# ML in FPGA
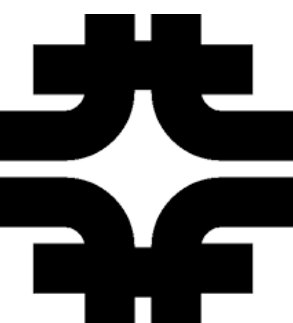## (& ASICs, etc)
## (for embedded systems)
## (for science, particularly the EIC)

Nhan Tran, Fermilab
October 12, 2022

# Outline

- Motivation

- State-of-the-art workflow for FPGA/ASIC

  - Towards a sustainable and robust ecosystem
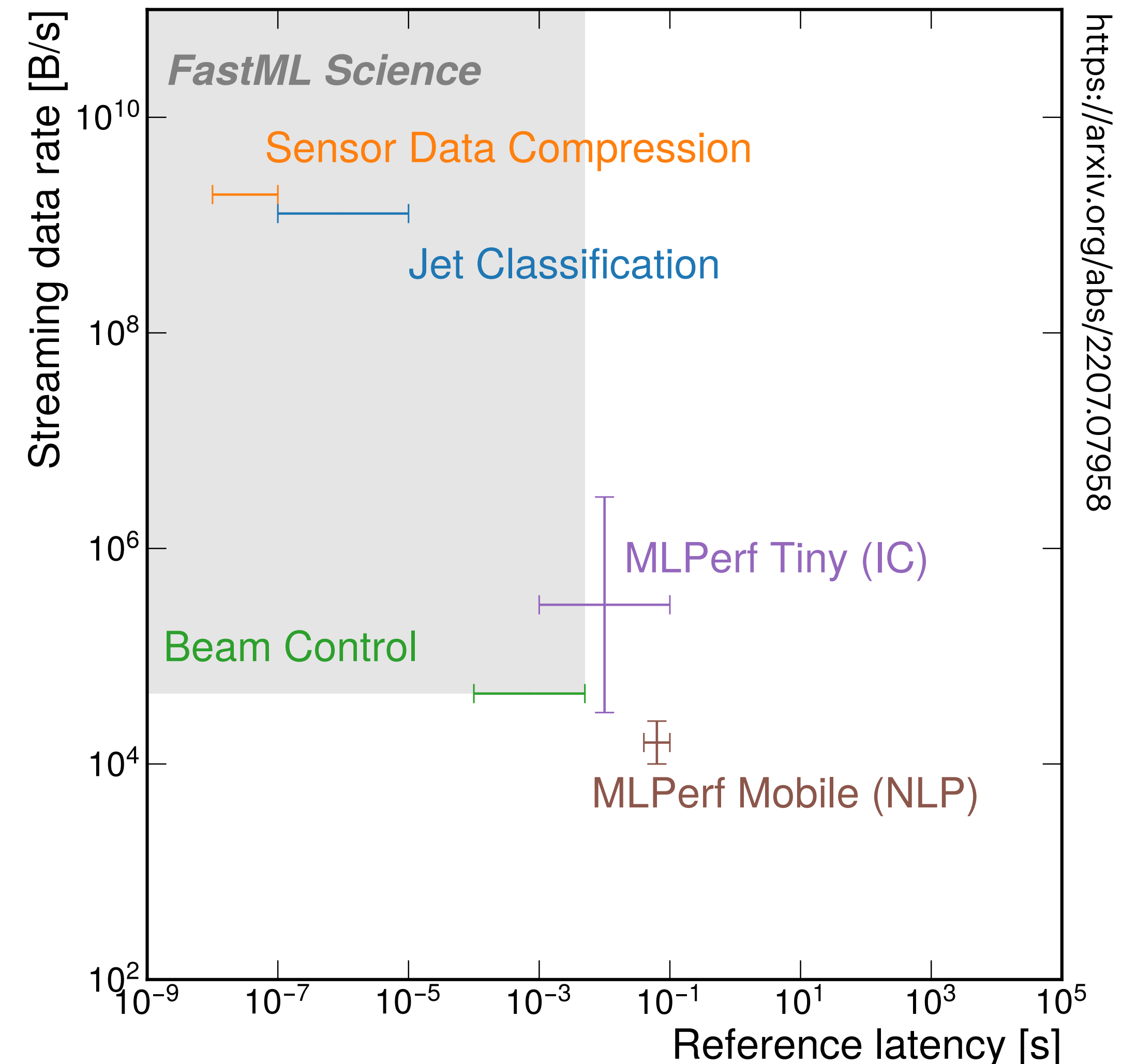
  - Emerging technologies

*This is a big area!*
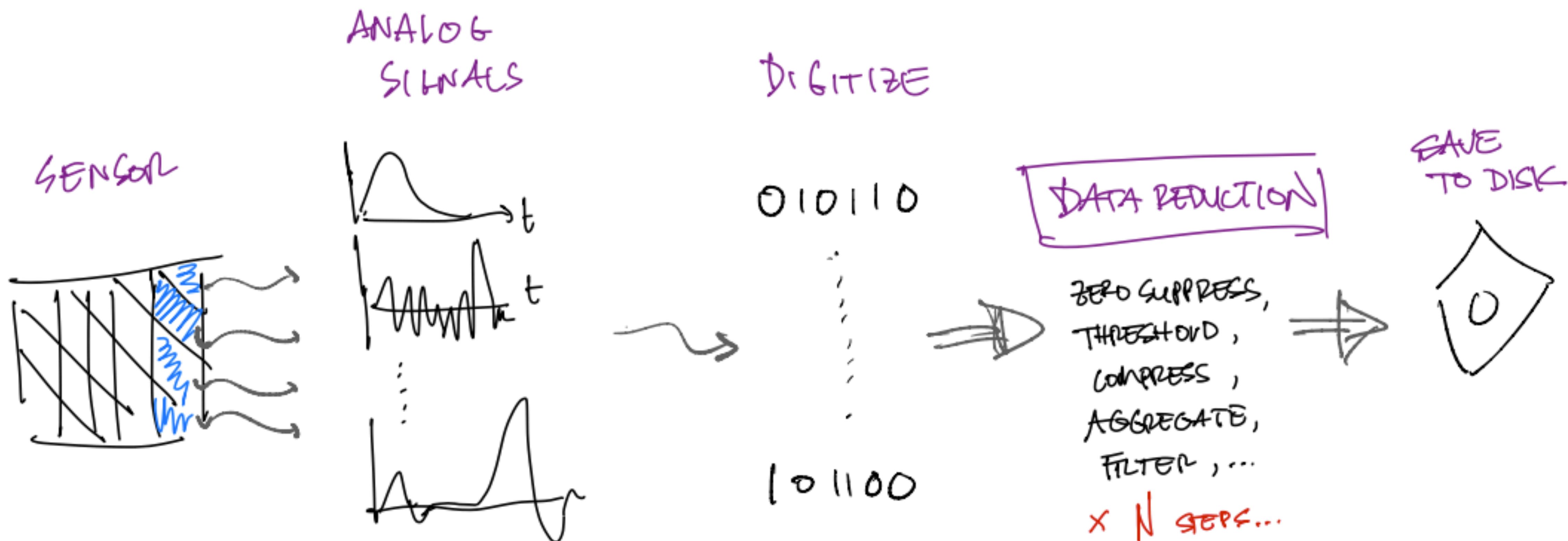My approach — present key important topics and provide a lot of references.
Follow references if you are interested in learning more;
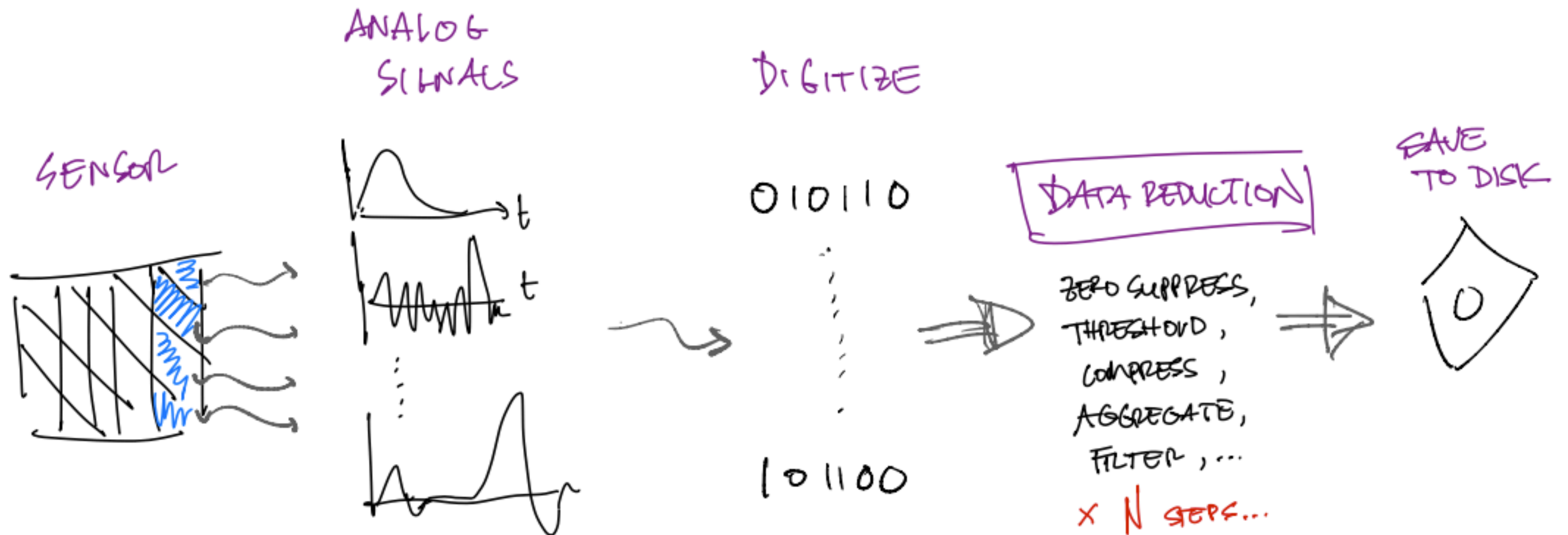reach out if you are even more interested after following the references.

