# Summary of Fast ML for Science Workshop

https://indico.cern.ch/event/924283/

Shih-Chieh Hsu
University of Washington
Dec 16 2020
HSF Trigger Reco Working Group

# Fast Machine Learning For Science 2020



History
- ML for LHC (Apr 2018) at MIT [URL]
- UltraFast DNN (Feb 2019) Zurich [URL]
- 1st Edition (Sep 2019) at FNAL [URL]
- 2nd Edition (30 Nov - 3 Dec, 2020) Virtual (this workshop)
  - Allison Deiana (LOC chair)
  - 3 days of workshops (581 registrations)
    - 47 plenary talks
  - 1 day of FastML tutorials (60 participants)

Mini-workshop on Portable Inference (Dec 4 2020)
- An IRIS-HEP Blueprint Workshop (32 participants)
- Mark Neubauer (Host)
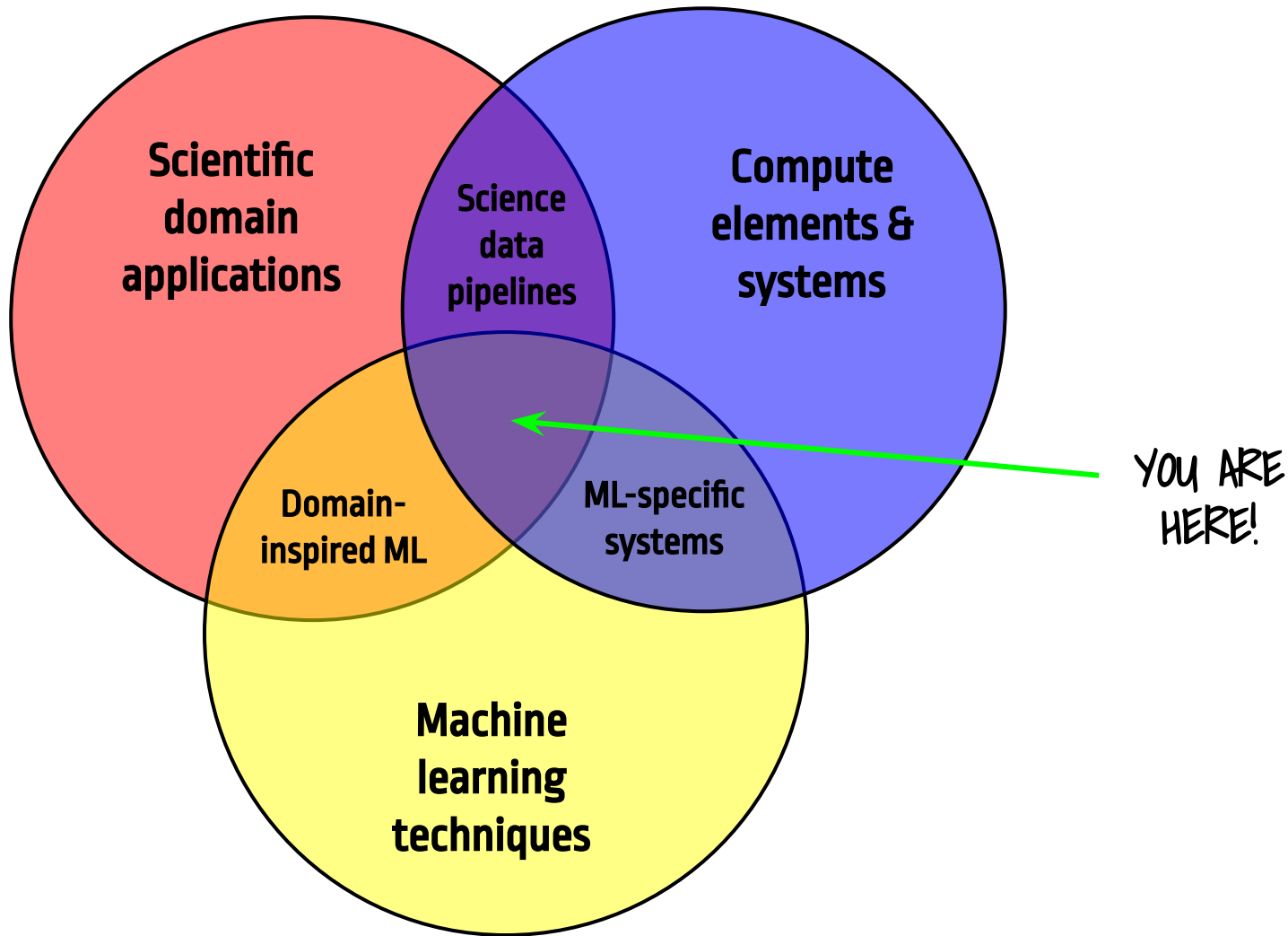- https://indico.cern.ch/event/972791

Nhan Tran

# Why we're here...

Cutting edge science requires:

**Faster, more precise, bigger, more granular,...**

Real-time, accelerated machine learning can greatly accelerate *time to science*, allowing us to:

‣ test hypotheses significantly faster
‣ enhance and automate performance of detectors/accelerators
‣ save and maximize potentially lost data

