



DESIGNING ARTIFICIAL INTELLIGENCE FOR EMBEDDED SYSTEMS

Introduction to Embedded AI

WELCOME



INTRODUCTION

Name

Home institution

Research

First time at IEEE RTC?



SCHEDULE

Time	Name	Instructor
9:30-12:30	Designing Artificial Intelligence for Embedded Systems	Audrey Corbeil Therrien
14:30-18:00	hls4ml	Benjamin Ramhorst
9:30-12:30	Coyote v2: Open-source Abstractions and Infrastructure for FPGAs	Benjamin Ramhorst
14:30-18:00	Super Neural Architecture Codesign Package (SNAC-Pack)	Dmitri Demler

DESIGNING AI FOR EMBEDDED SYSTEMS

- What is AI?
- Hardware
- Co-design
- Squish
- Validation
- Overview of tools

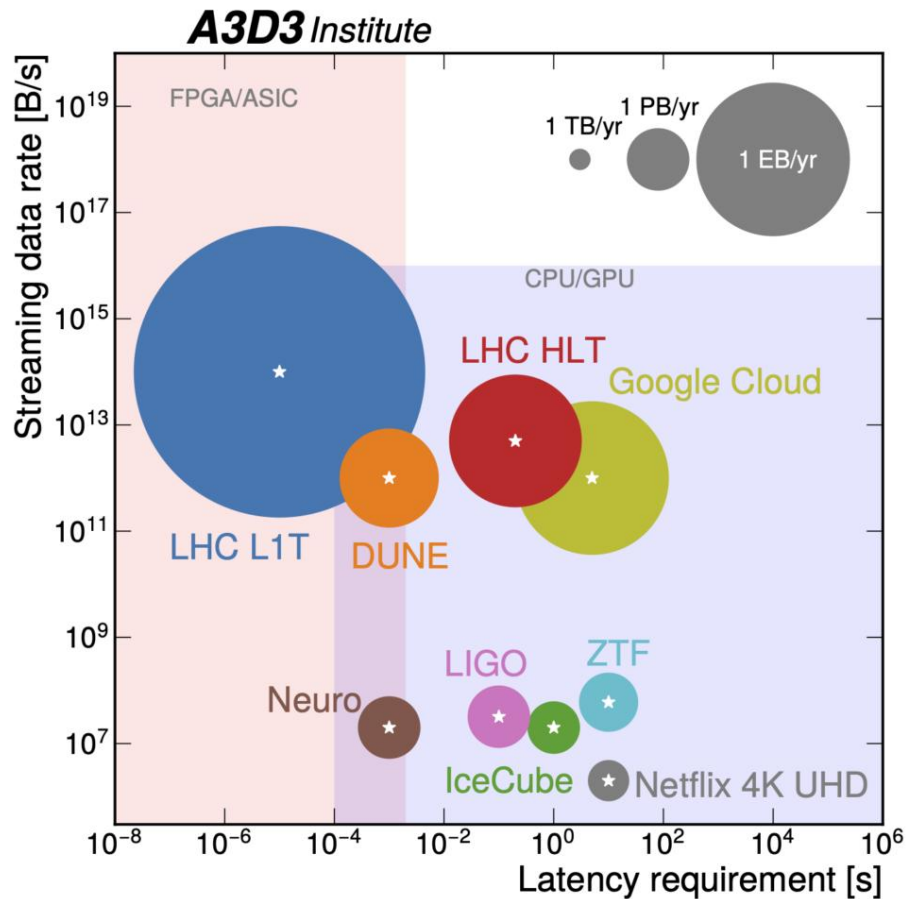
Univers en soi

UNIVERS EN SOI

WHAT IS AI

Univers en soi

UNIVERS EN SOI



Artificial Intelligence

Machine Learning

Neural Networks

MLP

Deep Learning

Transformers

Autoencoders

GNN

RNN

Generative

LLM

AI

Genetic Algorithms

K-Means

Planning

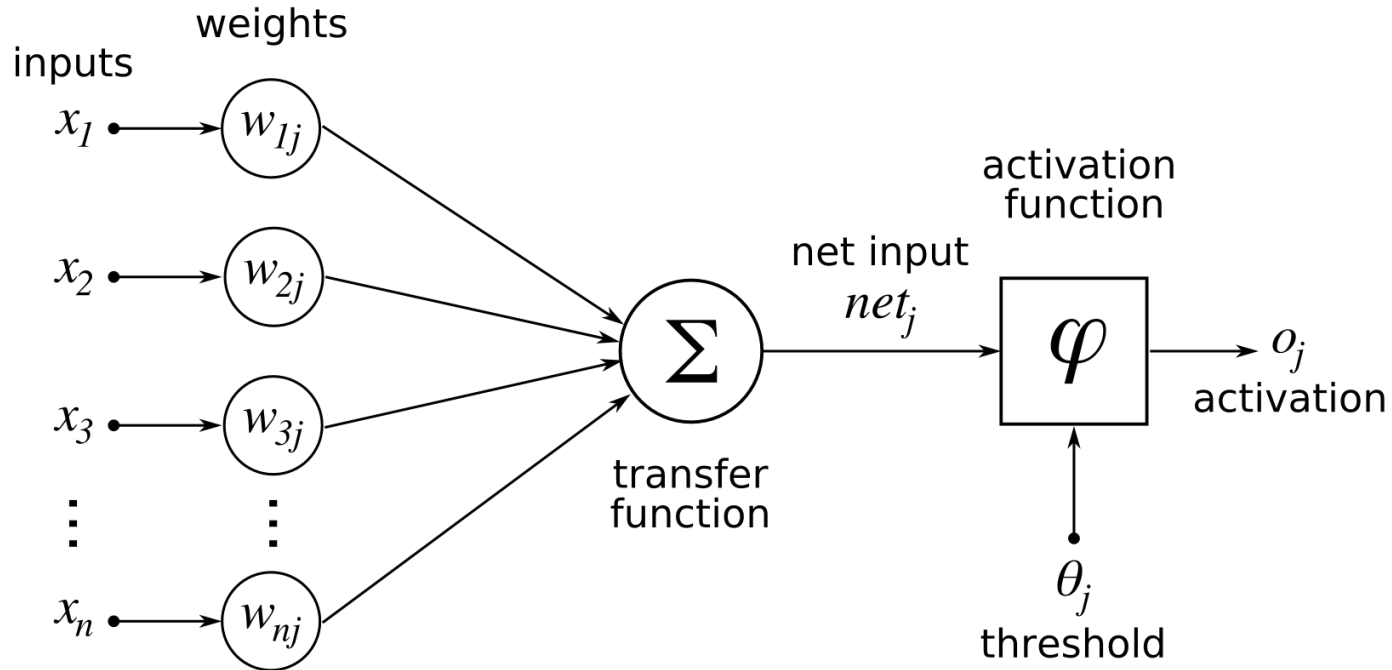
First Order Logic

Expert Systems

