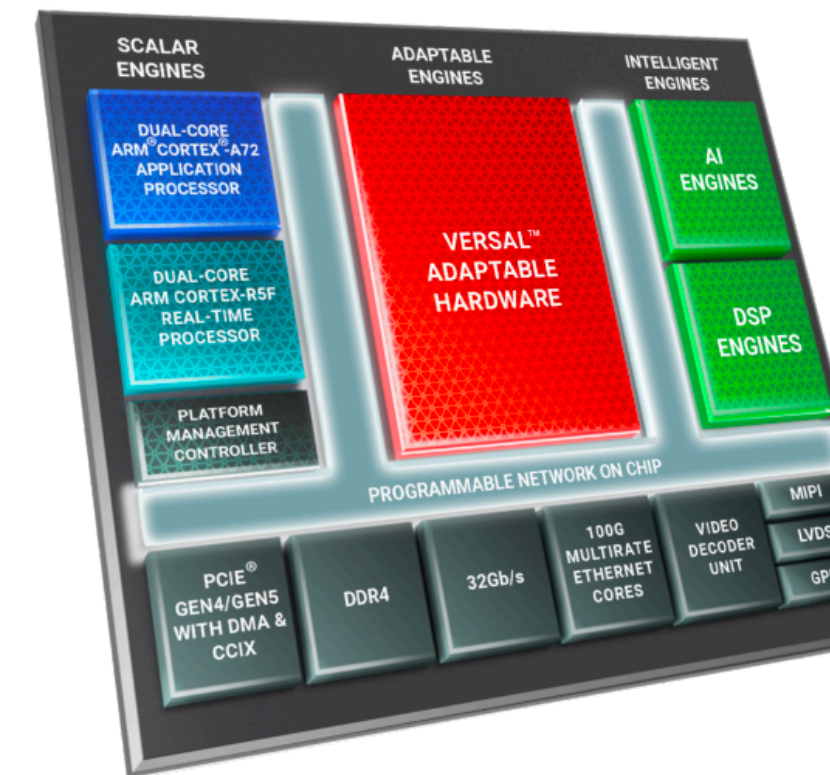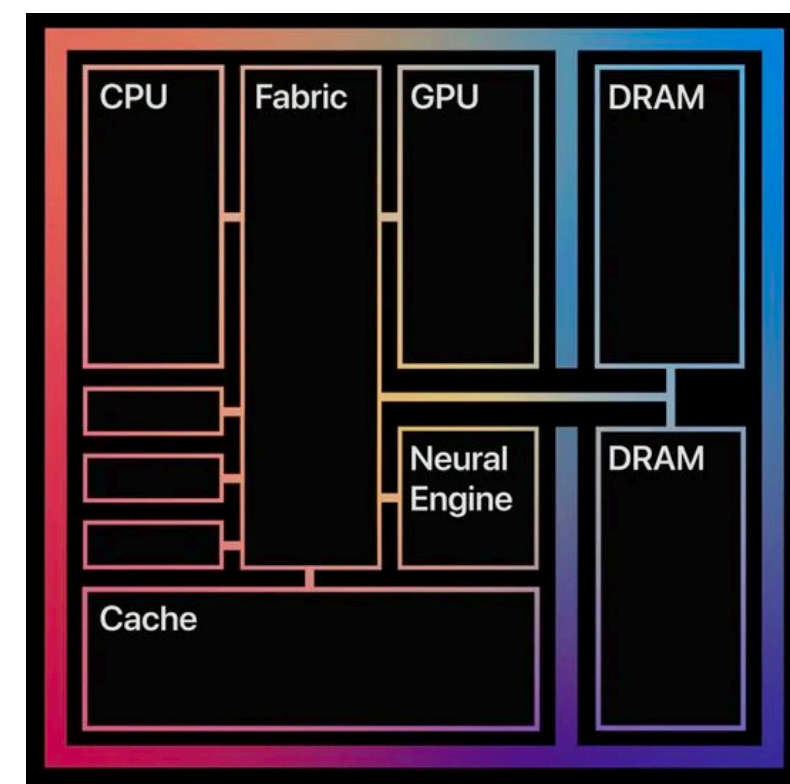
Graphcore

AMD / Xilinx

Apple

Cerebras

Meta

Groq

Cerebras WSE
1.2 Trillion transistors
46,225 mm² silicon

AMD / Xilinx

can be inefficient:

sage



GPU connection:

Narrow "sweet spot" in
terms of models or cores

or cores =
GPU

Also: workflows can only take advantage of
acceleration if they run on a machine with a
coprocessor – expensive at large scales!

Fermilab

cessor

COPROCESSOR
(GPU,FPGA,ASIC)

COPROCESSOR
(GPU,FPGA,ASIC)

COPROCESSOR
(GPU,FPGA,ASIC)

Inflexible & Expensive            Complex, Requires R&D

# Since 2018

GPU-as-a-service

*Hardware platforms*

FPGA-as-a-service Toolkit
('homegrown' gPRC server)