A Multi-Disciplinary Workshop

Scientific domain applications

Compute elements & systems

Science data pipelines

Domain-inspired ML

ML-specific systems

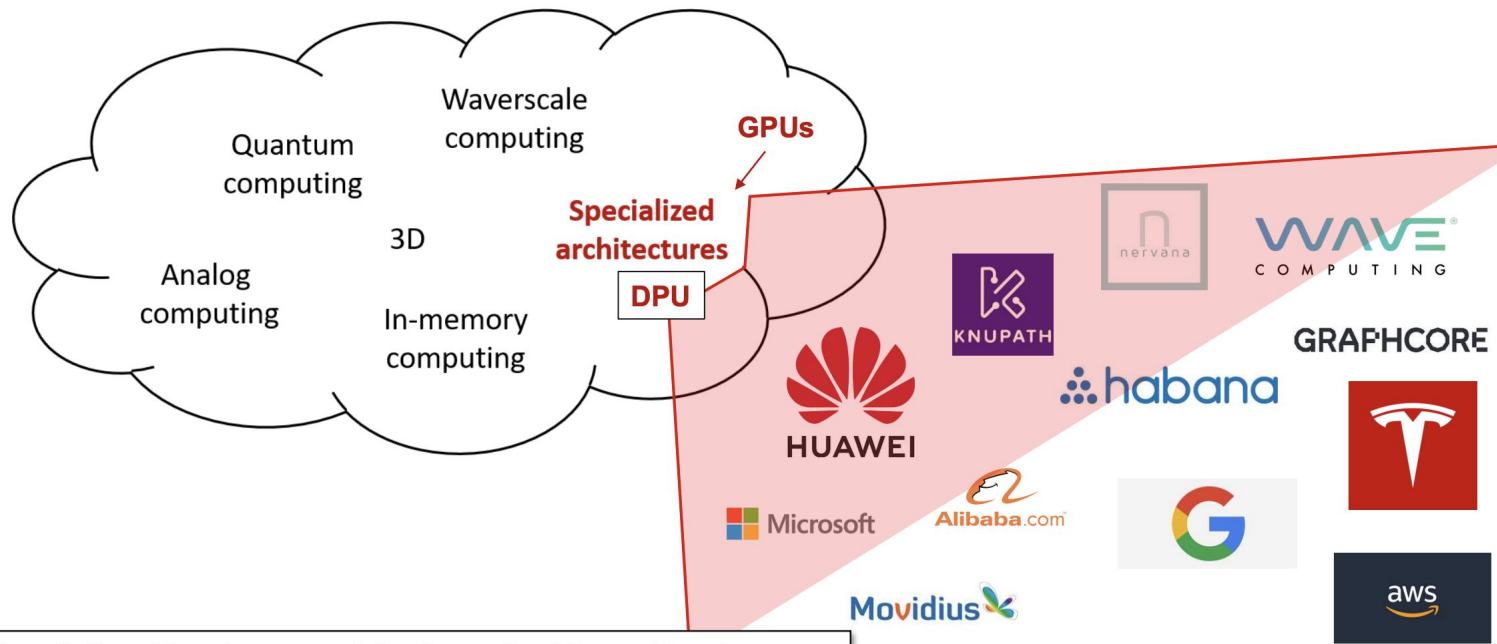Machine learning techniques

YOU ARE HERE!

# Outline

- The FastML workshop has a good mix between coprocessor and low latency
- This talk will highlight coprocessor specific contents
  - Advances in heterogeneous computing
    - Compute trends and tools
  - Applications for accelerated ML
    - Challenges in LHC physics
    - Challenges in adjacent science domains
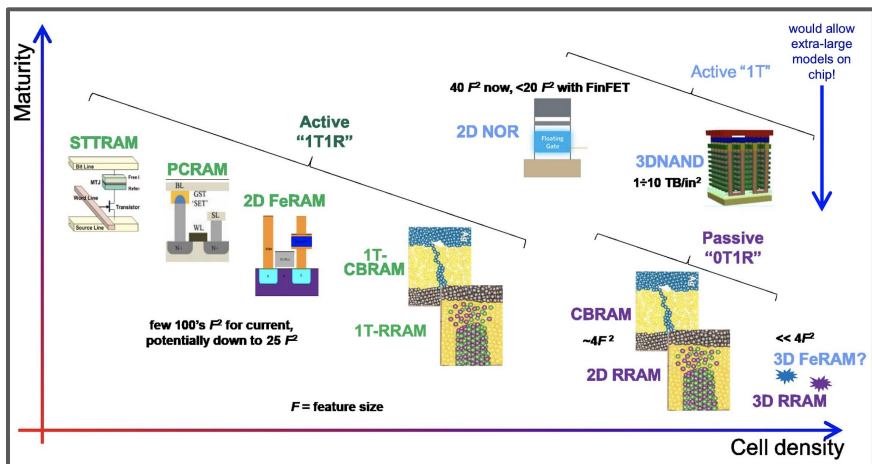
# Trends in computing architectures

Michaela Blott

**Limits of semiconductor chip development** Moore's law, Denard scaling, High manufacturing cost of 3 nm/ Trends driven by big data explosion & computationally expensive ML methods.
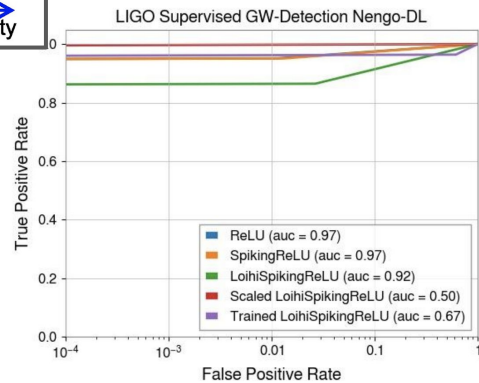


**Specialized** hardware architectures for ML workloads

# More types of hardware



Beyond CMOS tech maturity, Dmitri Strukov

First exploration on neuromorphic chips,
Bartlomiej Borzyszkowski, Eric Moreno



LIGO Supervised GW-Detection Nengo-DL

- ReLU (auc = 0.97)
- SpikingReLU (auc = 0.97)
- LoihiSpikingReLU (auc = 0.92)
- Scaled LoihiSpikingReLU (auc = 0.50)
- Trained LoihiSpikingReLU (auc = 0.67)

| Layer (type) | Output Shape | Param # | |
|---|---|---|---|
| input (InputLayer) | [(None, 2048)] | 0 | |
| reshape (Reshape) | (None, 2048, 1, 1) | 0 | |
| conv2d (Conv2D) | (None, 2045, 1, 16) | 64 | → Off-chip |
| conv2d_1 (Conv2D) | (None, 511, 1, 16) | 1024 | |
| conv2d_2 (Conv2D) | (None, 127, 1, 32) | 2048 | |
| conv2d_3 (Conv2D) | (None, 31, 1, 64) | 8192 | |
| conv2d_4 (Conv2D) | (None, 6, 1, 128) | 65536 | ON-chip |
| flatten (Flatten) | (None, 768) | 0 | |
| dense (Dense) | (None, 128) | 98304 | |
| dense_1 (Dense) | (None, 64) | 8192 | |
| output (Dense) | (None, 2) | 130 | |