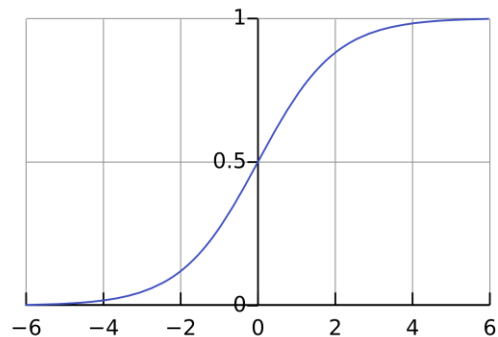
Bayes Classification

Fuzzy Logic

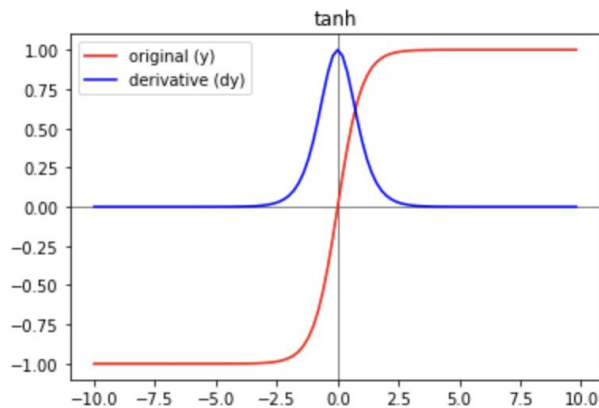
NEURON



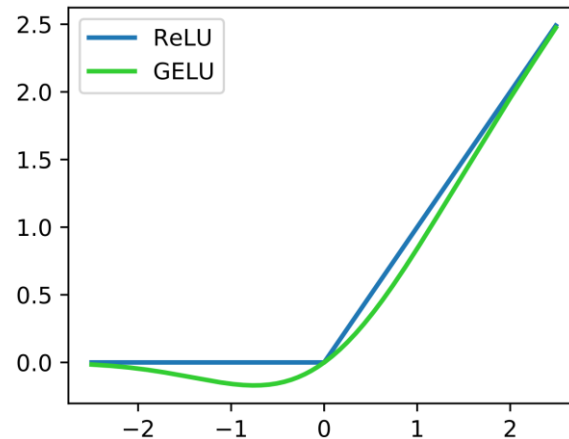
ACTIVATION FUNCTIONS



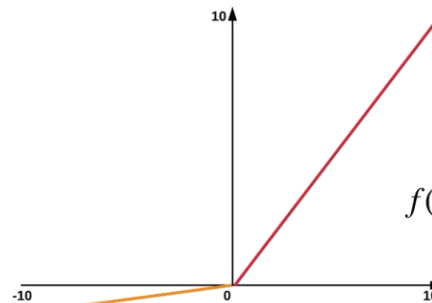
Logistic/sigmoid



Nonlinearities

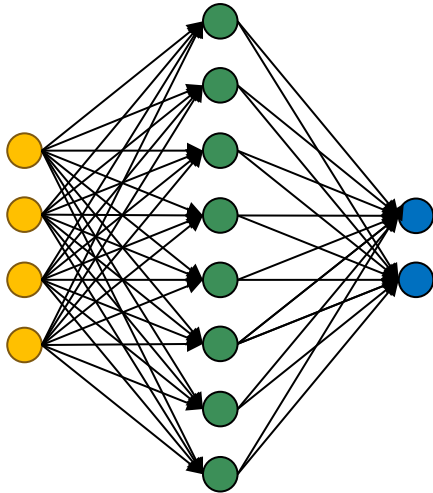


Leaky ReLU Activation Function



$$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

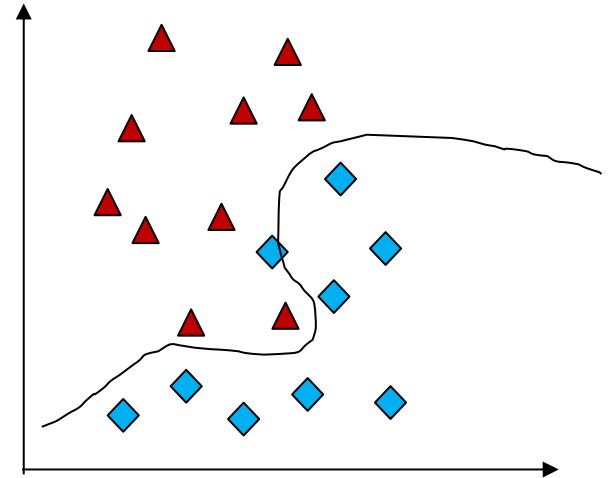
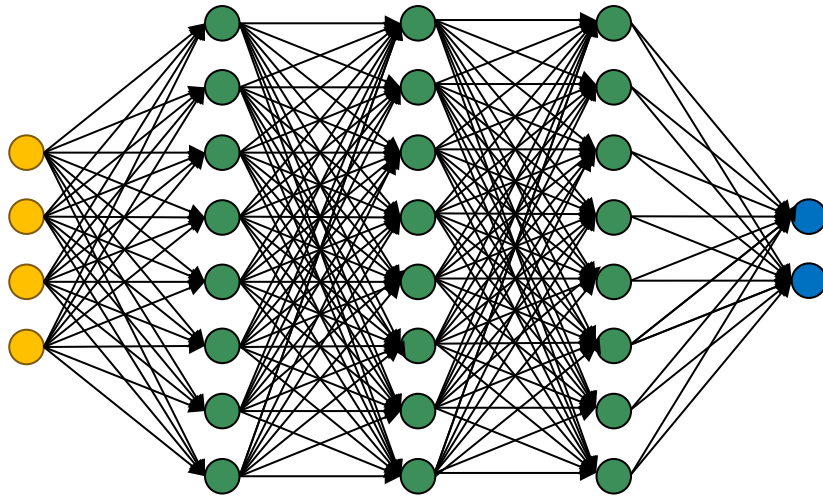
NEURAL NETWORK