# Motivation

High intensity collider experiments explore nature at the **finest temporal and spatial scales** Leads to data rates far surpassing industry — requires developing **innovative techniques**
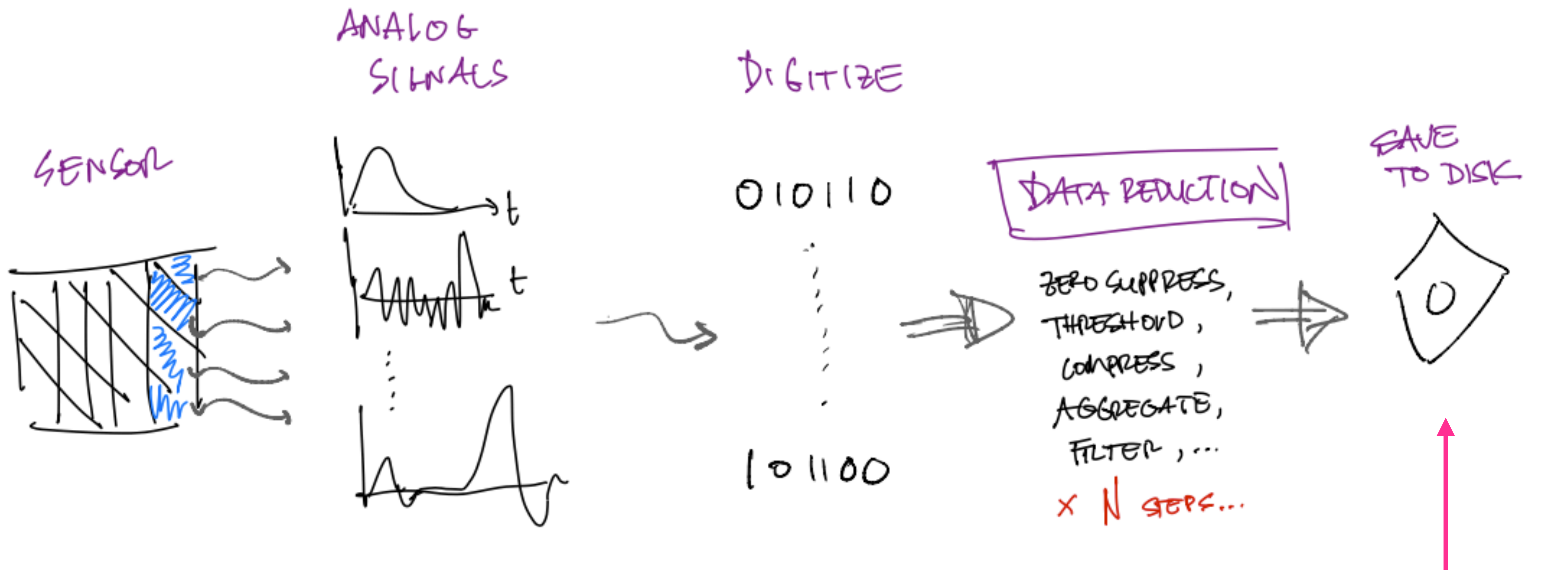
- ML in specialized embedded architectures require in ***real-time*** to reduce and filter data

- Optimal data selection enables **more efficient operation, saves lost data, and accelerates time-to-discovery**

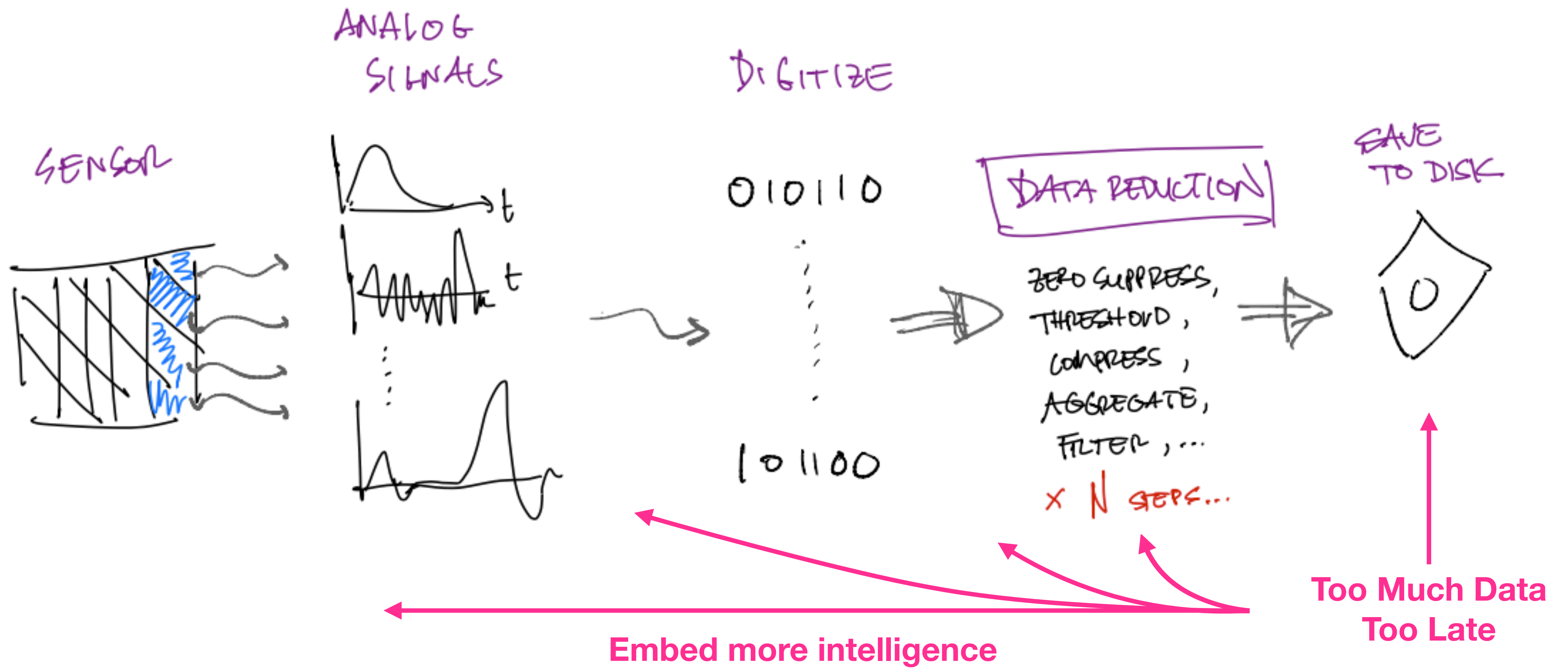SENSOR

ANALOG SIGNALS

DIGITIZE

010110

101100

DATA REDUCTION

ZERO SUPPRESS,
THRESHOLD,
COMPRESS,
AGGREGATE,
FILTER, ...
× N STEPS...