Total params: 183,490
Trainable params: 183,490
Non-trainable params: 0
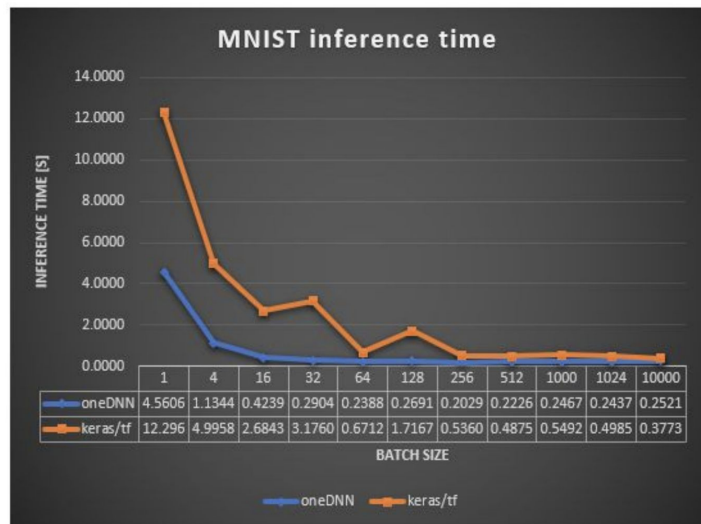
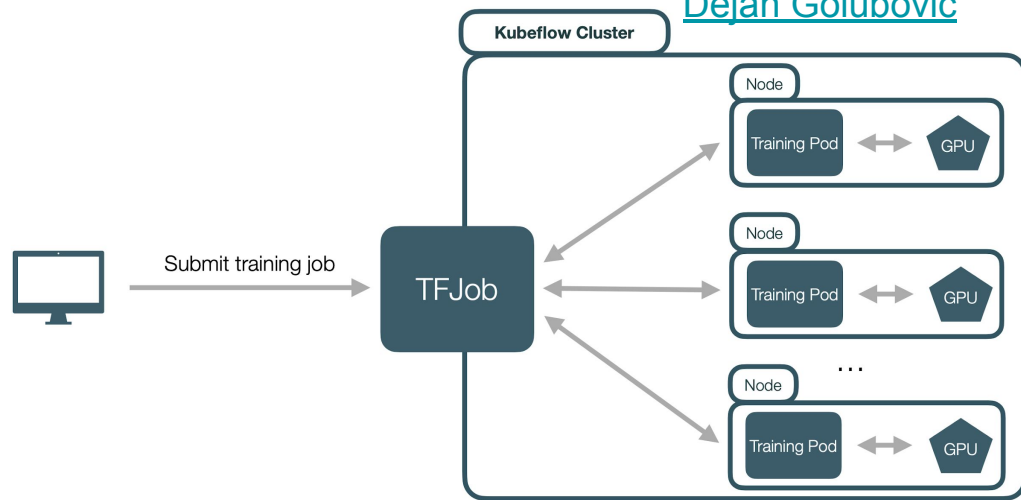Executed on Loihi chip

# Tools for rapid processing of ML algorithms

Intel OneAPI leads to fast CPU inference
Vladimir Loncar

Kubeflow tool for fast distributed training
Large scale control of the system
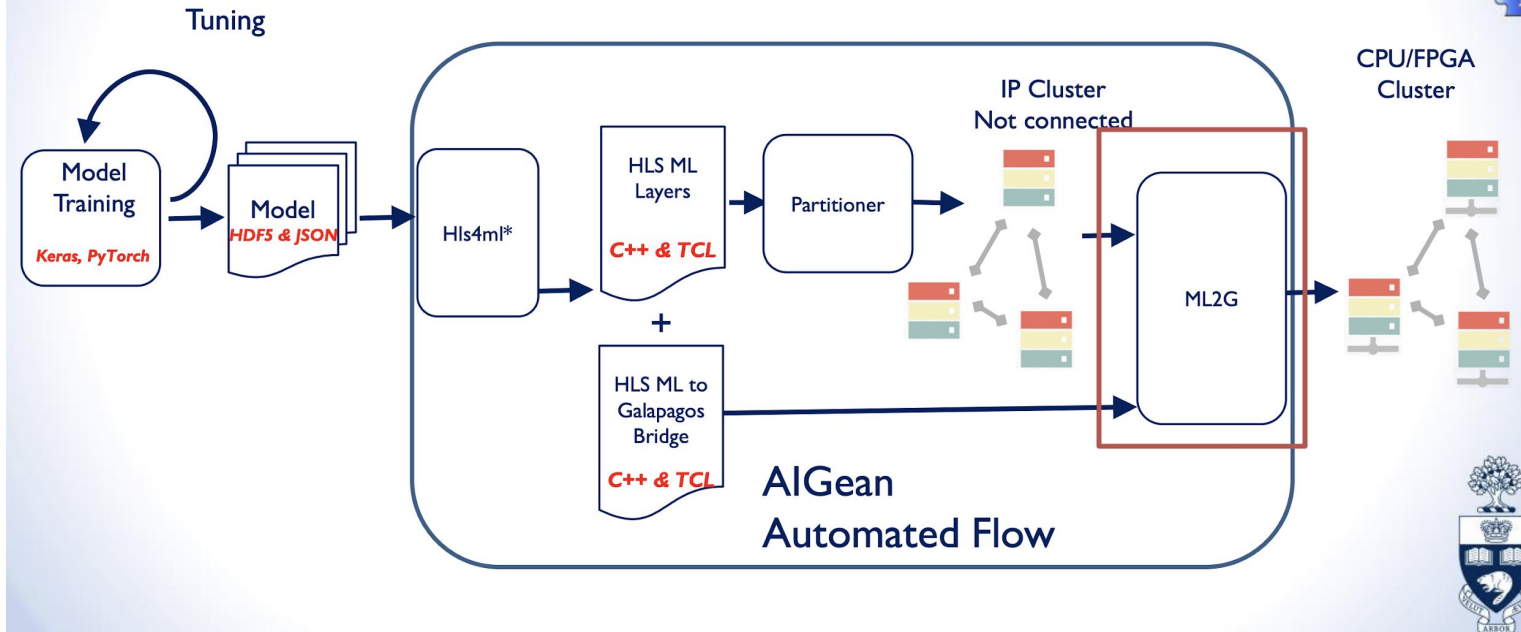deployed and available at CERN