**hls 4 ml**

*Open source tools*

GPU-as-a-service for DUNE

*Deployment in experiments*

**Groq**

**This talk: Portable Acceleration of CMS Production Workflow with Coprocessors as a Service**
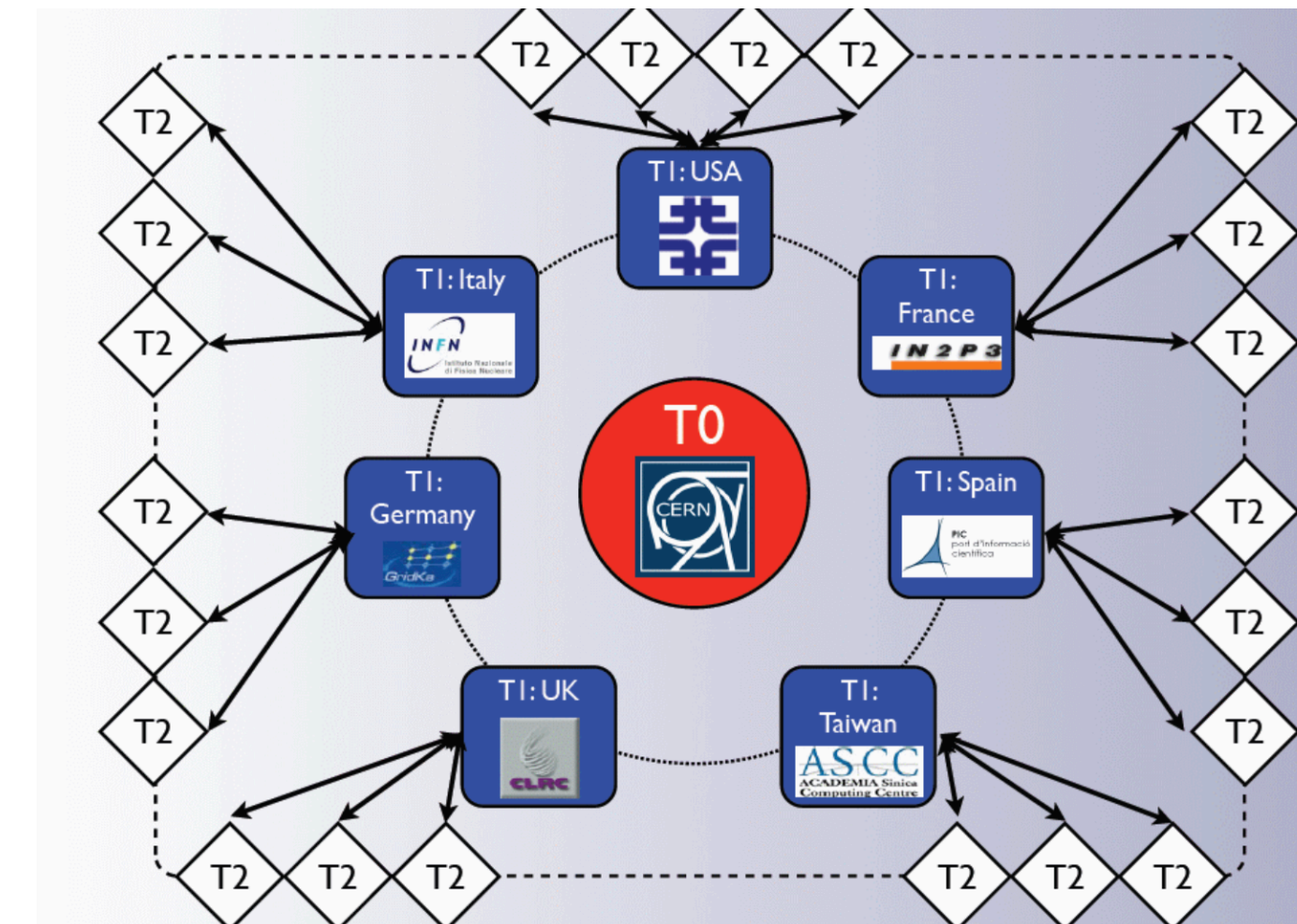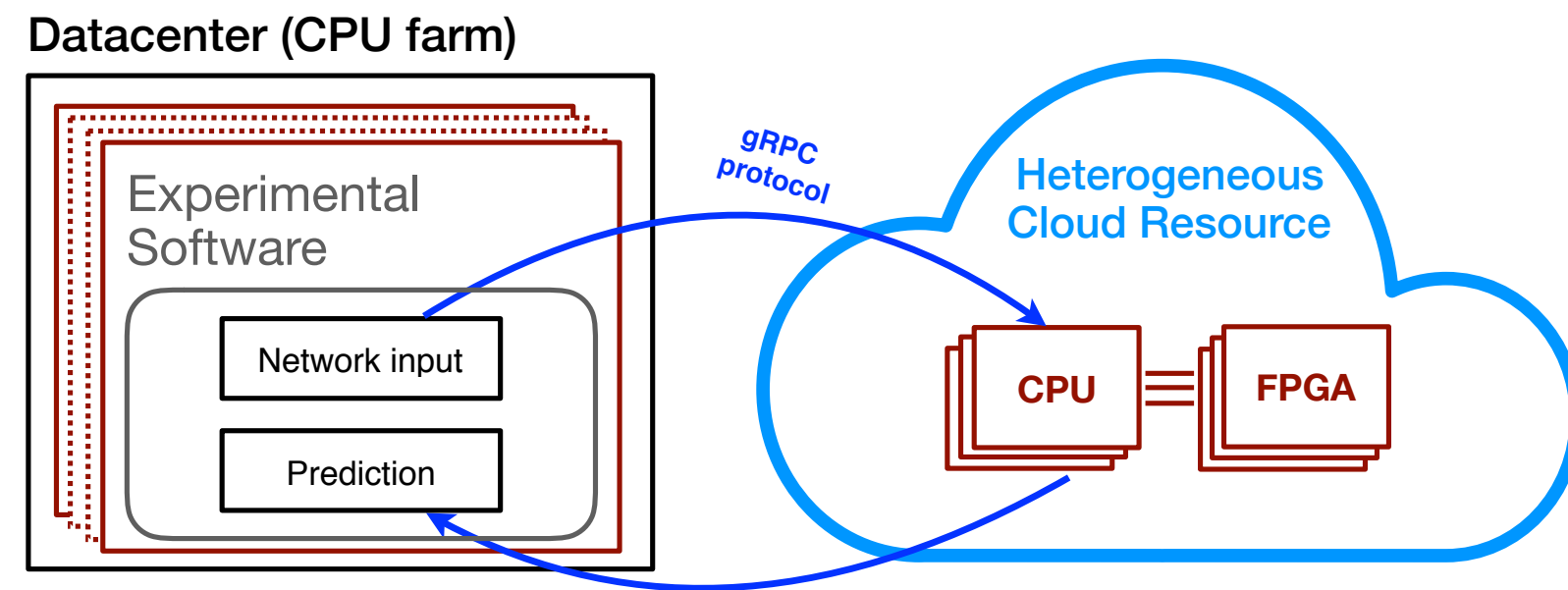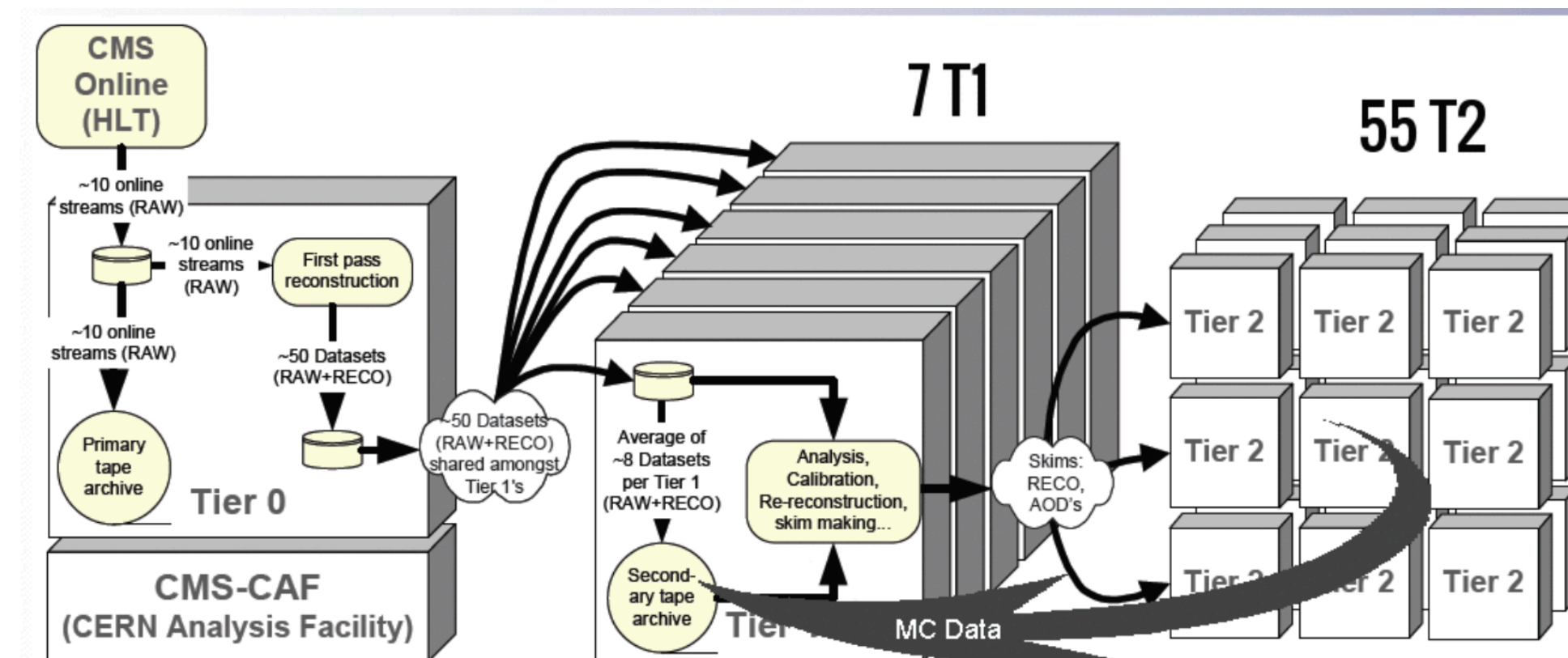
NVIDIA
Triton Inference Server

**Graphcore**

# How to deploy SONIC in CMS

## Cloud vs. Edge

CPU farm

Heterogeneous Cloud Resource

gRPC

Datacenter (CPU farm)

Experimental Software

gRPC protocol

Heterogeneous Cloud Resource

Network input

CPU ═ FPGA

Heterogeneous "Edge" Resource

CPU

Experimental software

gRPC protocol

...oud service has latency

...n CMSSW on Azure cloud machine
...te local installation of FPGAs
..." or "edge")

...test of "HLT-like" performance

NVIDIA Triton Inference Server

CMS Online (HLT)

~10 online streams (RAW)

~10 online streams (RAW)

9 First pass reconstruction

~50 Datasets (RAW+RECO)

~10 online streams (RAW)

Primary tape archive

Tier 0

~50 Datasets (RAW+RECO) shared amongst Tier 1's

CMS-CAF (CERN Analysis Facility)

Average of ~8 Datasets per Tier 1 (RAW+RECO)

Analysis, Calibration, Re-reconstruction, skim making...

Second-ary tape archive

7 T1

Skims: RECO, AOD's

MC Data

55 T2

Tier 2  Tier 2  Tier 2

Tier 2  Tier 2  Tier 2

Tier 2  Tier 2  Tier 2

T2 (various)