SAVE TO DISK

**1 channel ~ 10b**
**1 channel, 1 MHz rate ~ 10 Mb/s**
**100k channels, 1 MHz rate ~ 1 Tb/s**

1 channel ~ 10b
1 channel, 1 MHz rate ~ 10 Mb/s
100k channels, 1 MHz rate ~ 1 Tb/s

# Applications in nuclear physics and beyond

https://indico.cern.ch/e/fml2022



A workshop dedicated to real-time applications of ML across the sciences

See also:
Applications and Technique in
Fast Machine Learning for Science
https://www.frontiersin.org/articles/10.3389/fdata.2022.787421/full

5

# Applications in nuclear physics and beyond

https://indico.cern.ch/e/fml2022



A workshop dedicated to real-time applications of ML across the sciences

See also:
Applications and Technique in
Fast Machine Learning for Science
https://www.frontiersin.org/articles/10.3389/fdata.2022.787421/full

See Jin Huang's talk for a great overview of exciting real-time applications



I will be referencing other talks from the Fast ML workshop to point to other examples of state-of-the-art studies

# State-of-the-art



https://pypi.org/project/hls4ml/
682 Github stars,
580 downloads last month

QKeras (Google)
Brevitas (AMD)
HAWQ (UC Berkeley)
QONNX (Microsoft/AMD)

# State-of-the-art



Keras
PyTorch
ONNX

hls4ml

FPGAs
ASICs

https://pypi.org/project/hls4ml/
682 Github stars,
580 downloads last month

Model → Quantized model → **hls4ml** → HLS project → Hardware

**QKeras** (Google)
**Brevitas** (AMD)
**HAWQ** (UC Berkeley)
**QONNX** (Microsoft/AMD)

VIVADO Mentor
A Siemens Business

Quartus
Prime

# State-of-the-art

**Physics requirements**

**Data representation
→ ML architecture**

**Neural architecture search/
Hyperparameter optimization**

hls4ml

https://pypi.org/project/hls4ml/
682 Github stars,
580 downloads last month

Model → Quantized model → **hls4ml** → HLS project → Hardware

**QKeras** (Google)
**Brevitas** (AMD)
**HAWQ** (UC Berkeley)
**QONNX** (Microsoft/AMD)

FPGAs    ASICs

VIVADO    Mentor

# State-of-the-art

**Physics requirements**

**Data representation**
**→ ML architecture**

**Neural architecture search/**
**Hyperparameter optimization**

**What kind of platform?**

**Latency?**
**Pipeline Interval?**

**How many**
**resources?**

**Area/power?**
**Radiation?**
**Cryo?**

https://pypi.org/project/hls4ml/
682 Github stars,
580 downloads last month

Model → Quantized model → hls4ml → HLS project → Hardware

QKeras (Google)
Brevitas (AMD)
HAWQ (UC Berkeley)
QONNX (Microsoft/AMD)

VIVADO  Mentor

# State-of-the-art

**Physics requirements**

**Data representation
→ ML architecture**

**Neural architecture search/
Hyperparameter optimization**

**What kind of platform?**

**Latency?
Pipeline Interval?**

**How many
resources?**

**Area/power?
Radiation?
Cryo?**



https://pypi.org/project/hls4ml/
682 Github stars,
580 downloads last month

Model   Quantized model   hls4ml   HLS project   ware

QKeras (Google)
Brevitas (AMD)
HAWQ (UC Berkeley)
QONNX (Microsoft/AMD)

VIVADO

AI circuit for ultrafast inference on FPGA

Inference time: 280 ns
Throughput: 104 Gb/s

Dense Network
23 → 30 → 25 → 20
→ momentum & classifier

# State-of-the-art

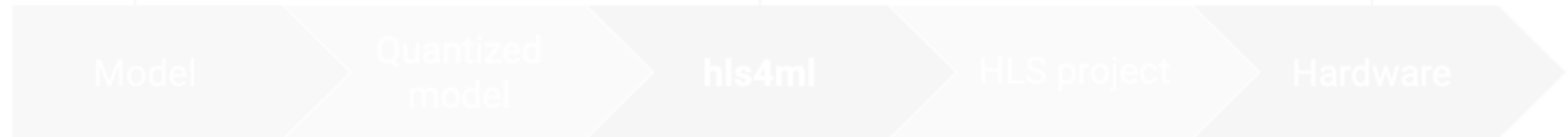**Physics requirements**

**Data representation → ML architecture**

**Neural architecture search/ Hyperparameter optimization**

**Quantize network**

**What kind of platform?**

**Latency? Pipeline Interval?**

**How many resources?**

**Area/power? Radiation? Cryo?**

**See tools like: QKeras HAWQ Brevitas**

| Relative Energy Cost | | |
|---|---|---|
| Operation: | Energy (pJ) | |
| 8b Add | 0.03 | |
| 16b Add | 0.05 | |
| 32b Add | 0.1 | |
| 16b FP Add | 0.4 | |
| 32b FP Add | 0.9 | |
| 8b Mult | 0.2 | |
| 32b Mult | 3.1 | |
| 16b FP Mult | 1.1 | |
| 32b FP Mult | 3.7 | |
| 32b SRAM Read (8KB) | 5 | |
| 32b DRAM Read | 640 | |

*Adapted from Horowitz*

1   10   100   1000   10000

**Roughly quadratic**

9

# State-of-the-art

**Physics requirements**

**Data representation → ML architecture**

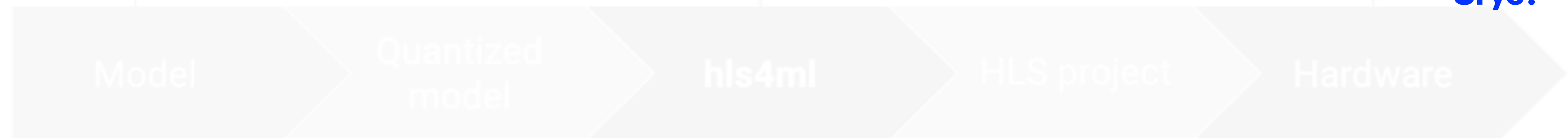**Neural architecture search/ Hyperparameter optimization**

**Quantize network**

**Intermediate (quantized) representations**

**What kind of platform?**

**Latency? Pipeline Interval?**

**How many resources?**

**Area/power? Radiation? Cryo?**

hls4ml

HLS project

Hardware



QKeras (Google)
Brevitas (AMD)
HAWQ (UC Berkeley)
QONNX (Microsoft/AMD)