Dejan Golubovic

# Tools for large distribution of networks on FPGAs

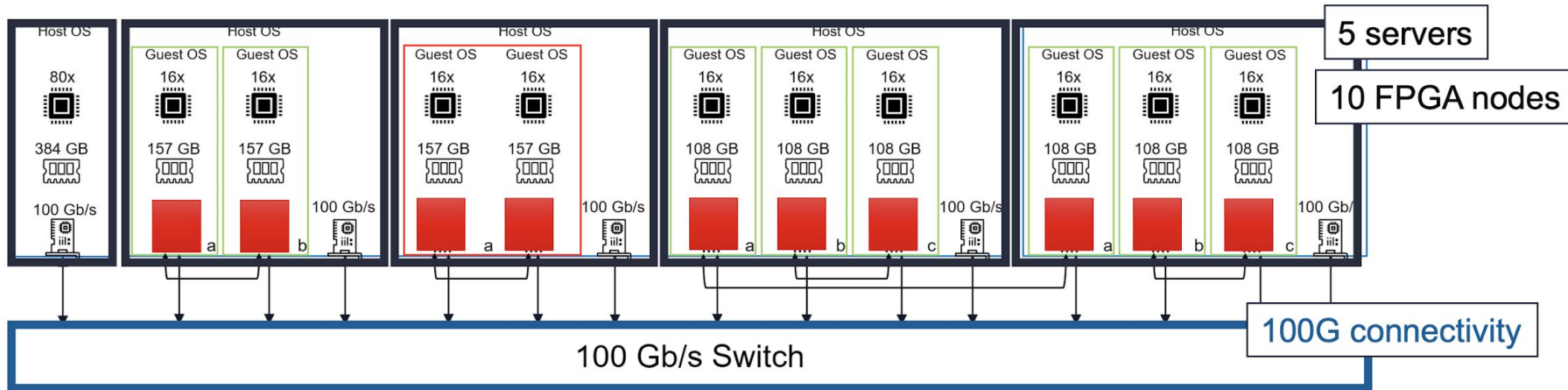# Tools for large distribution of networks on FPGAs

A large FPGA cluster to study algorithms across many FPGAS
FPGAs are linked with high speed connection
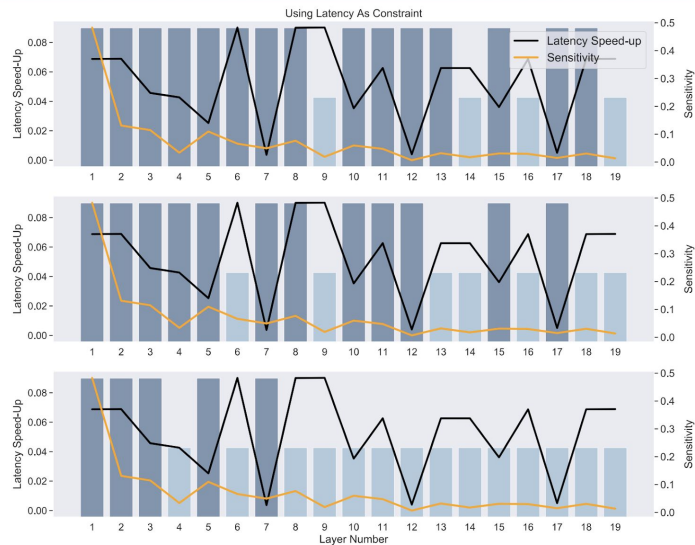Leads to a very different use of HPC cluster

© Copyright 2020 Xilinx

Michaela Blott, Xilinx

# Efficient Design at the Edge

- What optimizations in precision and design can lead to speed-ups and power reductions?
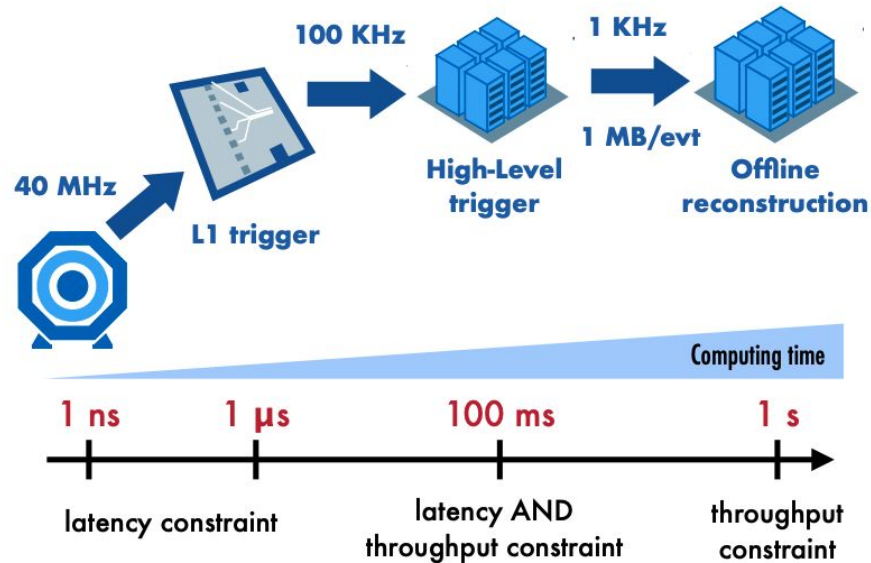
# Challenges in LHC Physics

# Offline Reconstruction (1 KHz, 1 s latency)

Jennifer Ngadiuba

# Offline Reconstruction (1 KHz, 1 s latency)



**1 KHz**

**1 MB/evt**

**Offline reconstruction**

## An example: tracking
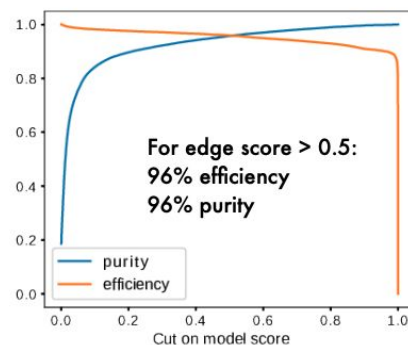
More in V. Razavimaleki talk
Wednesday @ 15:30

- New "faster" solutions being studied based on graph neural networks

- Graph construction:

  each tracker hit is a node
  edges built from geometrical consideration

- Edge classification:

  the GNN decides which edges connect
  hits from the same track

Grey edge: fake
Red edge : true

For edge score > 0.5:
96% efficiency
96% purity

purity
efficiency

Cut on model score

TrkX

**Exa.TrkX project**

NeurIPS 2019    CTD 2020

22

# High-Level Trigger (100 KHz, 100 ms latency)



**1 KHz**

**1 MB/evt**

**Offline reconstruction**

**100 KHz**
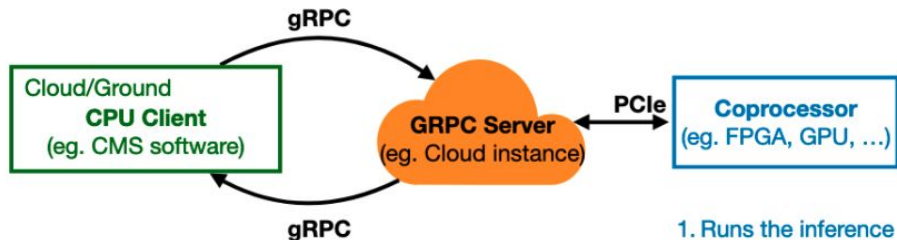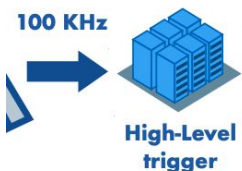
**High-Level trigger**

## MLaaS with SONIC

More in J. Krupa talk
Wednesday @ 15:40

- **Services for Optimized Network Inference on Coprocessors** (SONIC) enables inference as a service in experiment software frameworks

  - experiment software (C++) only has to handle converting inputs and outputs between event data format and inference server format
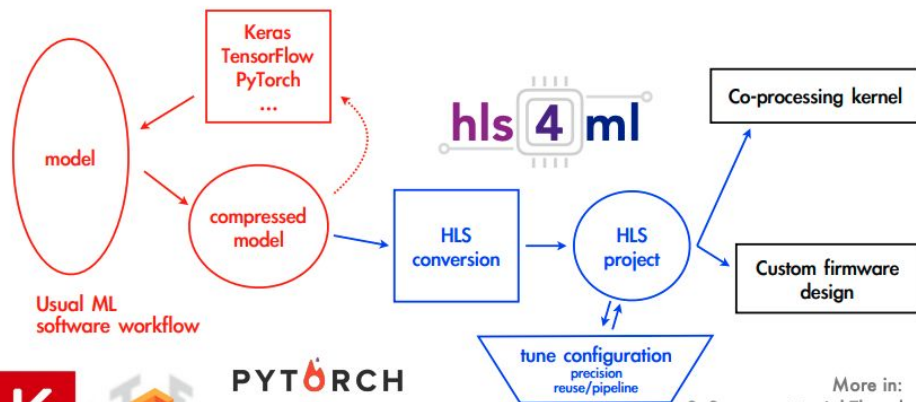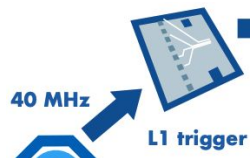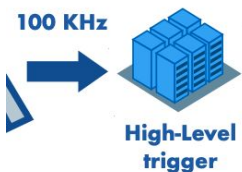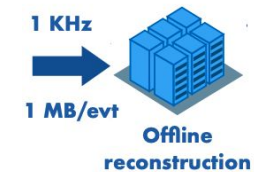
- Uses industry tools as gRPC communication and Nvidia Triton inference servers

- Interacts with cloud services: Azure, AWS, GCP

gRPC

Cloud/Ground
**CPU Client**
(eg. CMS software)

**GRPC Server**
(eg. Cloud instance)

PCIe

**Coprocessor**
(eg. FPGA, GPU, ...)

gRPC

1. Runs the inference

# Level 1 Trigger (40 MHz, 1 $\mu$s latency)



**high level synthesis for machine learning**

- A package for automatic translation of trained NN into HLS project
- Optimization for low-latency inference
- Very customazible: tunable precision and parallelization (area vs latency)
- Many architectures supported (dense NN, CNN, LSTM, Graph NN)
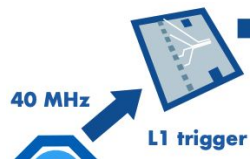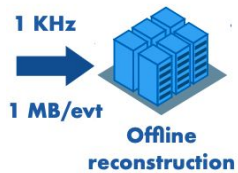- Includes an interface to the library and several utilities

More in:
S. Summers tutorial Thursday @ 9:00/13:00
V. Loncar talk Monday @ 15:00
H. Javed talk Monday @ 15:40
T. Aarrestad talk Tuesday @ 14:50
V. Razavimaleki talk Wednesday @ 15:30

# Ultra Fast ( 1ns)

**1 KHz**

**1 MB/evt**

Offline reconstruction

**100 KHz**
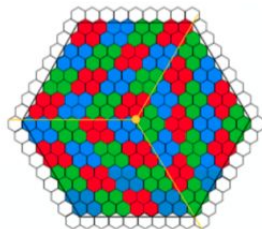
High-Level trigger

**40 MHz**

L1 trigger

## Example: CMS HG calorimeter

**Input**
48 "trigger cells"
7b floating point
(336b total)

ASIC

**Output:**
"Super trigger cell" algo
3[16 TC sum] x 13b = 39b

Can we do a better job of encoding the info in those bits w/o so much loss in granularity?

Encoder on ASIC ← → Decoder on L1 board

Encoder network

input shape    conv2D    Flatten  Dense    decoded shape

kernel    conv_y

n_filters  conv_x

(conv_y)x(conv_x)x(n_filters)  encoding dim

Really need quantized training here to optimize information encoding

**Use QKeras!**

Expansive part:
("volume" of conv. ouput)
x
(encoding dim)

More in G. Gluglielmo talk
Monday @ 14:50

# Challenges in adjacent science domains

# Gravitational Waves Challenges

## Inference-as-a-Service

LIGO deep clean algorithm, Alec Gunny

- Portable

Download

- Framework and architecture agnostic
- Critical for applications like DeepClean that need frequent retraining

- Co-locate downstream models for better resource allocation/autoscaling
- Manage and accelerate end-to-end latency of DeepClean + downstream algorithms to meet requirements