T1:USA

T1:Italy  INFN

T1:France  IN2P3

T1:Germany

T0  CERN

T1:Spain  PIC

T1:UK  GridPP

T1:Taiwan  ASCC

57

# Alternative solution : as-a-service

Flexible - task-based optimization; software abstraction; low software maintenance overhead
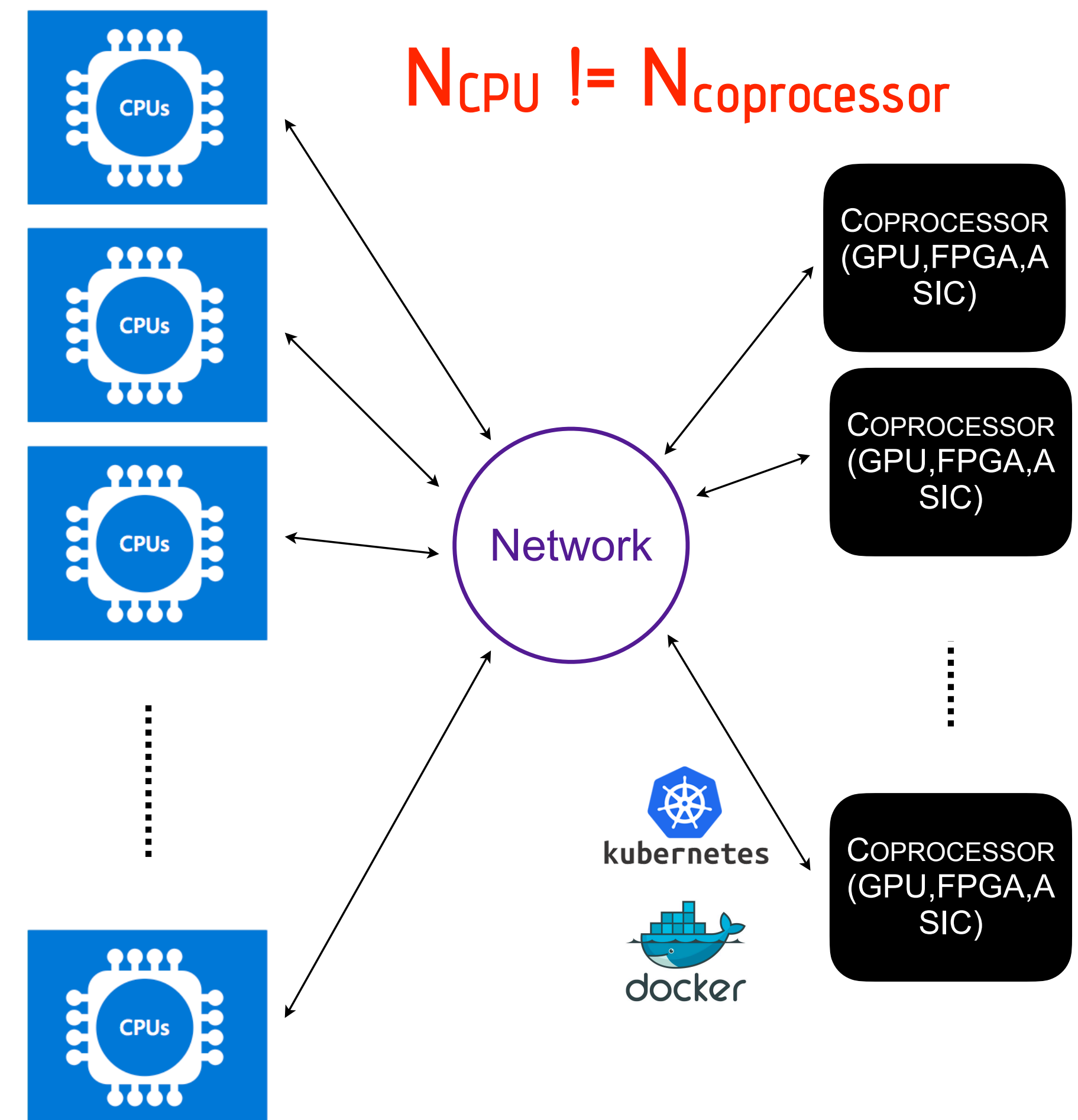
Adaptable - right-size the system based on compute needs, maximize e.g. GPU acceleration

Scalable - co-processor disassociated from existing CPU infrastructure; common software framework

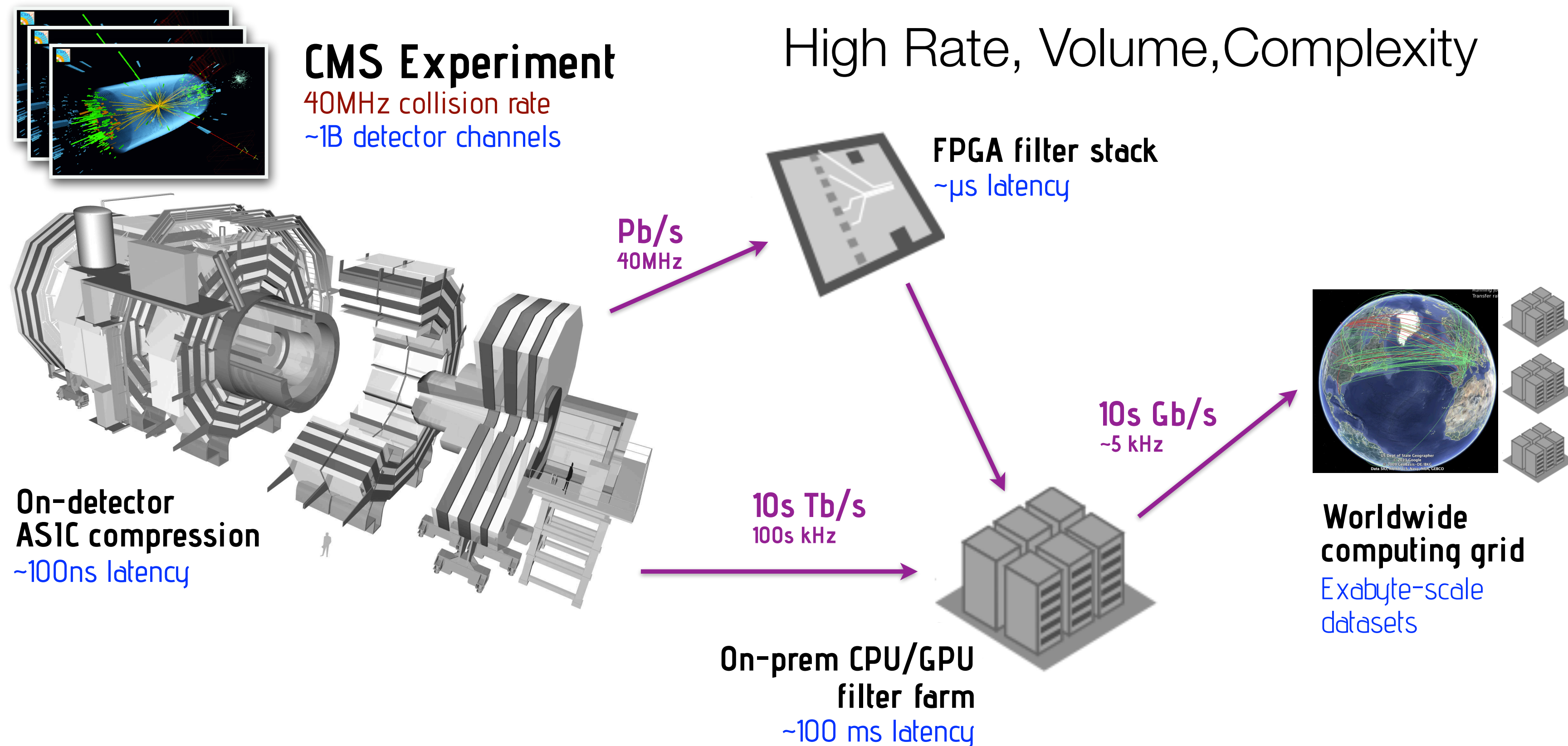Non-disruptive - maintain HEP computing paradigm, coprocessors as an enhancement

[First demonstration of integrating SONIC, with tests at Purdue CMS Tier-2 data center](#)

More details see [my talk at CPAD 2023](#).