**See proposal for QONNX**

10

# State-of-the-art

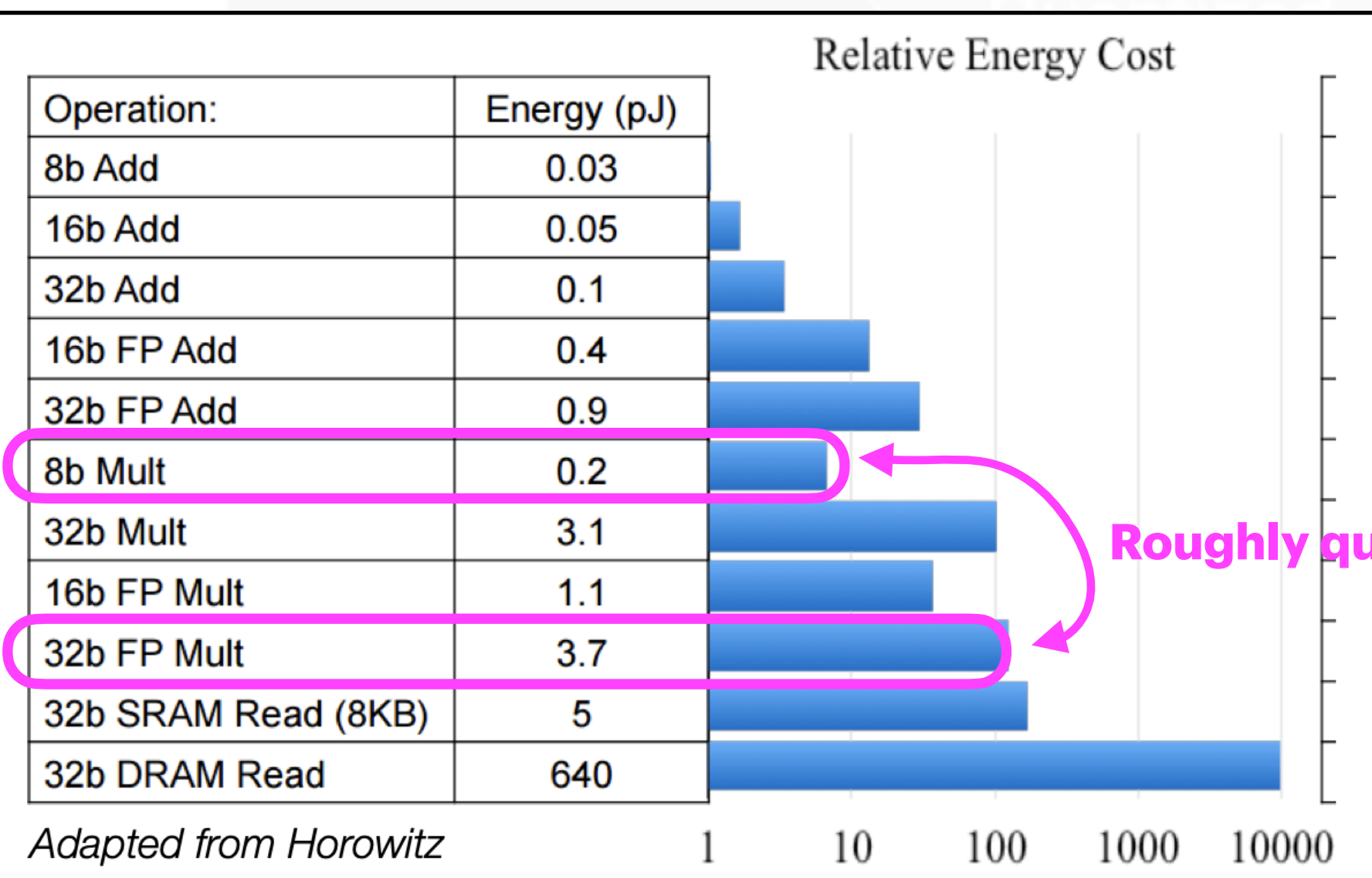**Physics requirements**

**Data representation → ML architecture**

**Neural architecture search/ Hyperparameter optimization**

**Quantize network**

**Intermediate (quantized) representations**

**Pruning/sparsity?**

https://pypi.org/project/hls4ml/
682 Github stars,
580 downloads last month

Model → Quantized Model → hls4ml → HLS project → Hardware

**What kind of platform?**

**Latency? Pipeline Interval?**

**How many resources?**

**Area/power? Radiation? Cryo?**



before pruning — pruning synapses — after pruning

pruning neurons

# State-of-the-art

**Physics requirements**

**Data representation → ML architecture**

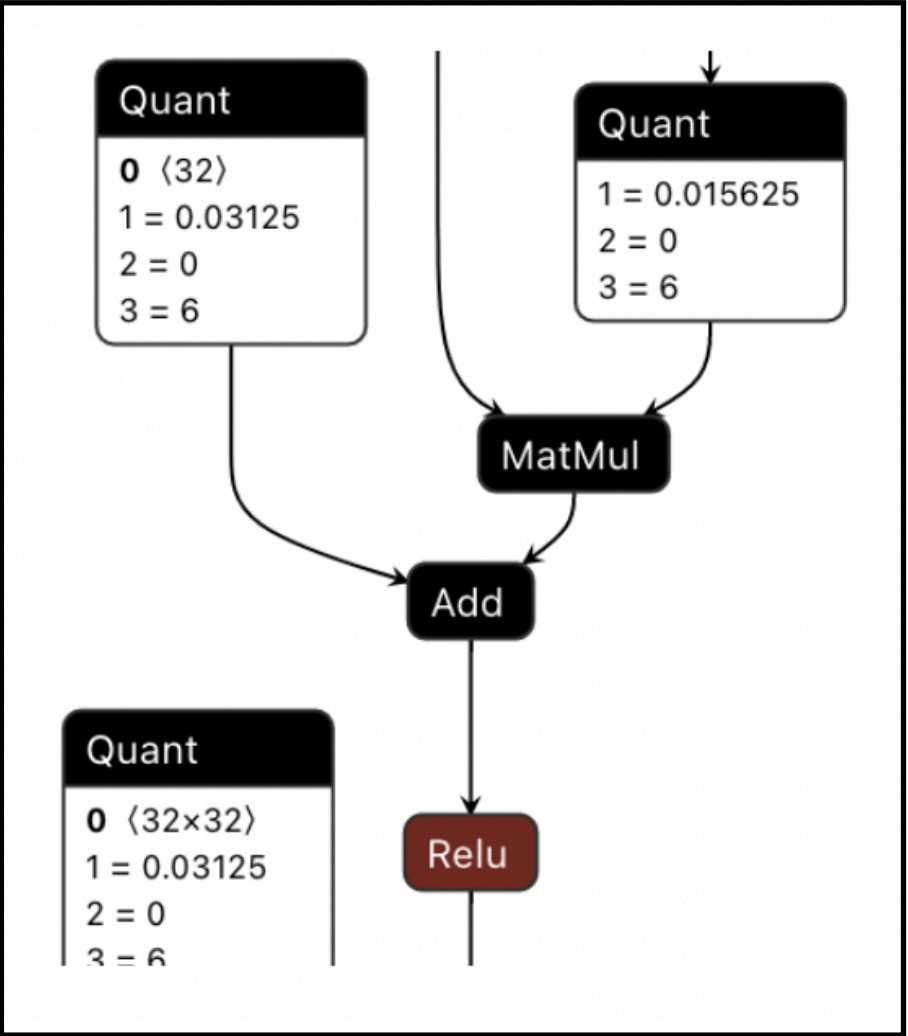**Neural architecture search/ Hyperparameter optimization**

**What kind of platform?**

**Latency? Pipeline Interval?**

**Microarchitecture**

**Quantize network**

**Intermediate (qu... representat...**

**Pru...**

QKeras (Google)
Brevitas (AMD)
HAWQ (UC Berkeley...
QONNX (Microsoft/AI...



DRAM

DPU (MPE)

DMA

Weight Buffer

Input & Activation Buffering

Compute Array
Matrix or Tensor Processing

Activation Functions/Pooling...

DPU (Spatial)

Dedicated Weight Buffers

Dedicated Compute Architecture

Activation Buffering

**Matrix of processing elements (Systolic Array)**

**Spatial Dataflow**

# State-of-the-art

**Physics requirements**

**Data representation → ML architecture**

**Neural architecture search/ Hyperparameter optimization**

**Quantize network**

**Intermediate (quantized) representations**

**Pruning/sparsity?**

QKeras (Google)
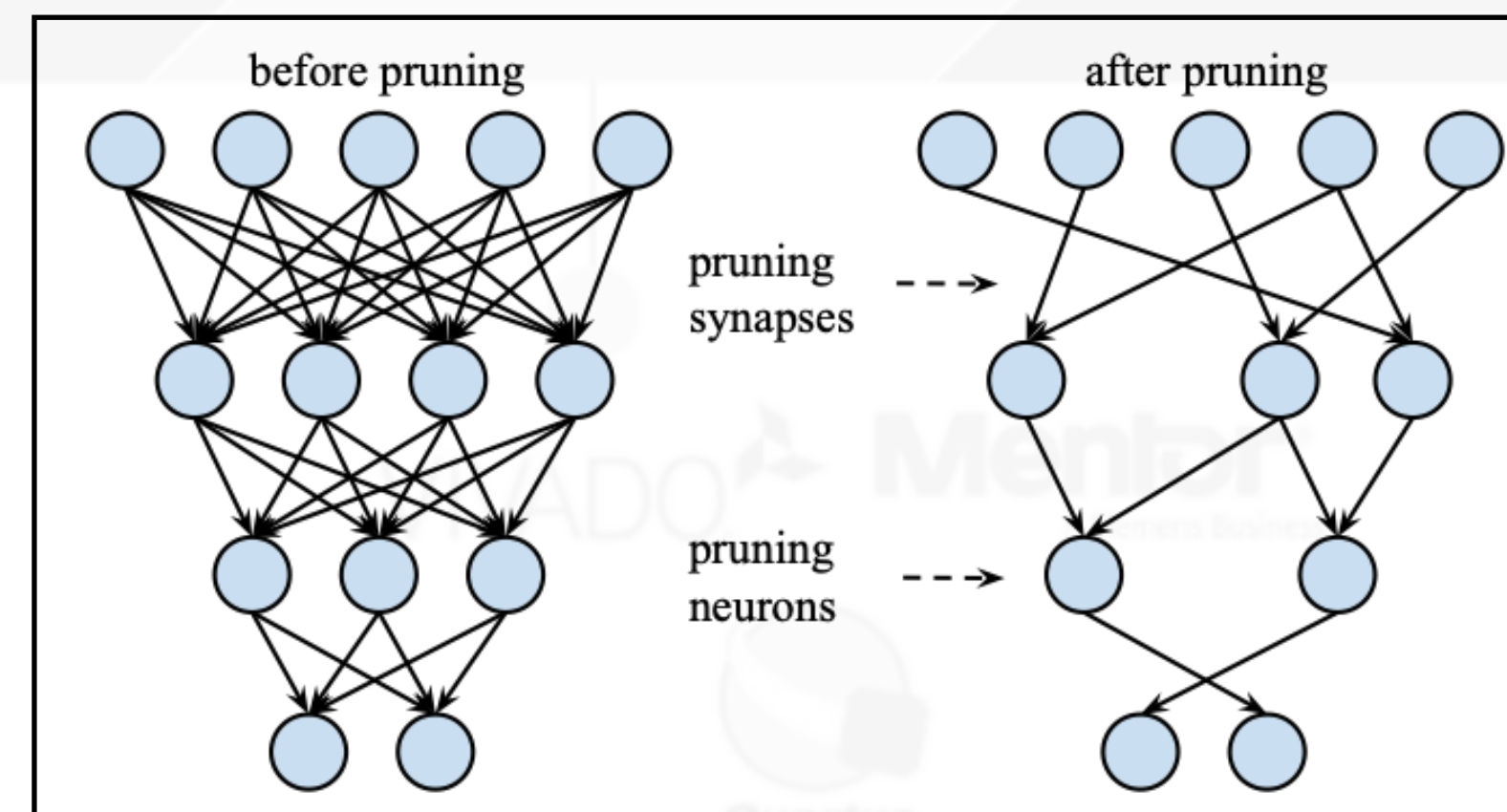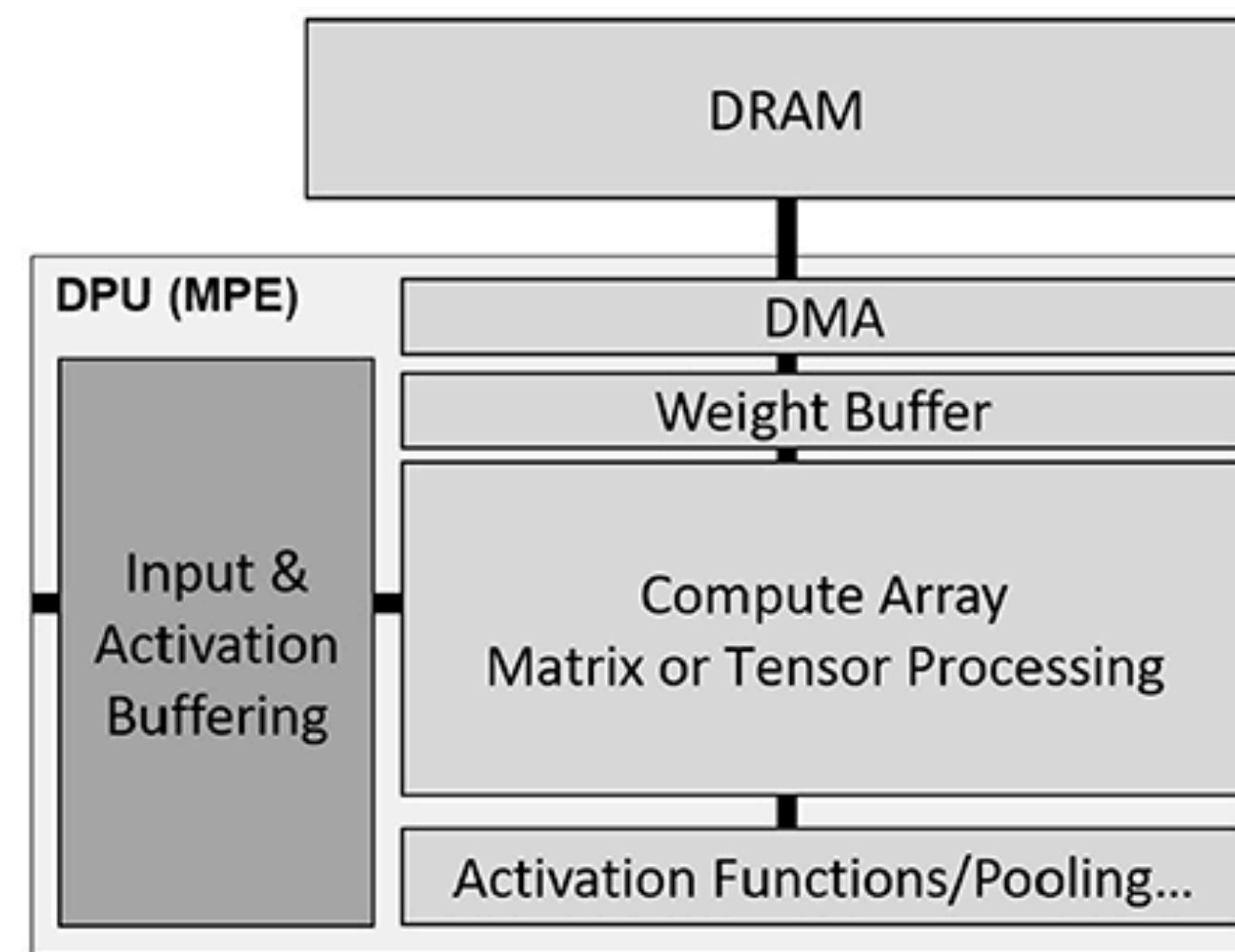Brevitas (AMD)
HAWQ (UC Berkeley)
QONNX (Microsoft/AMD)