# Neutrino Astrophysics challenge

Kate Scholberg

**SNEWS Alert Latency**

From A. Habig, M. Strait



Now, the example where FastML may help!
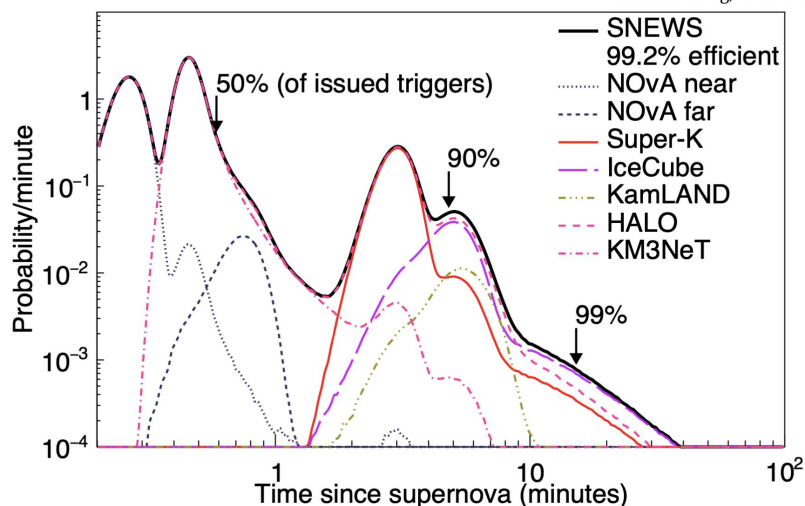
## Methods for Pointing Using Neutrinos
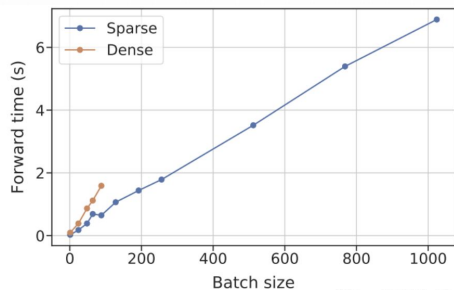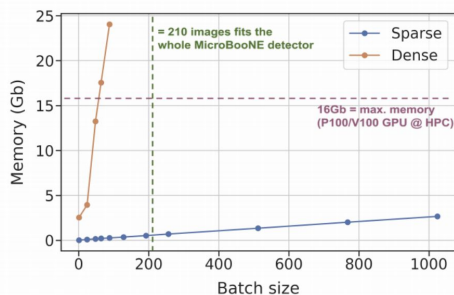
- ❑ **Anisotropic neutrino interactions**
  combined with detector technology that can exploit it,
  using the burst neutrino signal
- ❑ **Triangulation**
  using inter-detector timing
- ❑ **Oscillation pattern pointing**
  in high-energy resolution detectors
- ❑ **High-energy (~GeV) neutrino follow-on pointing**
  in directional detectors, using later neutrinos
- ❑ **All of the above!**

# Accelerator Neutrinos  reading out and processing

## Sparse Convolutions

CINCINNATI

Jeremy Hewes

- First paper investigating sparse CNNs within neutrino physics demonstrated significant improvements in inference time and memory usage for a MicroBooNE-equivalent detector.

- Sparse convolutions remove the need for ROI-finding in large detectors.

  - Scale of sparse pixel map set by **number of active pixels**.

  - Detector region can be arbitrarily large

- Sparse CNN approaches are being developed across the SBN and DUNE Near Detector by SLAC, and in NOvA and ProtoDUNE by University of Cincinnati.

Computational challenge  in Accelerator Neutrinos Are quickly becoming large



arXiv:1903.05663

# Challenges at Electron Ion Collider (EIC)

Markus Diefenthaler

**Streaming readout and its opportunities**

**Definition of streaming readout**

- data is read out in continuous parallel streams that are encoded with information about when and where the data was taken.
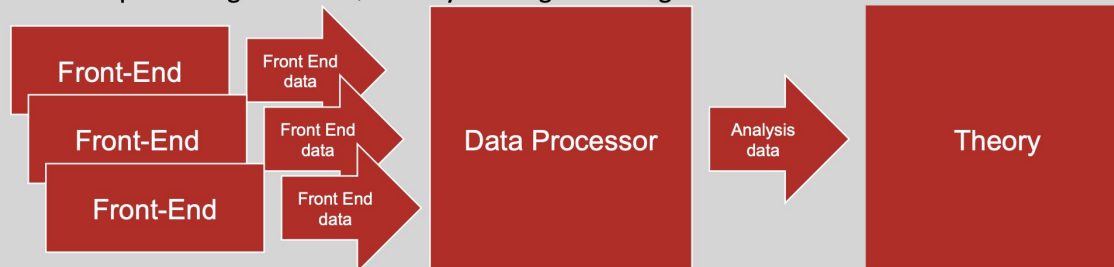
**Advantages of streaming readout**

- opportunity to streamline workflows
- take advantage of other emerging technologies, e.g. AI / ML