$N_{CPU} != N_{coprocessor}$

# What else?

# On-Chip?



Fermilab

**CMS Experiment**
40MHz collision rate
~1B detector channels

High Rate, Volume, Complexity

**FPGA filter stack**
~µs latency

**Pb/s**
40MHz

**On-detector
ASIC compression**
~100ns latency

**10s Tb/s**
100s kHz

**10s Gb/s**
~5 kHz

**Worldwide
computing grid**
Exabyte-scale
datasets

**On-prem CPU/GPU
filter farm**
~100 ms latency

Science with Big data: Multi-tier Data Processing

# hardware trigger capabilities

42M strips on 192m²

ls on 25m²

Designed to cope with high data rate, high radiation environment at the HL-LHC

Higher granularity, Low material budget, titled geometry

$P_T$ modules

1.6-4.0mm
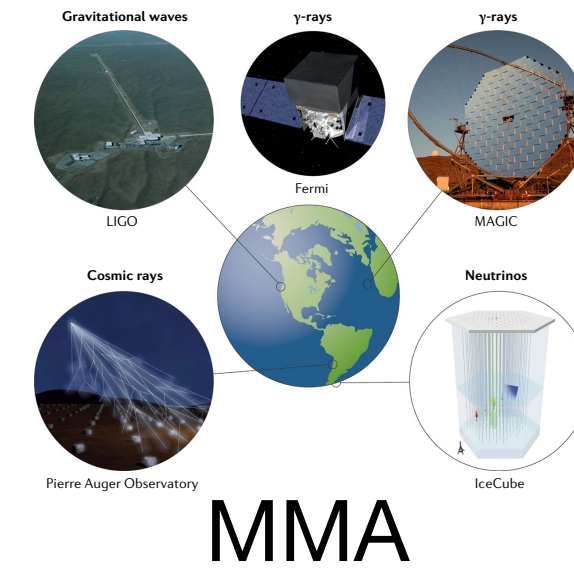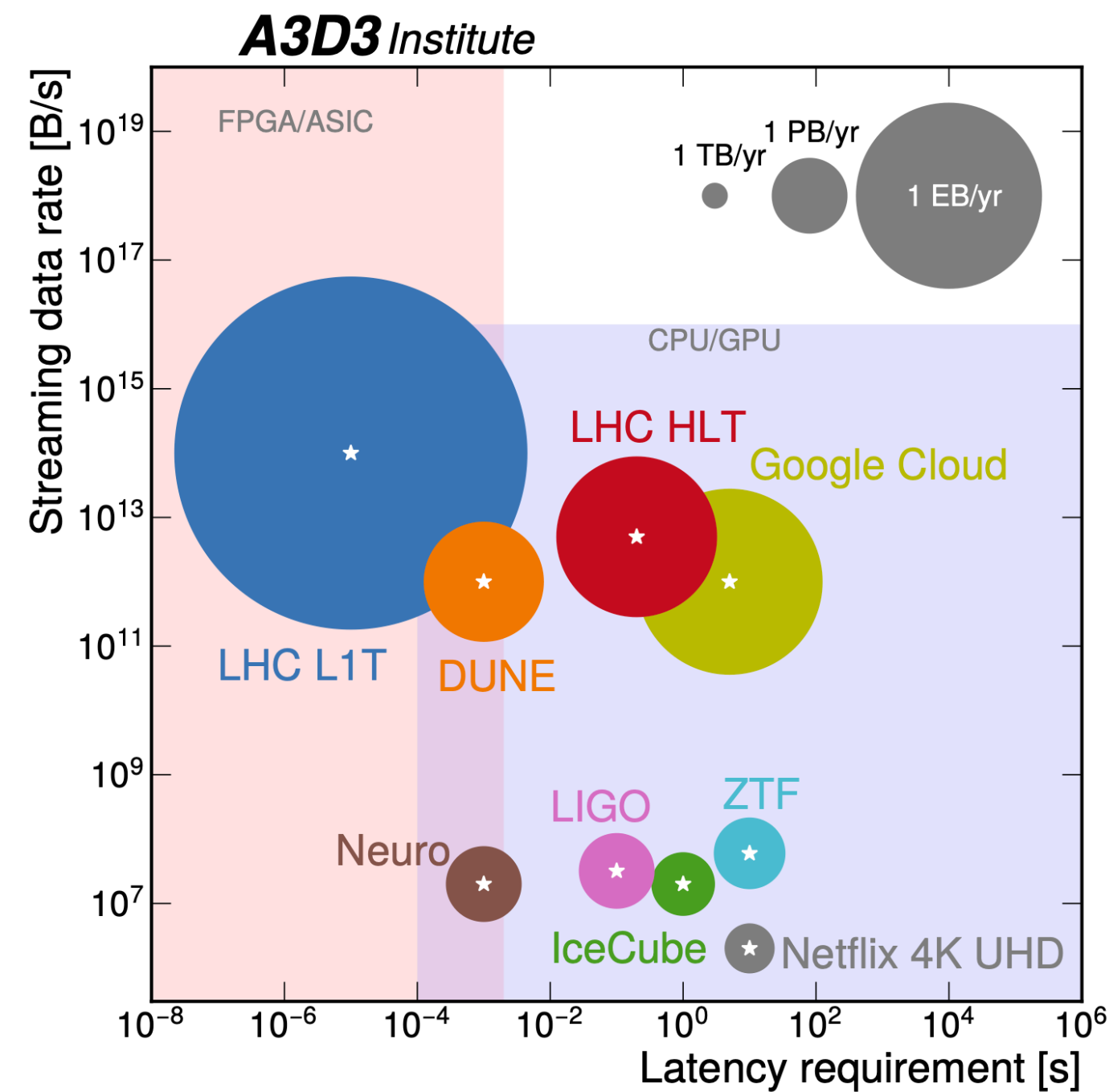
pass

programmable search window

fail

$\vec{B}$ ⊙

stub

Readout → Track Find → CMS Level-1

CMS DAQ

Tracker Back-end

CMS

Tracks per Event / 1GeV/c

$p_T$ [GeV/c]

# 'Pt modules' for Pixels?



**CMS Experiment**
40MHz collision rate
~1B detector channels

**FPGA filter stack**
~µs latency

Pb/s
40MHz

10s Gb/s
~5 kHz

**On-detector
ASIC compression**
~100ns latency

10s Tb/s
100s kHz

**Worldwide
computing grid**
Exabyte-scale
datasets

**On-prem CPU/GPU
filter farm**
~100 ms latency

Enabled by HLS4ML catapult

Use cluster information to infer particle kinematics

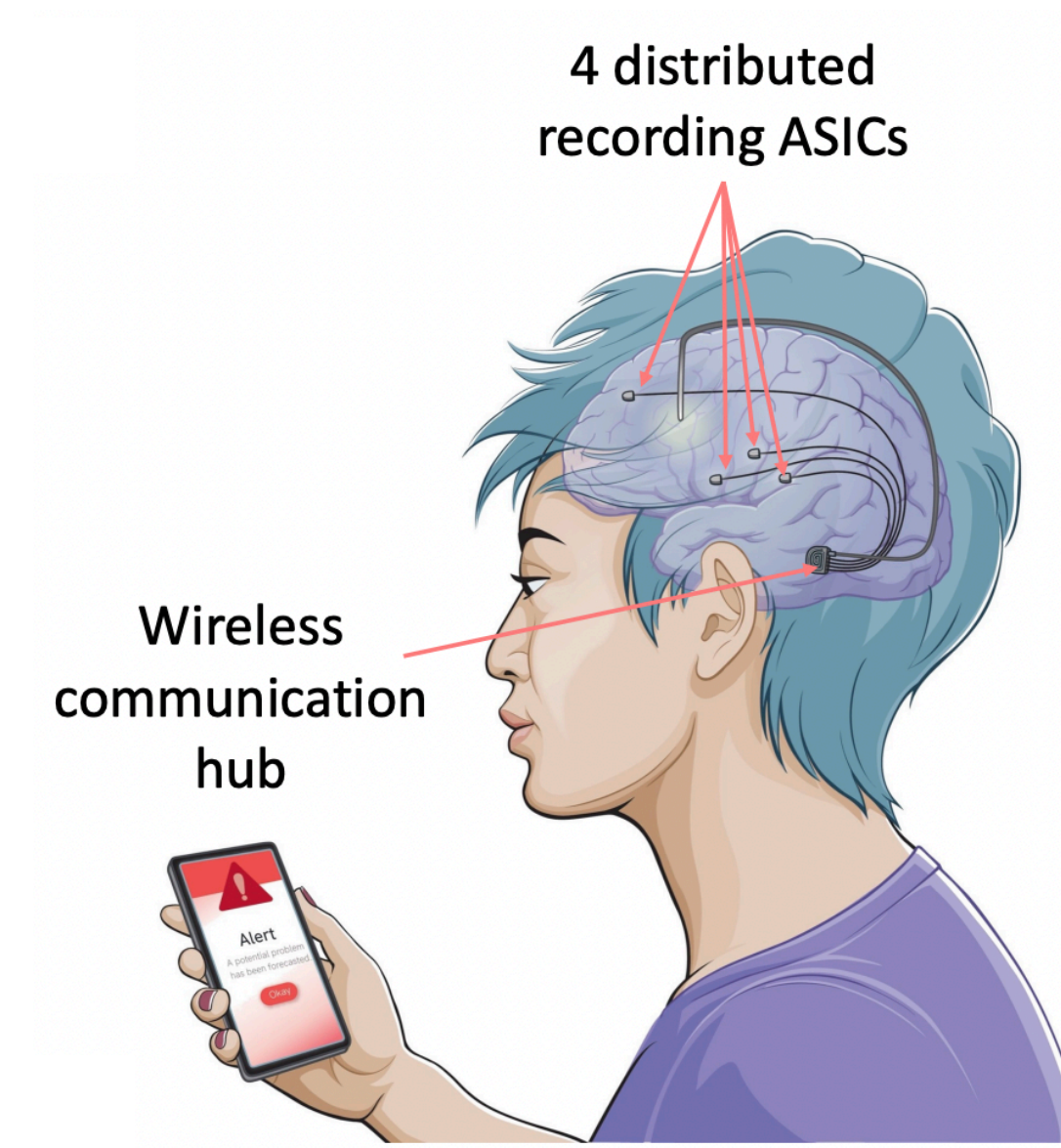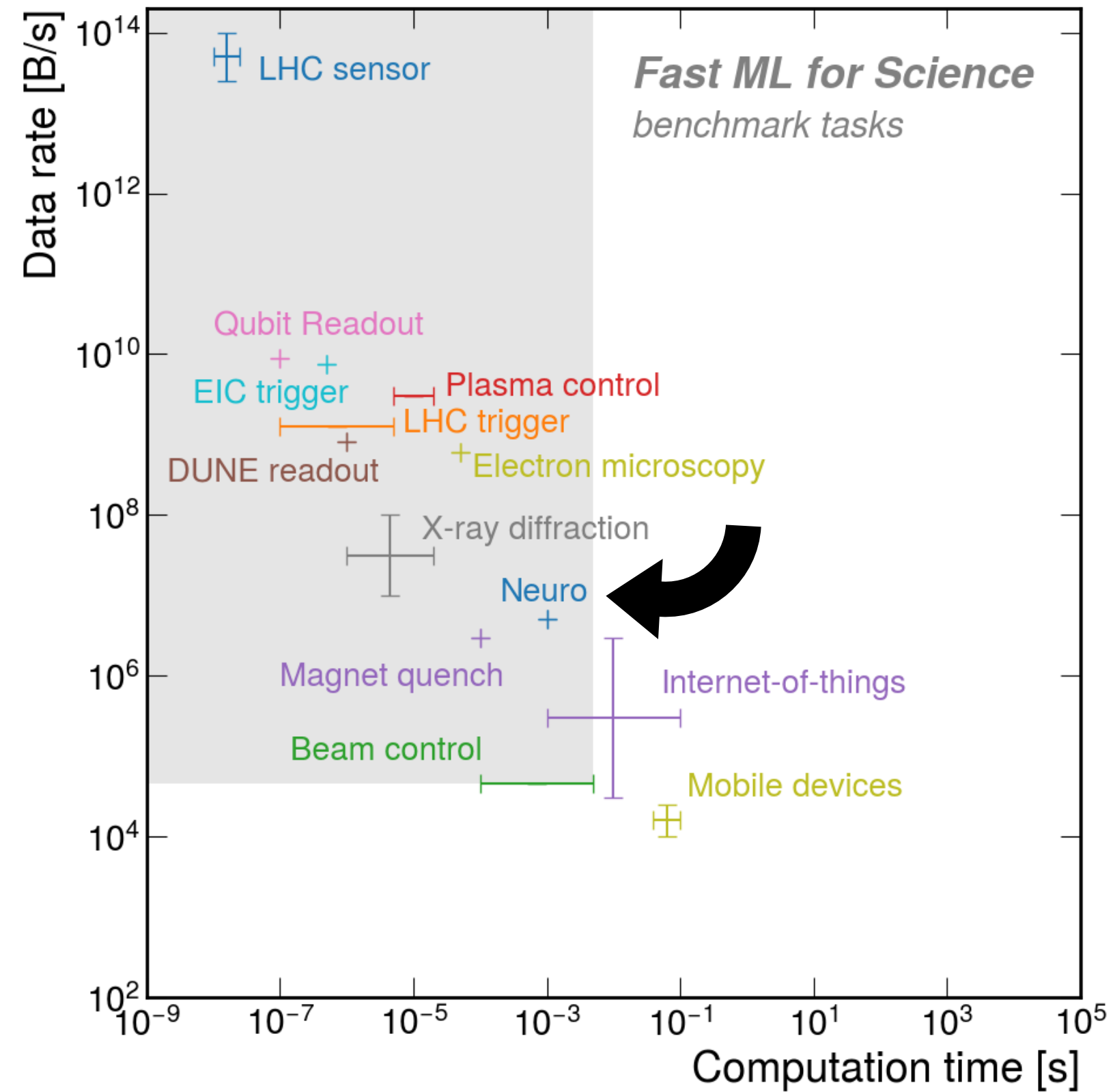# Accelerated AI Opportunities