**Microarchitecture**

**Parallelization**
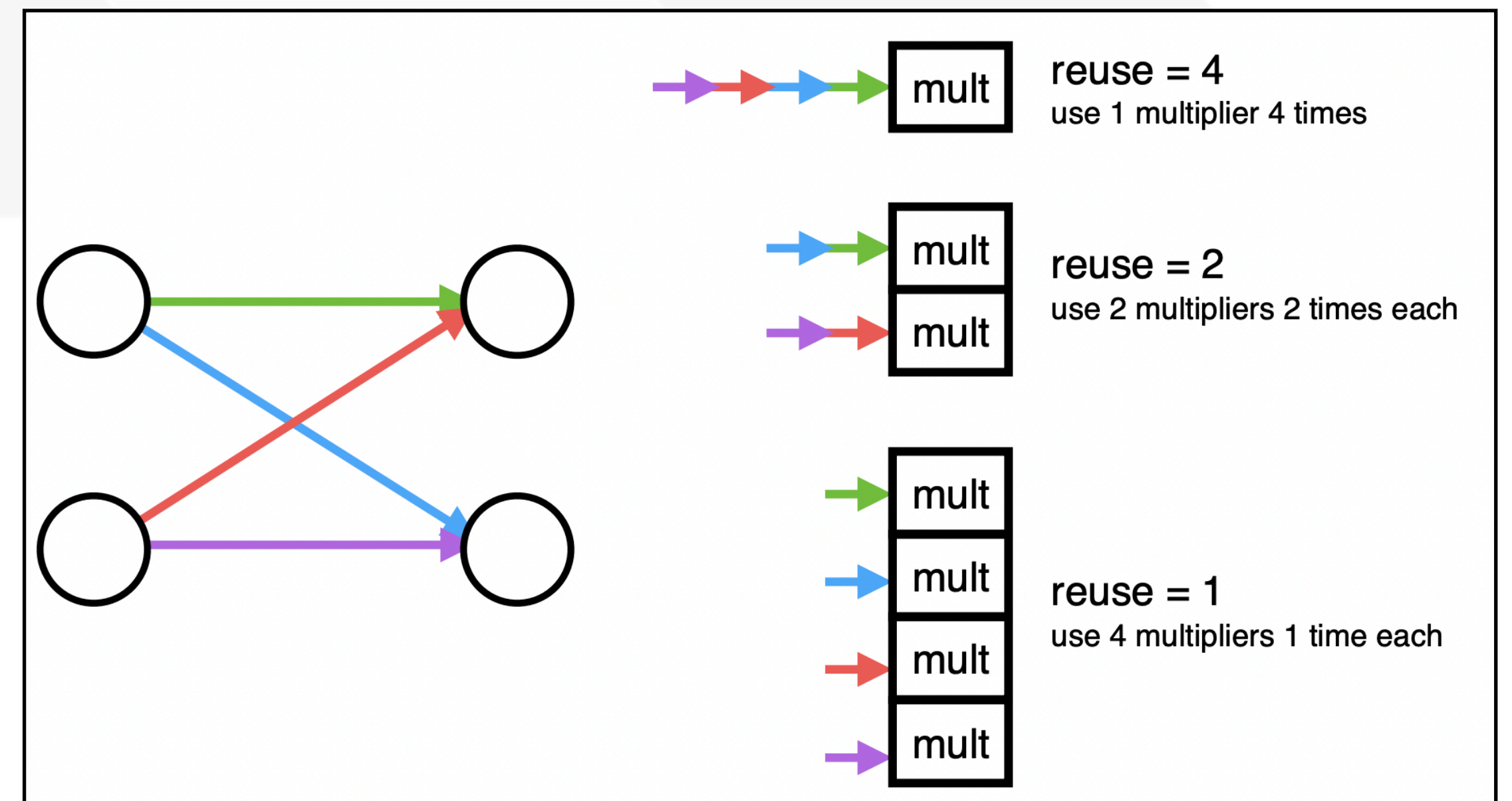
**What kind of platform?**

**Latency? Pipeline Interval?**

**How many resources?**

**Area/power? Radiation? Cryo?**



reuse = 4
use 1 multiplier 4 times

reuse = 2
use 2 multipliers 2 times each

reuse = 1
use 4 multipliers 1 time each

# State-of-the-art

Physics requirements

Data representation
→ ML architecture

Neural architecture search/
Hyperparameter optimization

Microarchitecture

What kind of platform?

Latency?
Pipeline Interval?

https://pypi.org/project/hls4ml/
682 Github stars,

Parallelization

How many
resources?

Area/power?
Radiation?
Cryo?