**Integration of DAQ, analysis and theory to optimize physics reach**

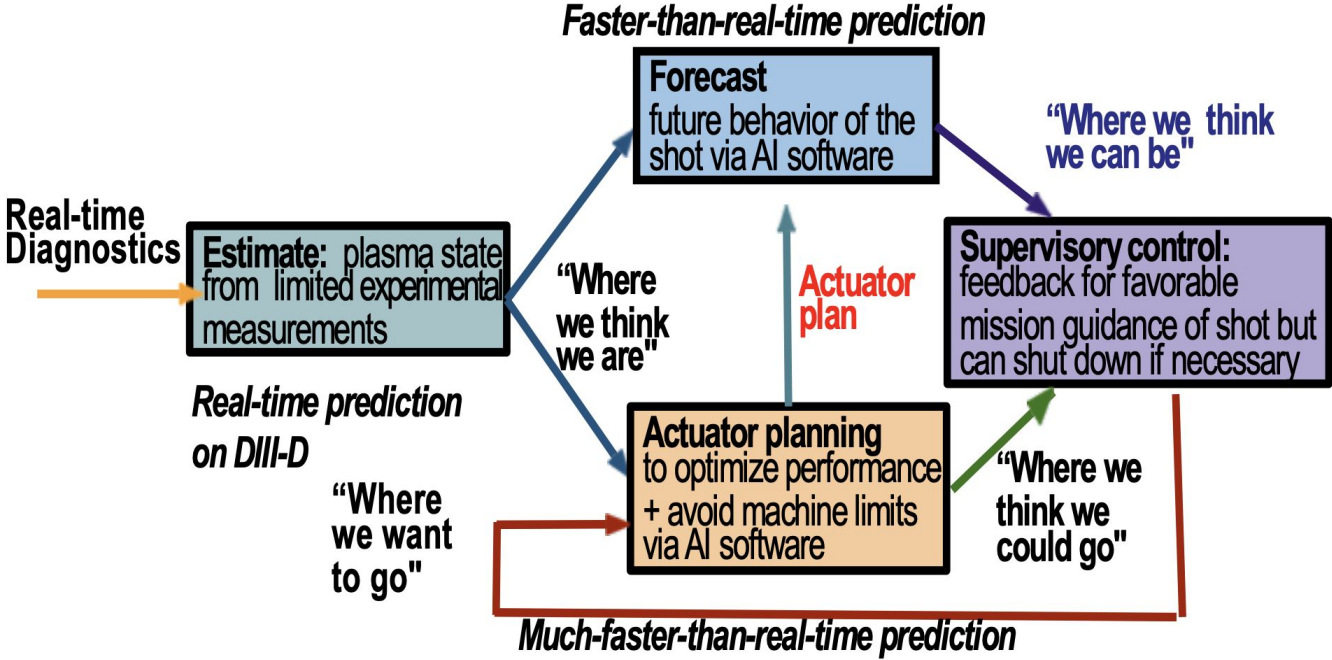- seamless data processing from DAQ to analysis using streaming readout



- opportunity for near real-time analysis using AI / ML
- opportunity to accelerate science (significantly faster access to physics results)

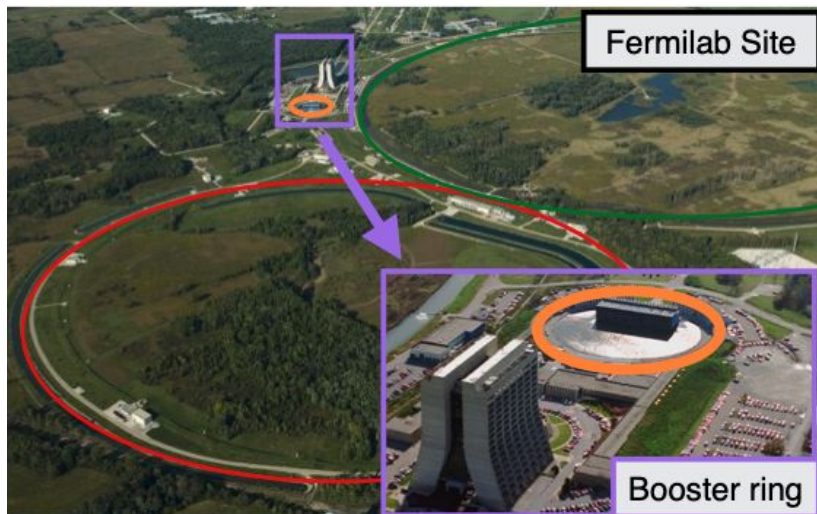Jefferson Lab

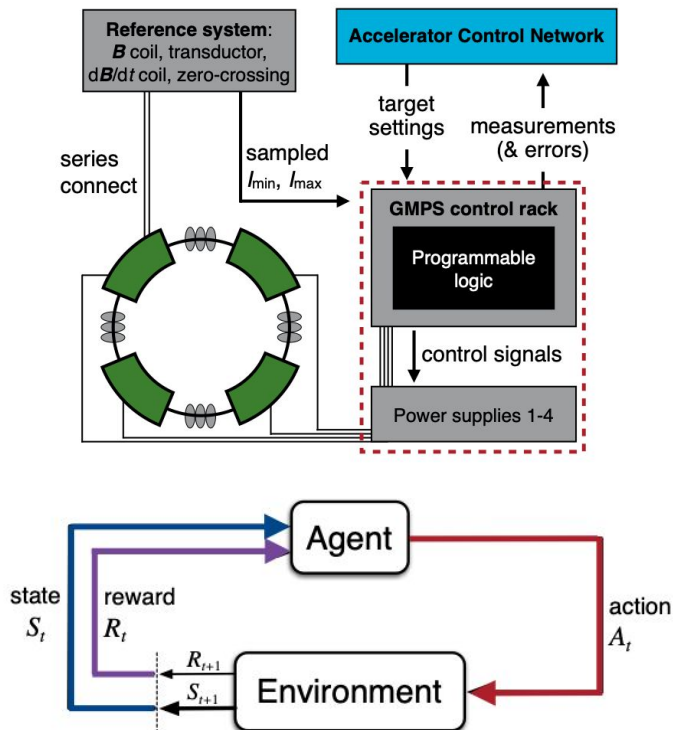# Plasma Physics: From Prediction to Control

Bill Tang

- How do we control a plasma?



*Faster-than-real-time prediction*

**Forecast**
future behavior of the shot via AI software

"Where we think we can be"

**Real-time Diagnostics**

**Estimate:** plasma state from limited experimental measurements

"Where we think we are"

**Actuator plan**

**Supervisory control:** feedback for favorable mission guidance of shot but can shut down if necessary

*Real-time prediction on DIII-D*

"Where we want to go"

**Actuator planning** to optimize performance + avoid machine limits via AI software

"Where we think we could go"

*Much-faster-than-real-time prediction*

23

# Accelerator Control

Christian Herwig

- How do we control proton accelerator?