NSF A3D3 institute: Domain Scientists, Computer Scientists and System Experts
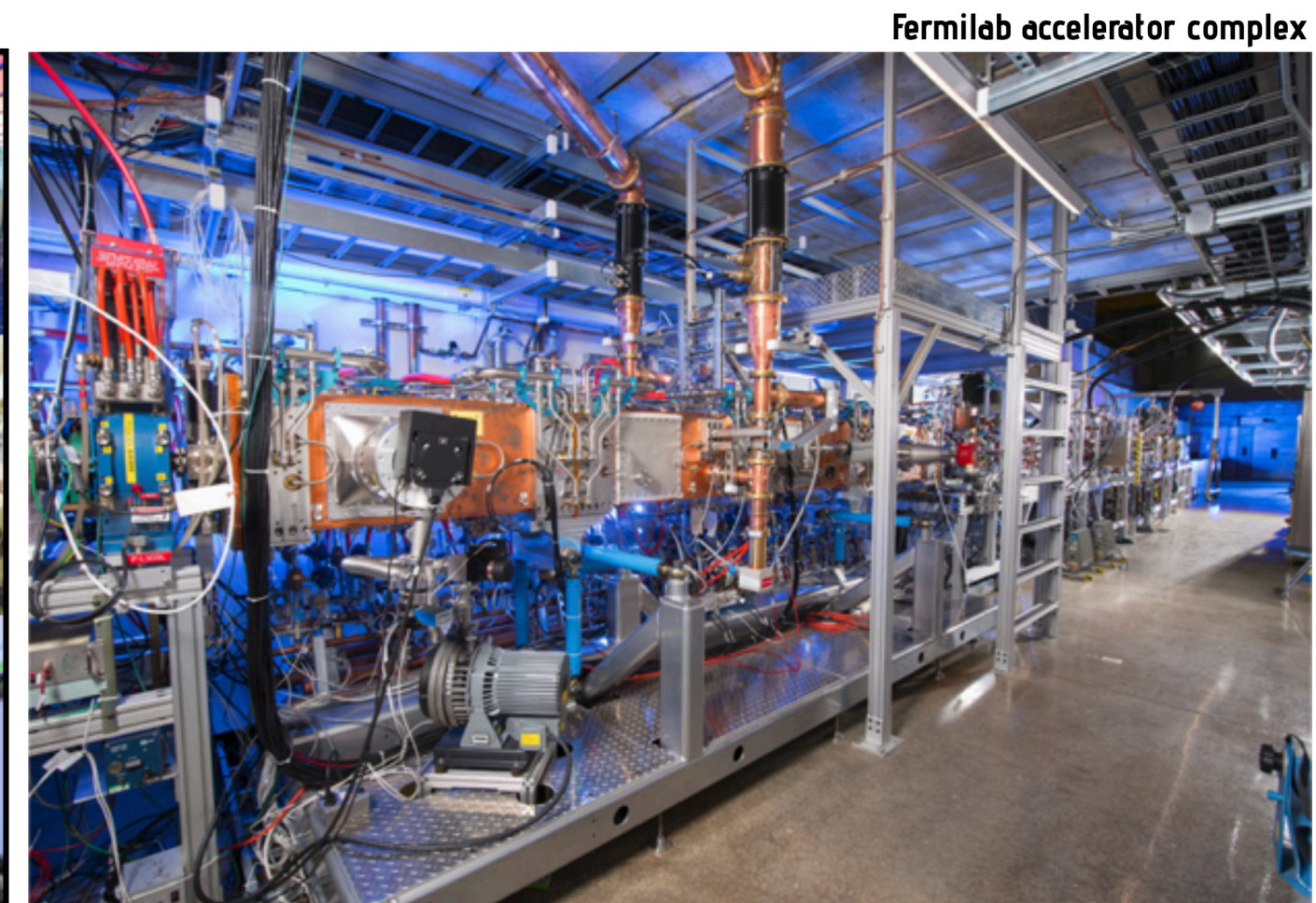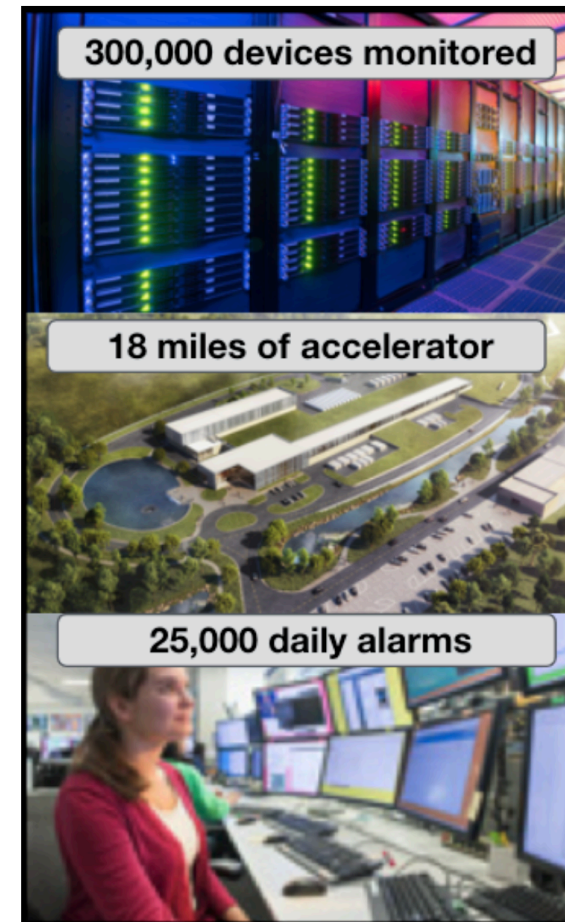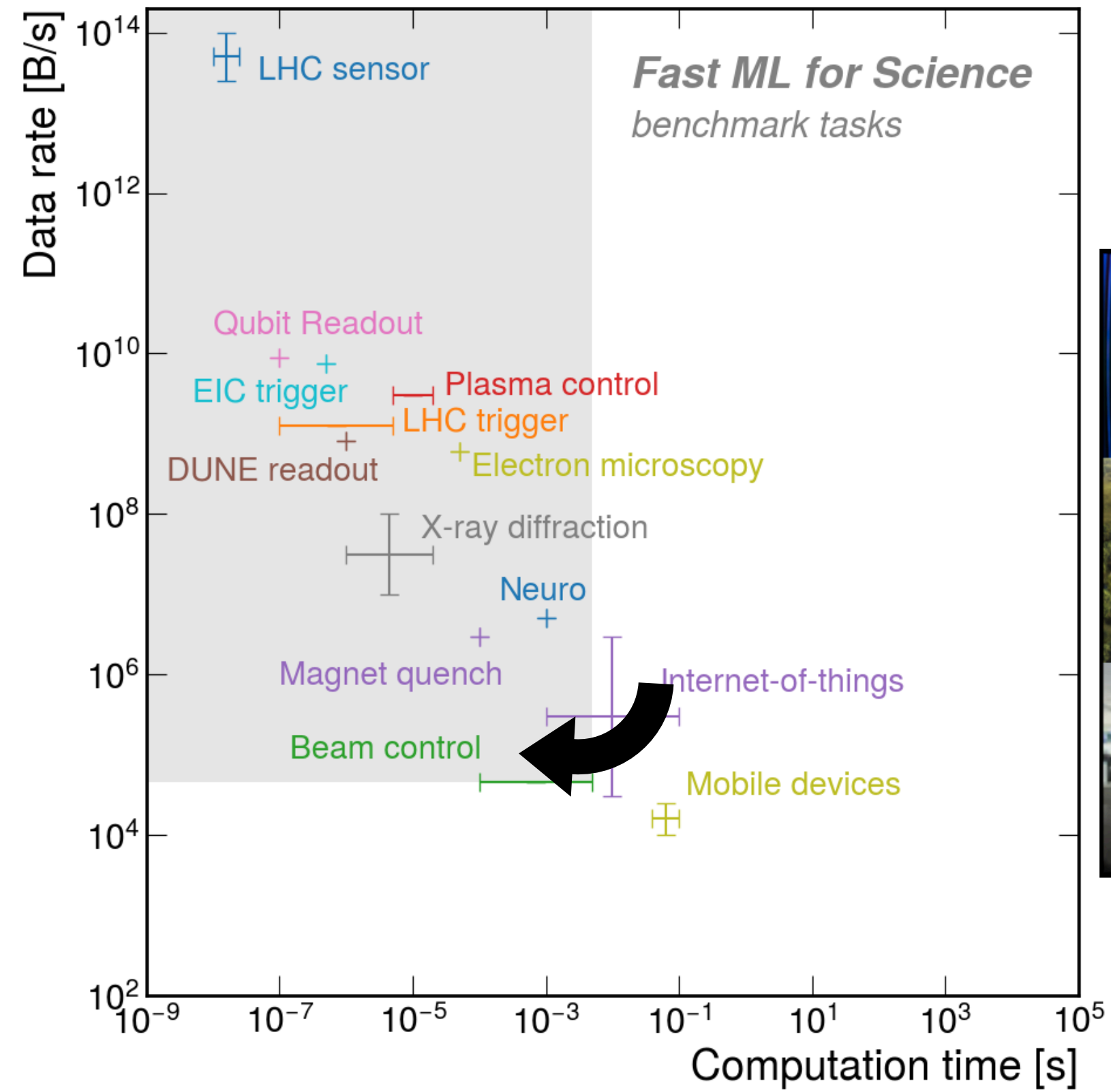Impact broader science domains Fast ML for Science Workshop