Quantize network

Model    Quantized    hls4ml    HLS project    Hardware
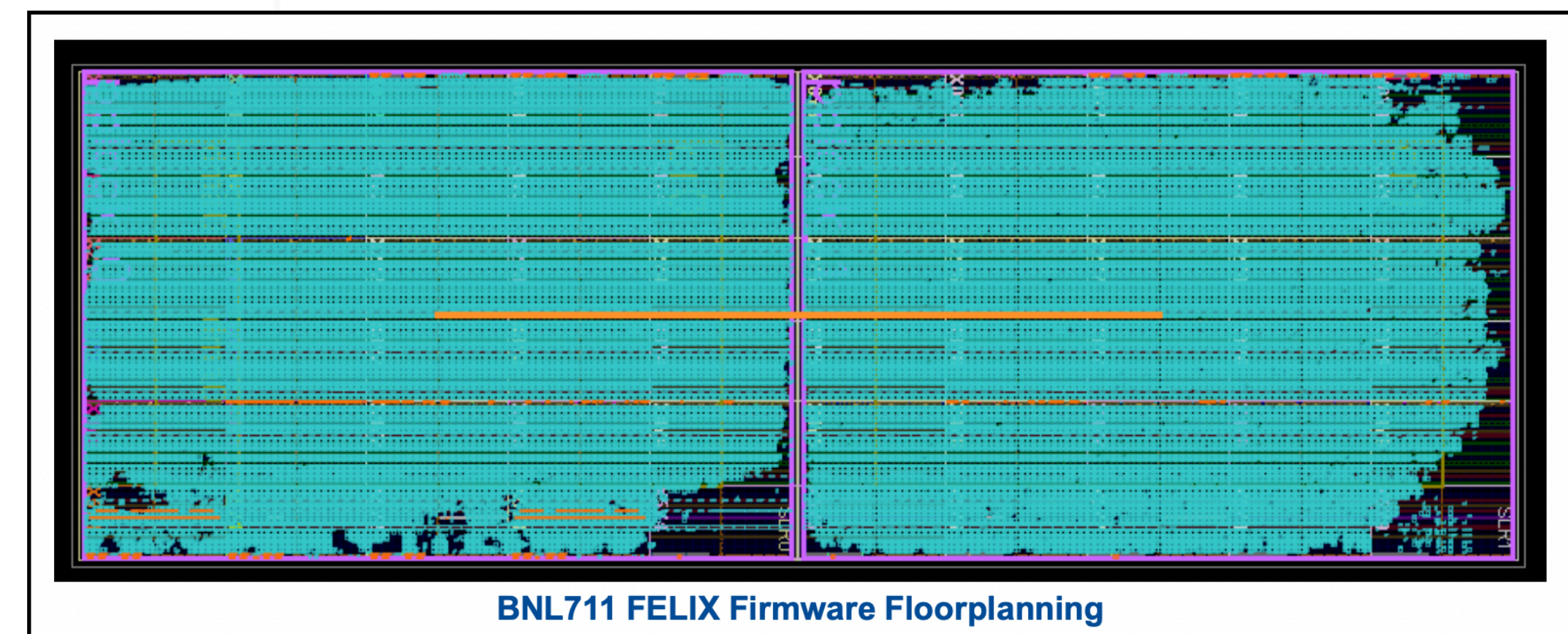
Intermediate (quantized)
representations

Synthesize, validate design,
satisfy design rules/timing

Pruning/sparsity?

QKeras (Google)
Brevitas (AMD)
HAWQ (UC Berkeley)
QONNX (Microsoft/AMD)



BNL711 FELIX Firmware Floorplanning

# State-of-the-art

**Physics requirements**

**Data representation → ML architecture**

**Neural architecture search/ Hyperparameter optimization**

**Quantize network**

**Intermediate (quantized) representations**

**Pruning/sparsity?**

**Microarchitecture**

**Parallelization**

**Synthesize, validate design, satisfy design rules/timing**

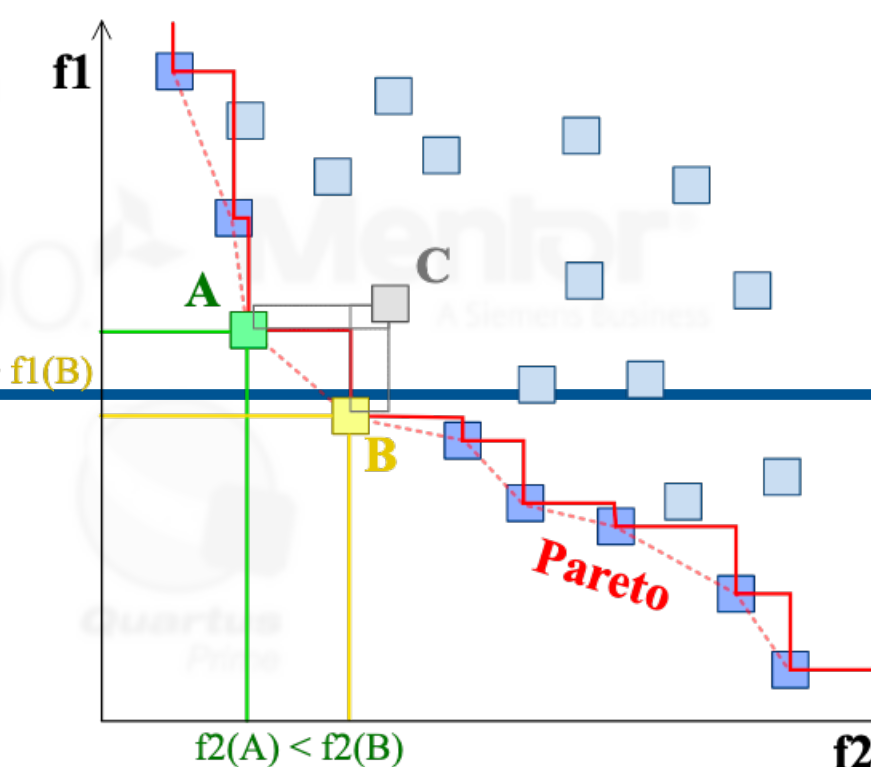**What kind of platform?**

**Latency? Pipeline Interval?**

**How many resources?**

**Area/power? Radiation? Cryo?**

**Multi-objective design space optimization**

# Towards sustainability and robustness

- Robustness and fault tolerance

  - DataPerf (V. Reddi)

  - FKeras (O. Weng)

  - Continual learning? (B. Radburn-Smith)

- Implementation within FW infrastructure, synthesize effectively

  - Issues, tricks, and tips
    (M. Rigatti, D. Hoang)

  - Emulating NNs in experimental SW?