Booster: 400 MeV—> 8 GeV

# Material Science
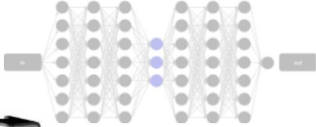
Processing Large amounts of data in near real-time

Materials Science, Josh Agar

Camera

Laser Energy

2 Gbps

<2 Gb/deposition

Database

FPGA Accelerator
<50 ms inference

Mass Flow Controller

Laser Optics

# Conclusions

New developments in computing allow for larger use of computing
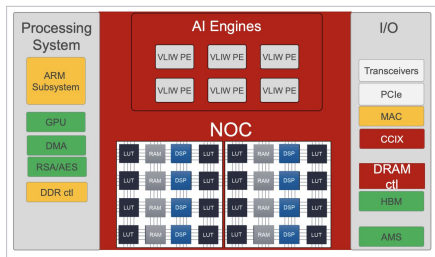
- New tools are emerging to allow for more advanced applications
  - Tools for large scale use of GPUs, FPGAs
  - Tools to improve algorithm design to reduce latency and resources
- New opportunities for use of fast ML in design
  - Real-time readout and processing in LHC collider physics
  - Real-time readout in neutrino physics, astrophysics and gravitational wave
  - Real-time control and operations in accelerator, plasma and materials science
- Supplement workshop with <span style="color:red">forward-looking white paper</span> to
  - Establish key applications and identify overlaps between domains scientific applications for resource-constrained ML (tentative timeline Feb 15 2021)

# Backup

# Trends in Computer Architecture

**Specialization = > Increasingly Heterogeneous**
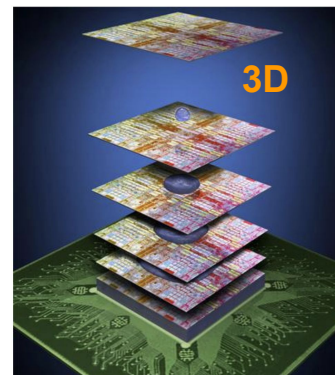
5 Series
7 Series
U+ Series
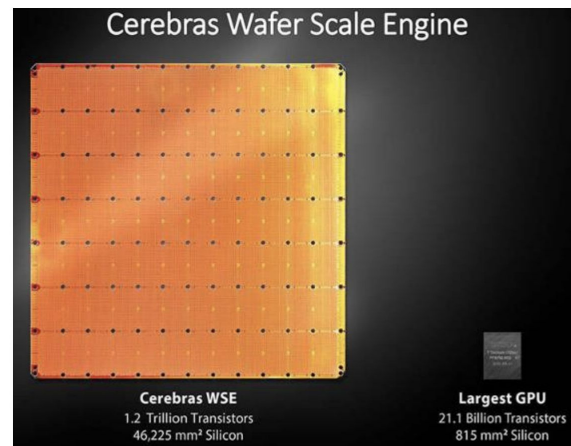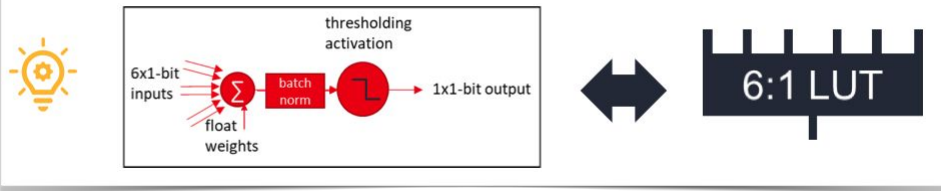Versal



Xilinx Example:
FPGA -> ACAP

More hardened functionality (=> heterogeneous)
to improve compute density and save power
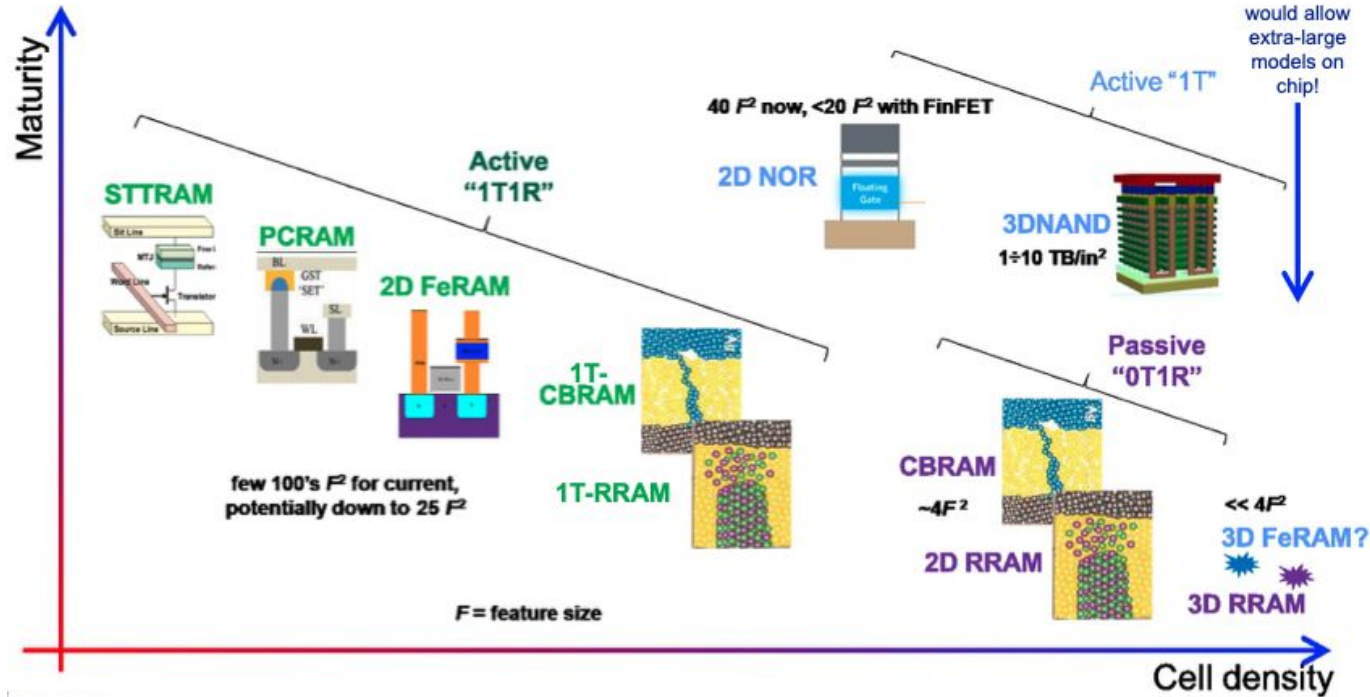
© Copyright 2020 Xilinx

**ΣXILINX**



3D



Cerebras Wafer Scale Engine

**Cerebras WSE**
1.2 Trillion Transistors
46,225 mm² Silicon

**Largest GPU**
21.1 Billion Transistors
815 mm² Silicon

28

*Source: HotChips2019*

# Beyond CMOS: Neuromophic Computer



Dmitri Strukov (UCSB)

Analog Circuit: natural fit for vector*matrix multiplications