# Fast ML for Science



**Fast ML for Science**
benchmark tasks

Data rate [B/s] vs Computation time [s]

LHC sensor, Qubit Readout, EIC trigger, Plasma control, LHC trigger, DUNE readout, Electron microscopy, X-ray diffraction, Neuro, Magnet quench, Internet-of-things, Beam control, Mobile devices



Credit: Michael Coughlin

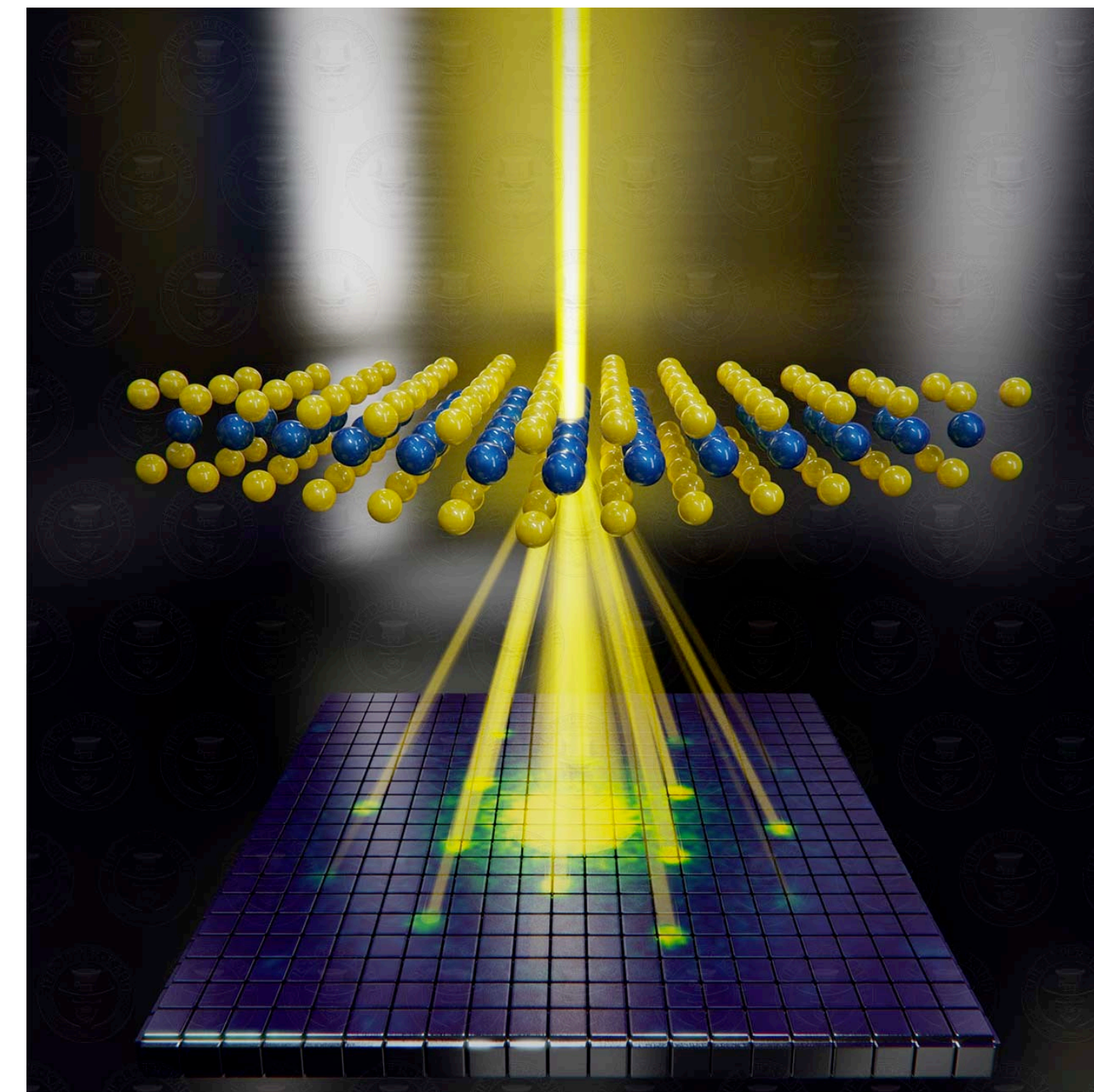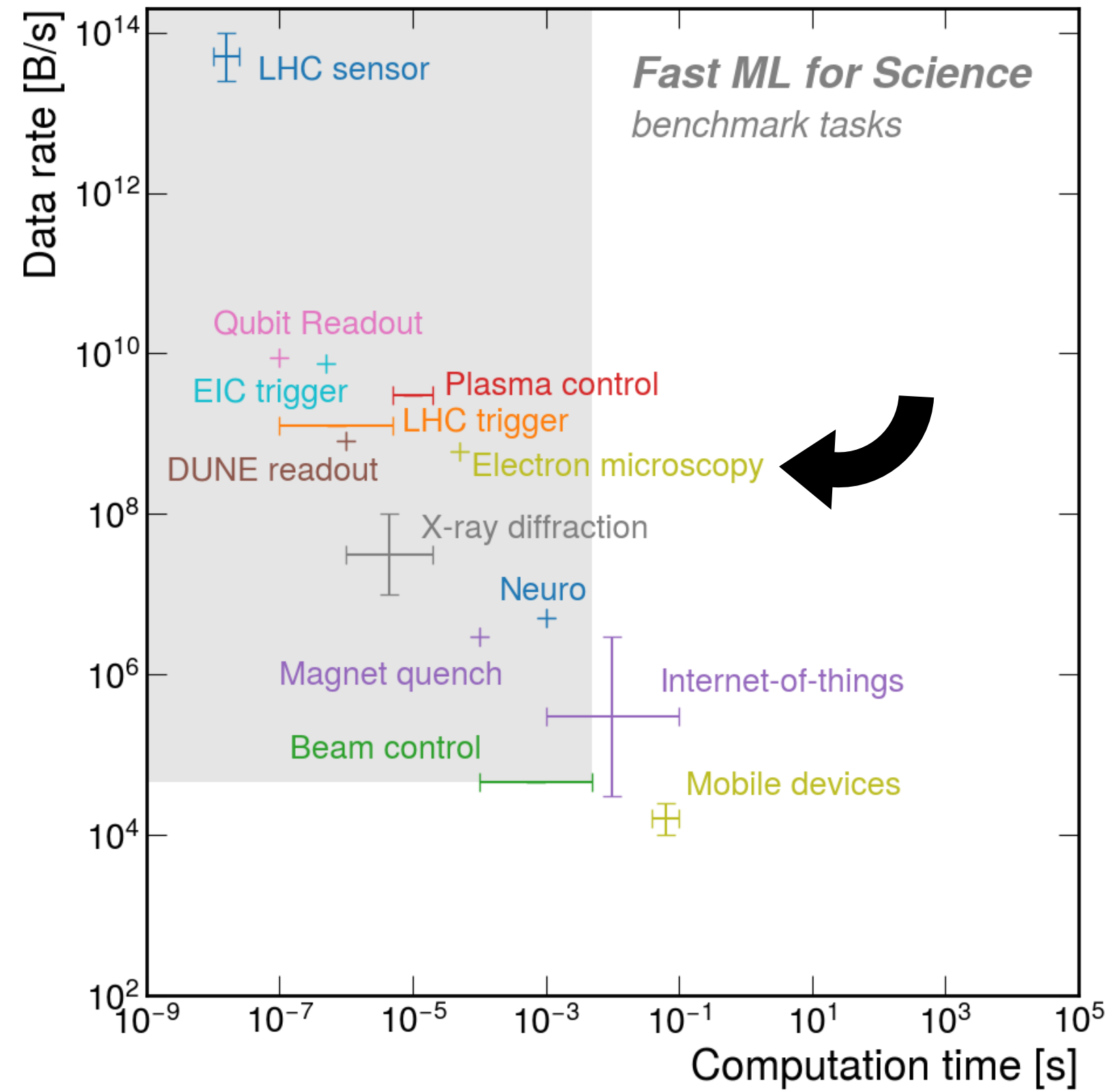**Supernova detection and multi-messenger astronomy**