Model-centric paradigm

Data-centric paradigm

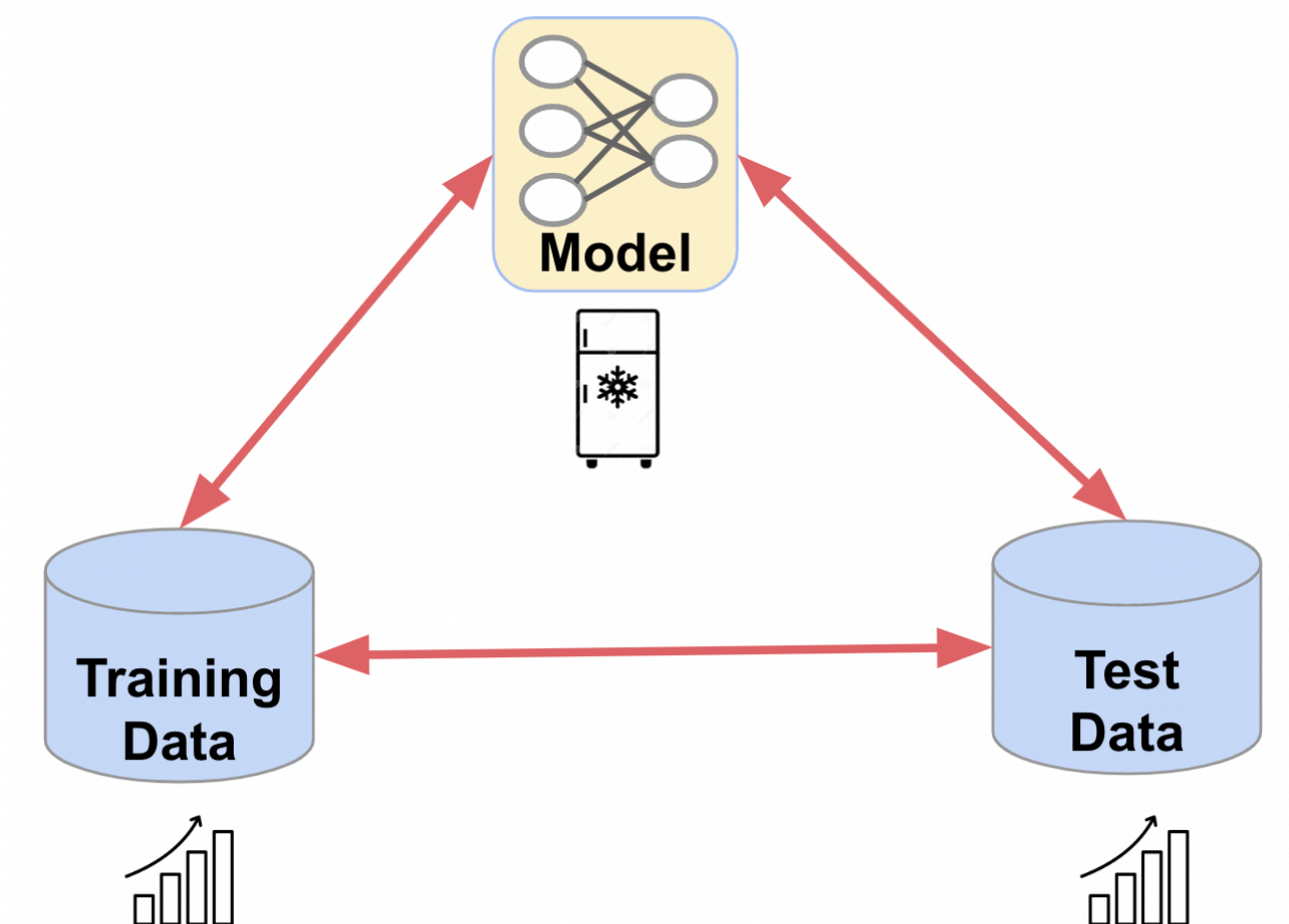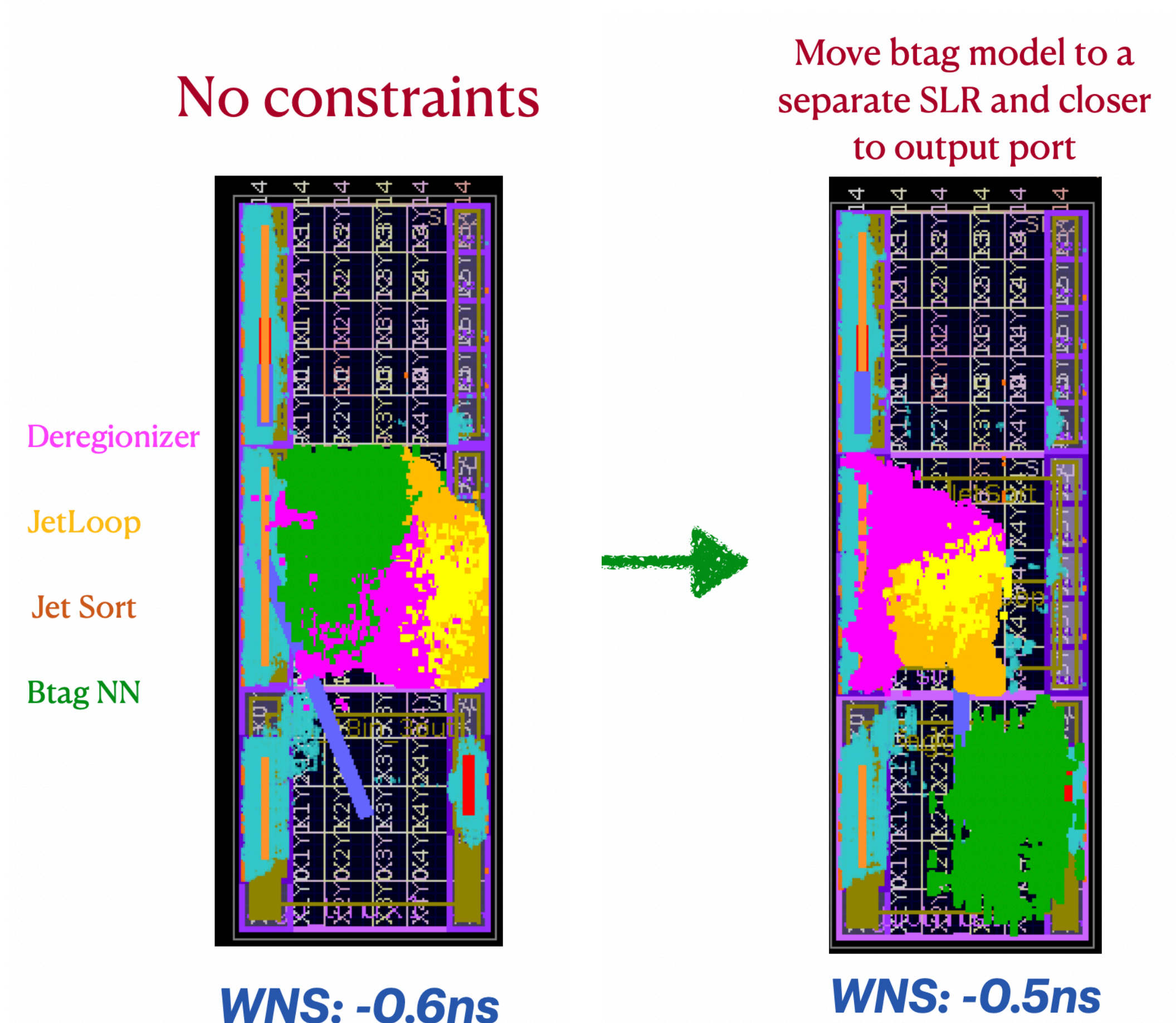# Towards sustainability and robustness

- Robustness and fault tolerance

  - DataPerf (V. Reddi)

  - FKeras (O. Weng)

  - Continual learning? (B. Radburn-Smith)

- Implementation within FW infrastructure, synthesize effectively

  - Issues, tricks, and tips
    (M. Rigatti, D. Hoang)

- Emulating NNs in experimental SW?

No constraints

Move btag model to a separate SLR and closer to output port

Deregionizer

JetLoop

Jet Sort

Btag NN

**WNS: -0.6ns**

**WNS: -0.5ns**

16

# Emerging technologies

- Emerging computing architectures

- Emerging neural architectures

  - Spiking, inductive bias, physics-inspired,...

- Emerging microelectronics technologies
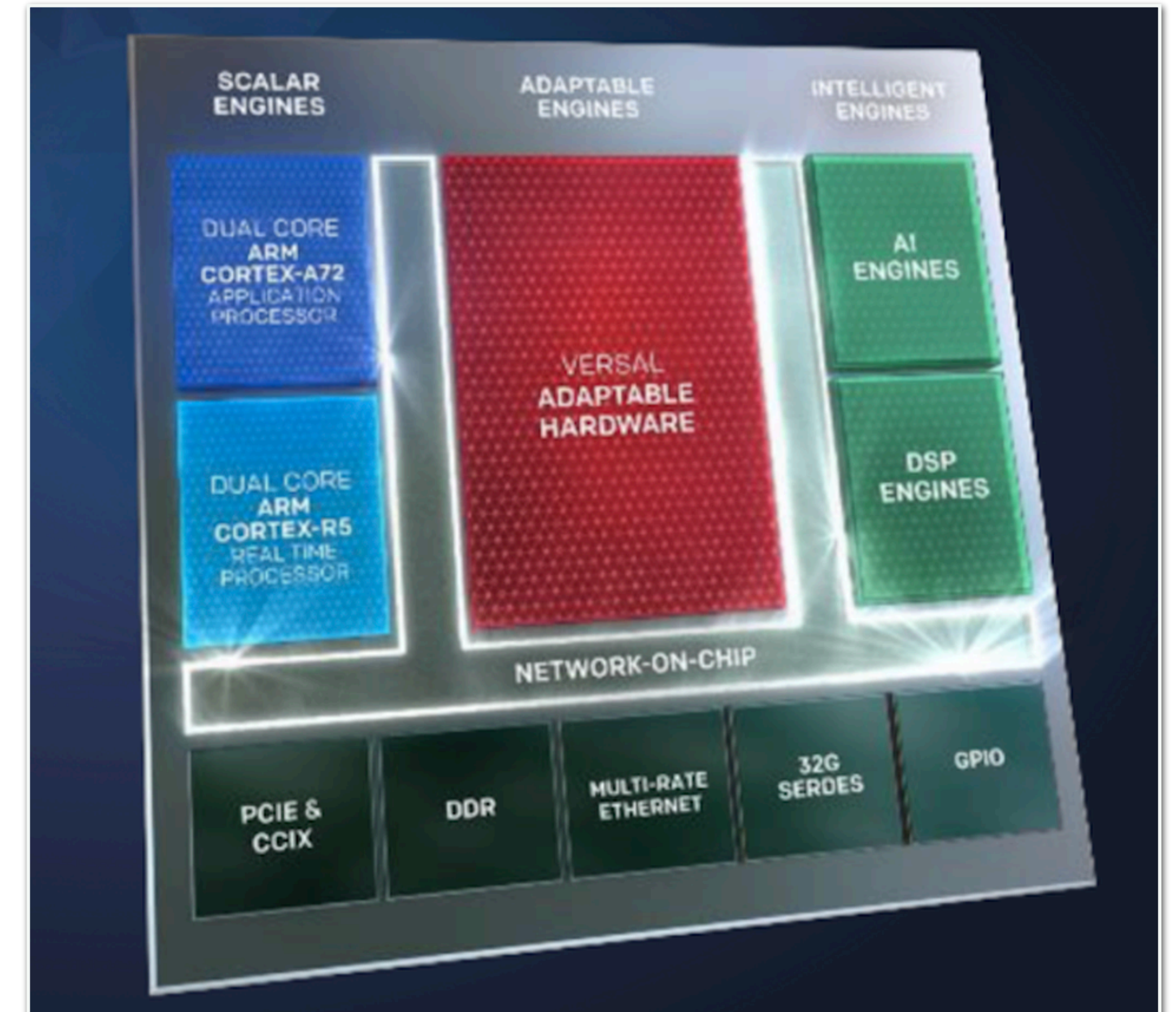
# Emerging technologies

- Emerging computing architectures

- Emerging neural architectures

  - Spiking, inductive bias, physics-inspired,...