64

# Fast ML for Science



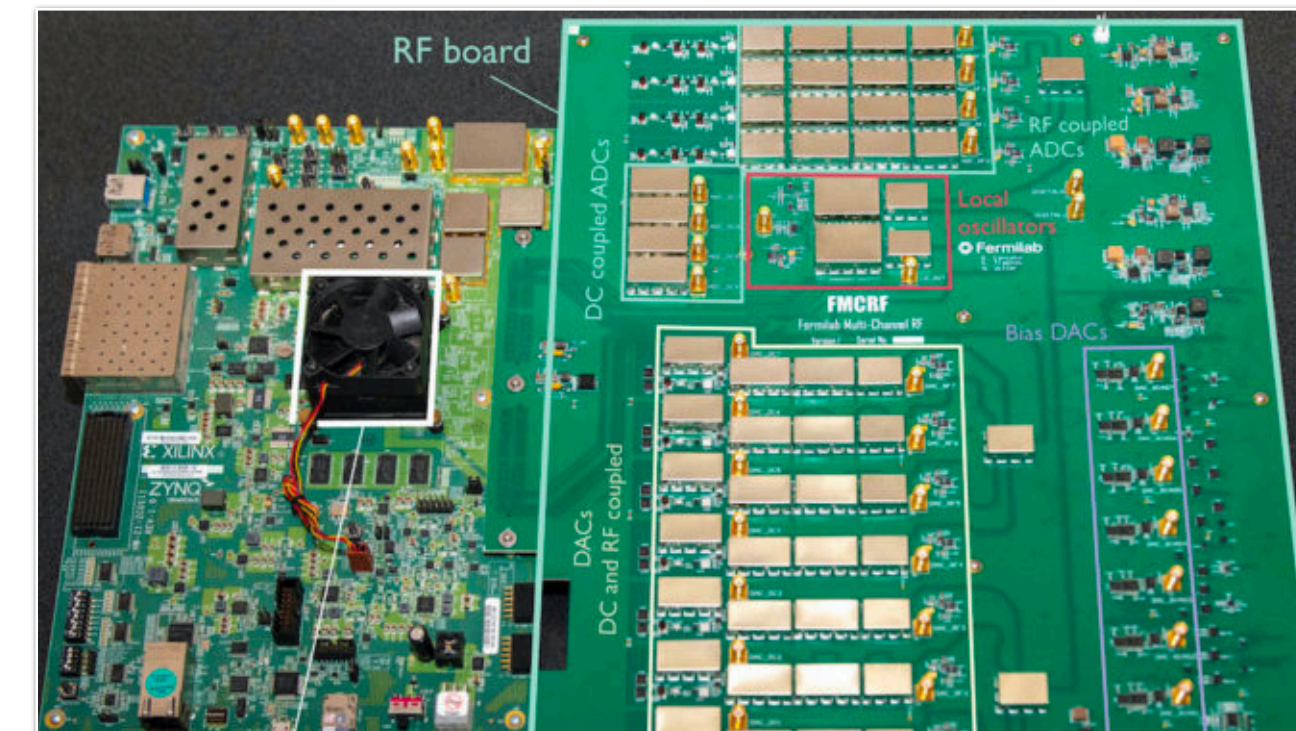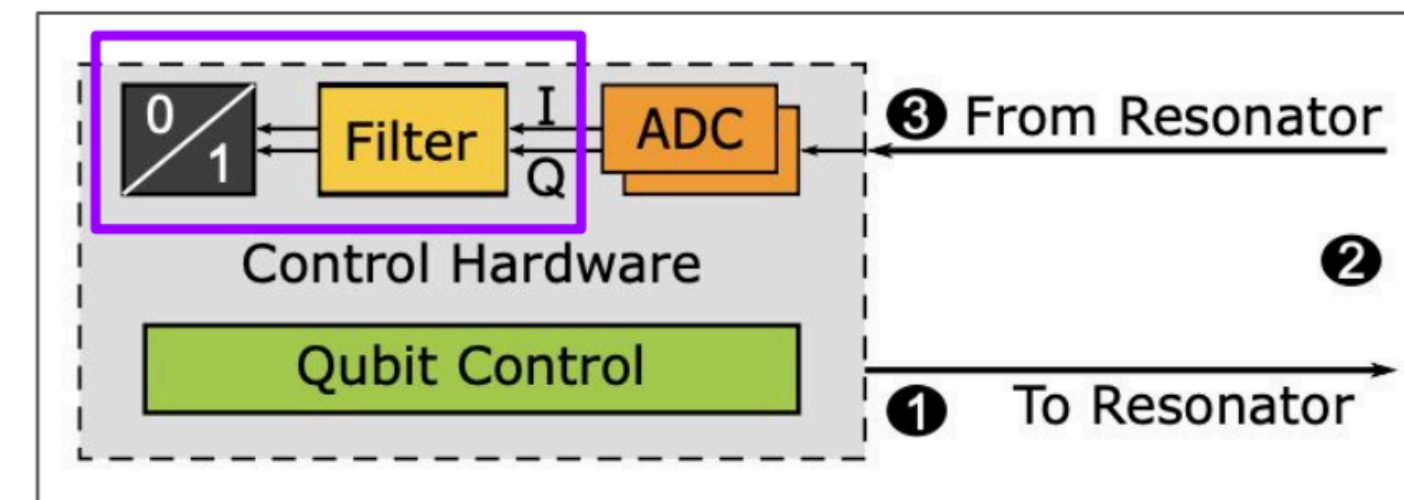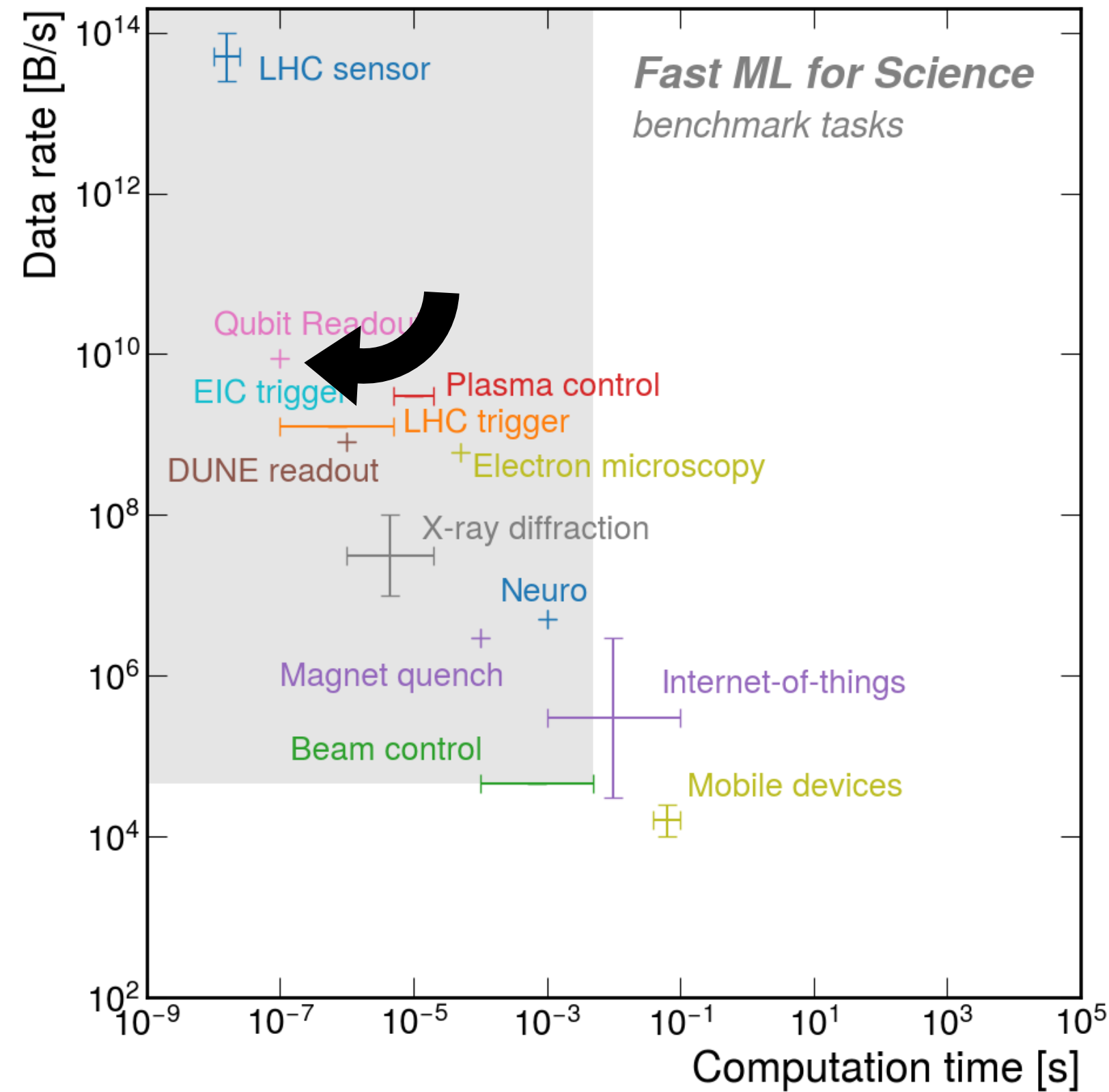**Real-time seizure detection**

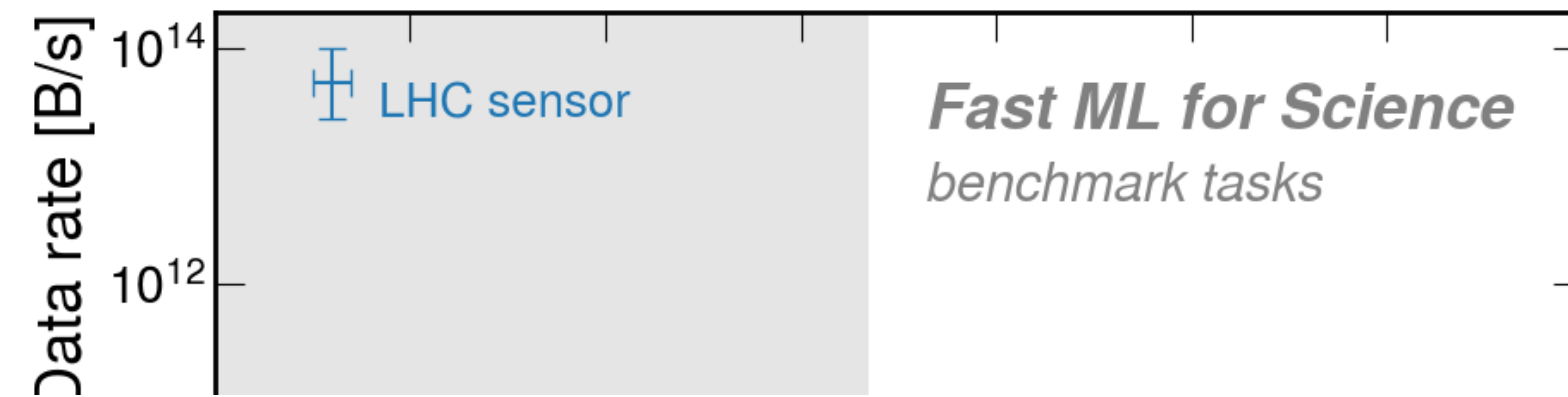# Fast ML for Science



**Particle accelerator controls**

# Fast ML for Science



**New materials for quantum and energy**

# Fast ML for Science



**Qubit readout and control**

# Fast ML for Science



See many more at the fast machine learning for science workshop2023





**Qubit readout and control**

# Final Remark

- Many existing and emerging opportunities in advancing our science results with Real-time ML/AI

  - New opportunities in searching for new physics

- Interesting research area touching overlaps of CS/AI, engineering and domain problems

- Lots of Fun



GNN based tagger can collect 5 times more new physics with exotic footprint in CMS detector

# Image detection network evolution

# Performance



Attention maps of performant ImageNet models (ResNets) are similar to each other, but the less performant model (NIN) has quite different attention maps.

Input

NIN (62% acc)

ResNet101 (77% acc)

$F_{\text{sum}}$  $F^2_{\text{sum}}$  $F^4_{\text{sum}}$  $F^2_{\text{max}}$  $F_{\text{sum}}$  $F^2_{\text{sum}}$  $F^4_{\text{sum}}$  $F^2_{\text{max}}$

# Towards Scalable, Flexible, Adaptable GNN/ transformer with HLS4ML

- **hls4ml: great support for MLP and CNN Keras models.**

- Support of parsing PyTorch models: this has been improved!
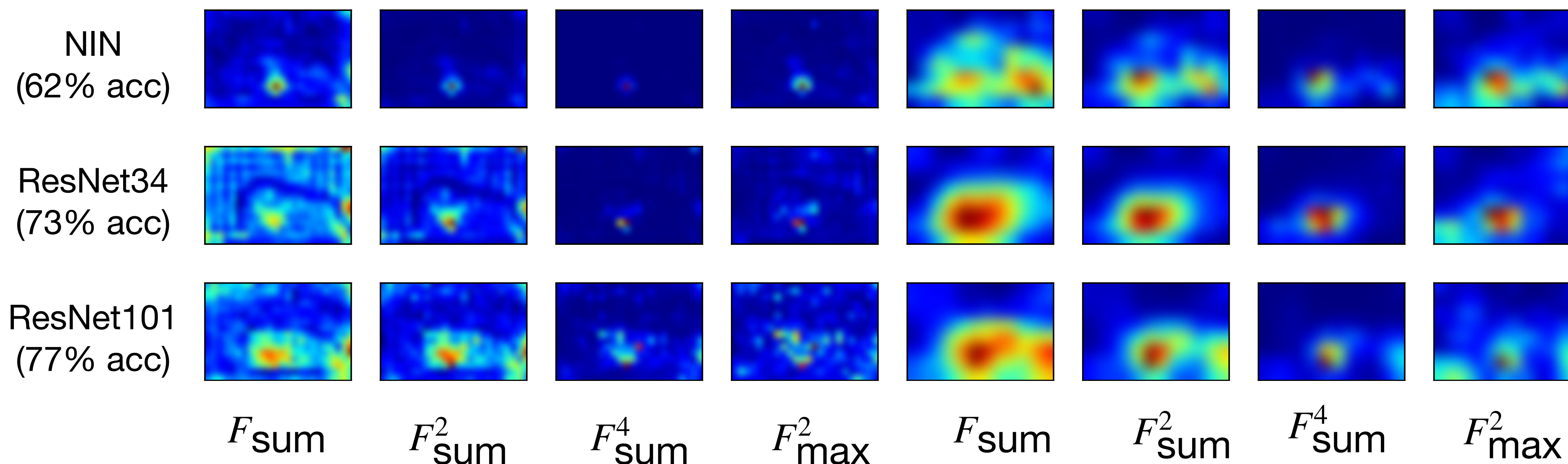
- **Some (non-trivial) engineering work to support GNN/transformers:**

- Tau3mu Detection: MessagePassing layers, and meet 100 ns latency!

- **Long term: need to improve hls4ml code generation**

  - Current code generation in hls4ml is based on naive string generation - i.e., it becomes a mess very fast for anything complex.

sPHENIX tracking GNN hls4ml synthesis results

- Network inputs: nodes=80, edges=100  **Extremely preliminary - DO NOT TRUST NUMBERS**
- Input network
  - Can be parallelized to be "nodes" times faster (i.e., 15ns)
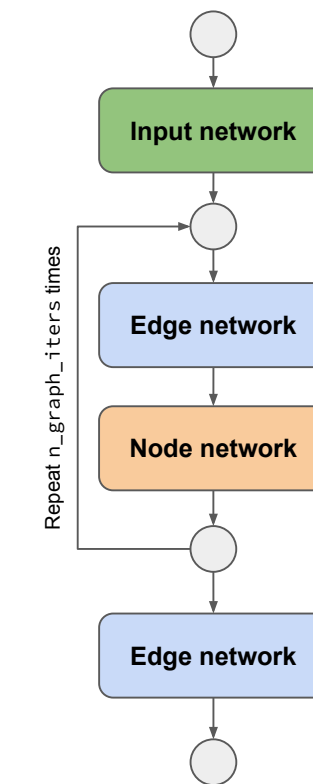
| Latency | BRAMs | DSPs | FFs | LUTs |
|---------|-------|------|-----|------|
| 1.2 us  | 6.5%  | 0.3% | 5%  | 7.5% |

- Edge network

| Latency | BRAMs | DSPs | FFs | LUTs |
|---------|-------|------|-----|------|
| 3 us    | 15%   | 2%   | 20% | 65%  |

- Node network (results from HLS synthesis, vivado synthesis OOM'd)
  - Neet to optimize the scatter_add function (expecting ~2us for the net)

| Latency | BRAMs | DSPs | FFs | LUTs |
|---------|-------|------|-----|------|
| 12 us   | 42%   | 7%   | -   | -    |

Input network

Edge network

Node network

Edge network

Repeat n_graph_iters times

**Example:** Extended operations supported in hls4ml to implement a GNN developed for track reconstruction in the sPhenix trigger

- Added missing operations for GNN: Scatter_* "getitem", "gather", "ones()" and "zeros()" etc

# A computing paradigm adaptive to changing hardware landscape

## Cloud vs. Edge

**CPU farm**

**Heterogeneous Cloud Resource**

**gRPC**

Datacenter (CPU farm)



**Microsoft** Catapult/Brainwave `FPGA`

Experimental Software

gRPC protocol

Heterogeneous Cloud Resource

Network input

CPU — FPGA

ed co-processor hardware
ine learning inference

`ASIC?`

1 Bionic neural engine

**S**ervices for **O**ptimize

https://arxiv.org/pdf/

Heterogeneous "Edge" Resource

CPU

Experimental software

gRPC protocol

FPGA

- Increasing demar
- Demonstrated off
  - FPGA co-proce
- CPU client softwa

ud service has latency

n CMSSW on Azure cloud machine
simulate local installation of FPGAs
n-prem" or "edge")

ovides test of "HLT-like" performance
wave FPGA services

vant for streaming)

not inference framework

`ASIC`

9