- Emerging microelectronics technologies

# Emerging technologies

- Emerging computing architectures

- Emerging neural architectures

  - Spiking, inductive bias, physics-inspired,…
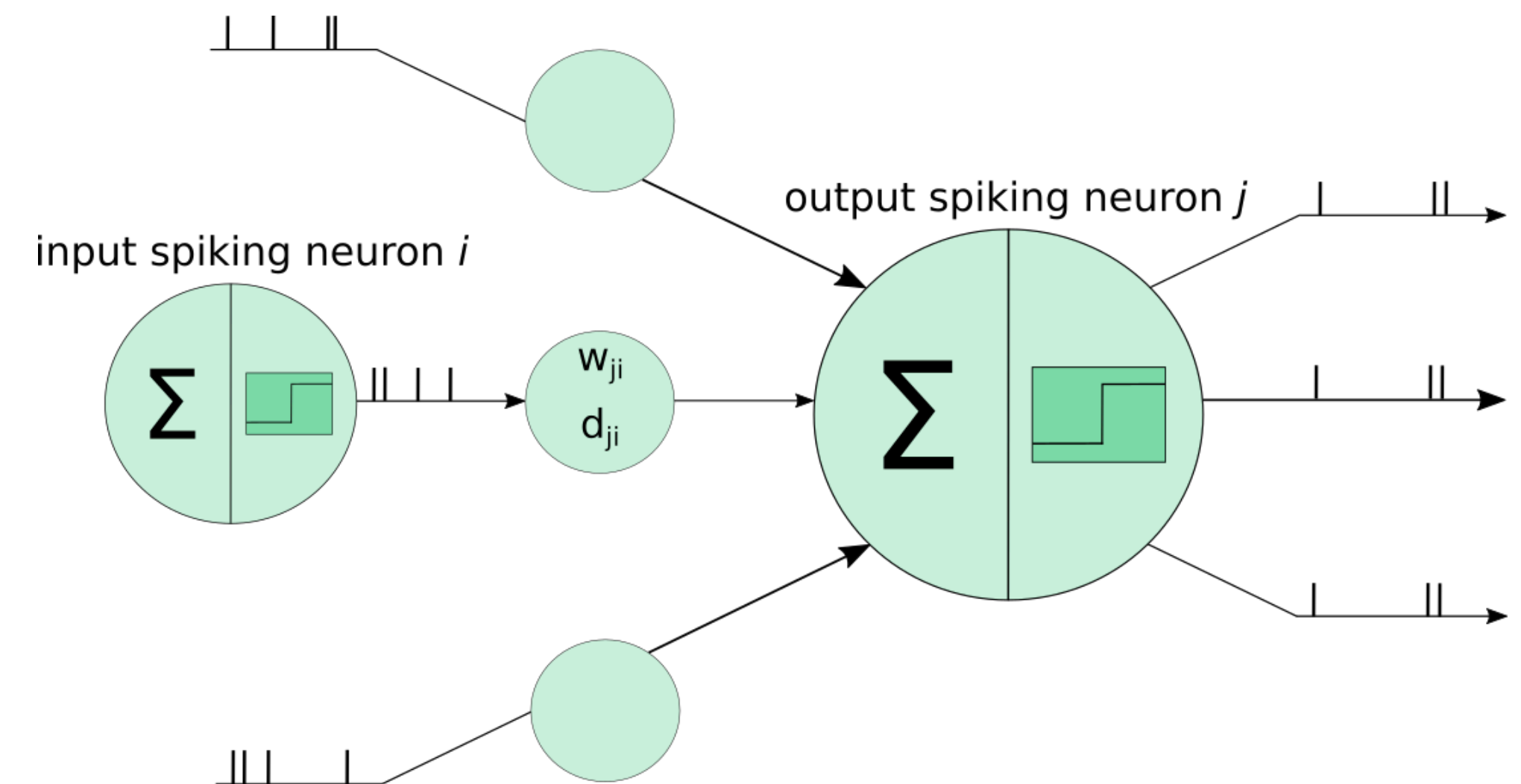
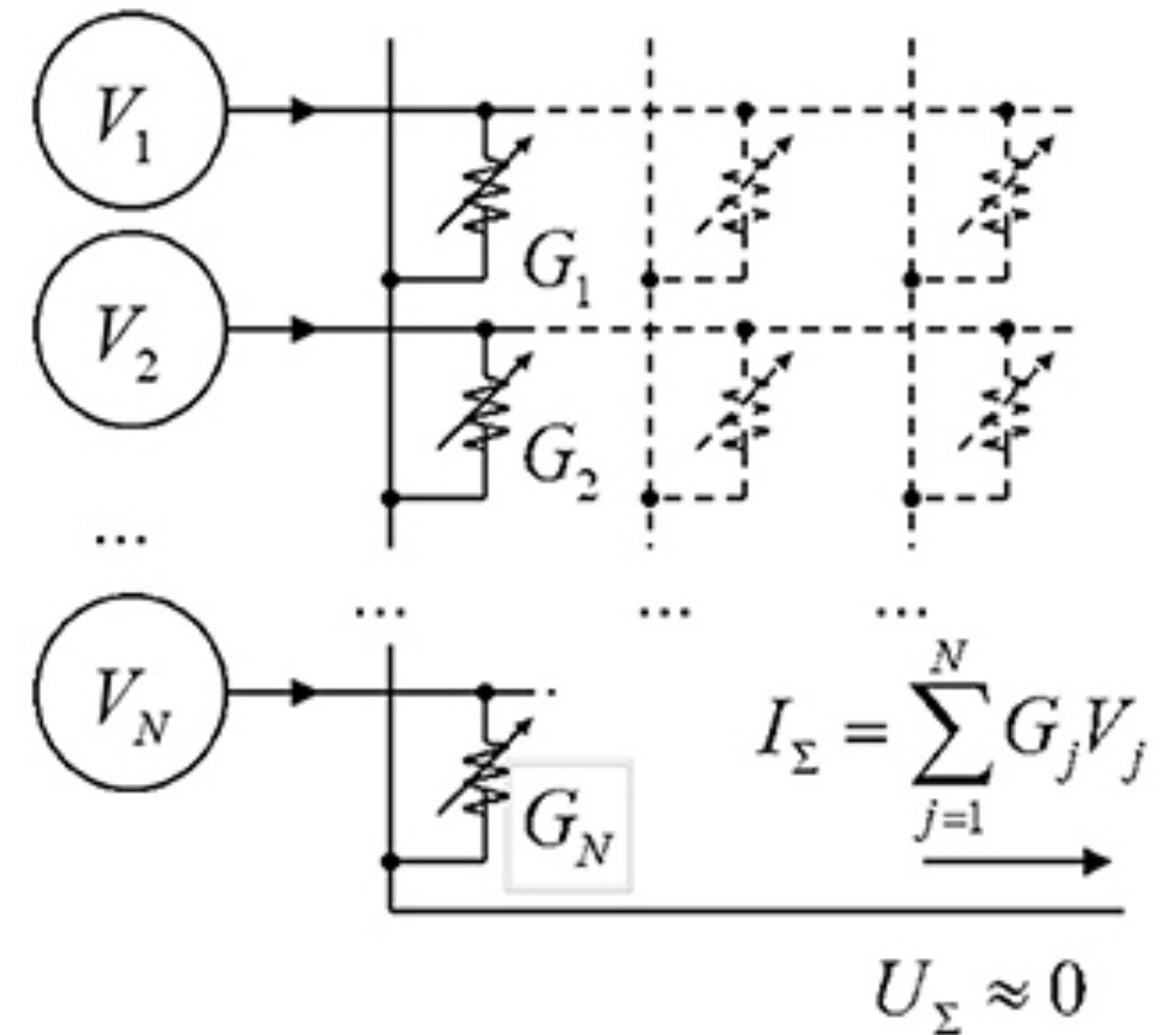- Emerging microelectronics technologies

# Emerging technologies

- Emerging computing architectures

- Emerging neural architectures

  - Spiking, inductive bias, physics-inspired,...

- Emerging microelectronics technologies



$$I_\Sigma = \sum_{j=1}^{N} G_j V_j$$

$$U_\Sigma \approx 0$$

# Summary

- A whirlwind tour through elements of developing embedded real-time ML!

- With hls4ml we try to make cutting edge techniques accessible to non-experts; open-source tools for scientific applications

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

- Powerful techniques exist

  - But there is still plenty of exciting research to do — ML techniques, computing architectures, microelectronics technologies

# Backup