$$f(x) = \varphi(\mathbf{w} \cdot \mathbf{x} + b)$$

$$f(x) = \varphi \left(\left(\sum_{i=1}^n w_i x_i \right) + b \right)$$

NEURAL NETWORK



CONVOLUTIONAL NEURAL NETWORK



chocolate cookie



fawn smooth Chihuahua



baked blueberry muffin



white chihuahua



fawn smooth Chihuahua



blueberry muffin



fawn smooth Chihuahua



muffin



[unknown]



brown coated Chihuahua



baked muffin



beige short coated puppy



tan smooth Chihuahua puppy



blueberry cupcakes



three smooth Chihuahua puppies



white and black muffin

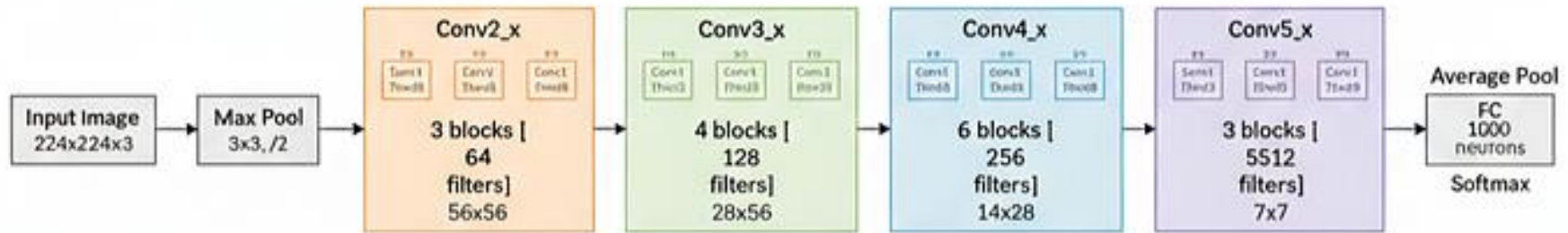
CONVOLUTIONAL NEURAL NETWORK

6	5	2	0	4	1
5	4	3	2	3	4
4	2	1	0	1	1
3	2	1	0	0	0
1	1	1	1	3	4

3	2	1
2	1	0
1	0	0

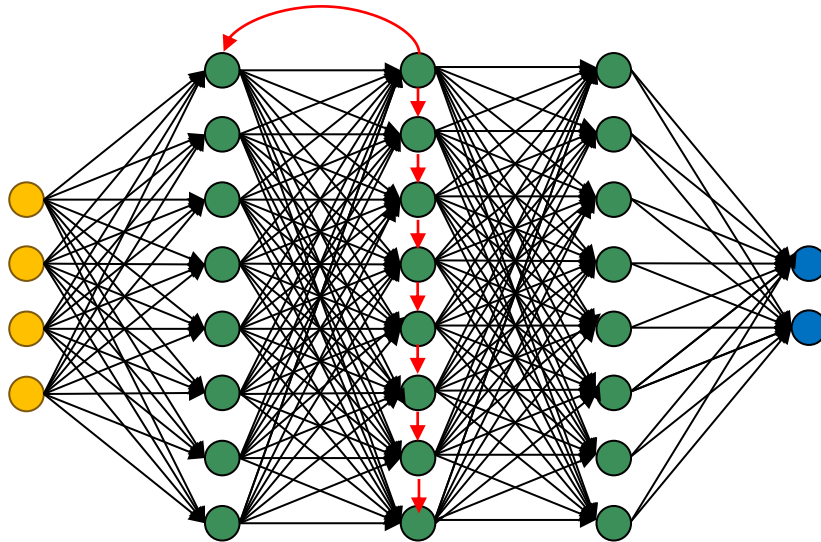
48	33	23	16
39	33	18	17
26	12	7	4

CONVOLUTIONAL NEURAL NETWORKS

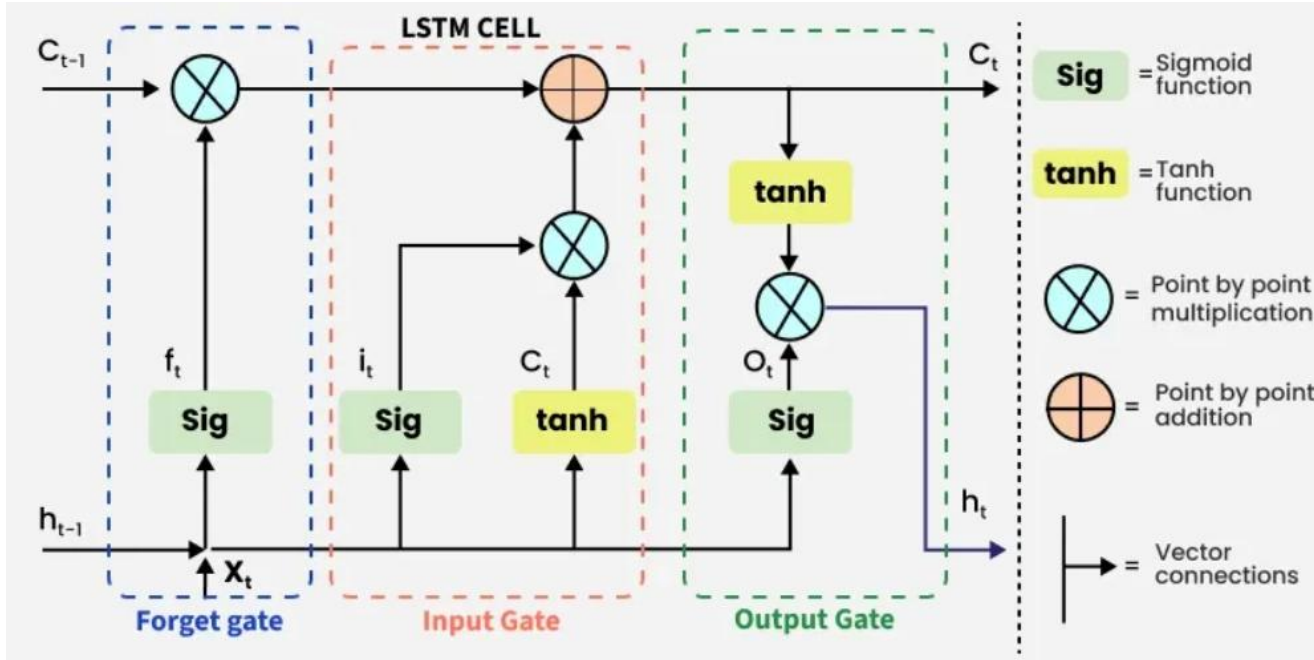


Resnet-50 - ~2016

RECURRENT NEURAL NETWORKS

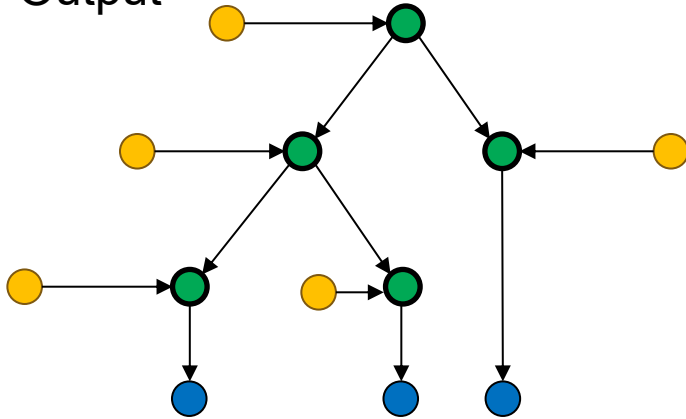


RECURRENT NEURAL NETWORKS

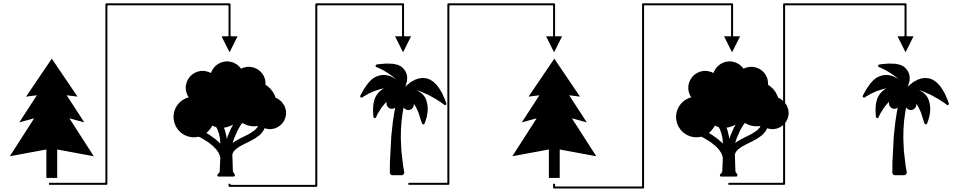


DECISION TREE

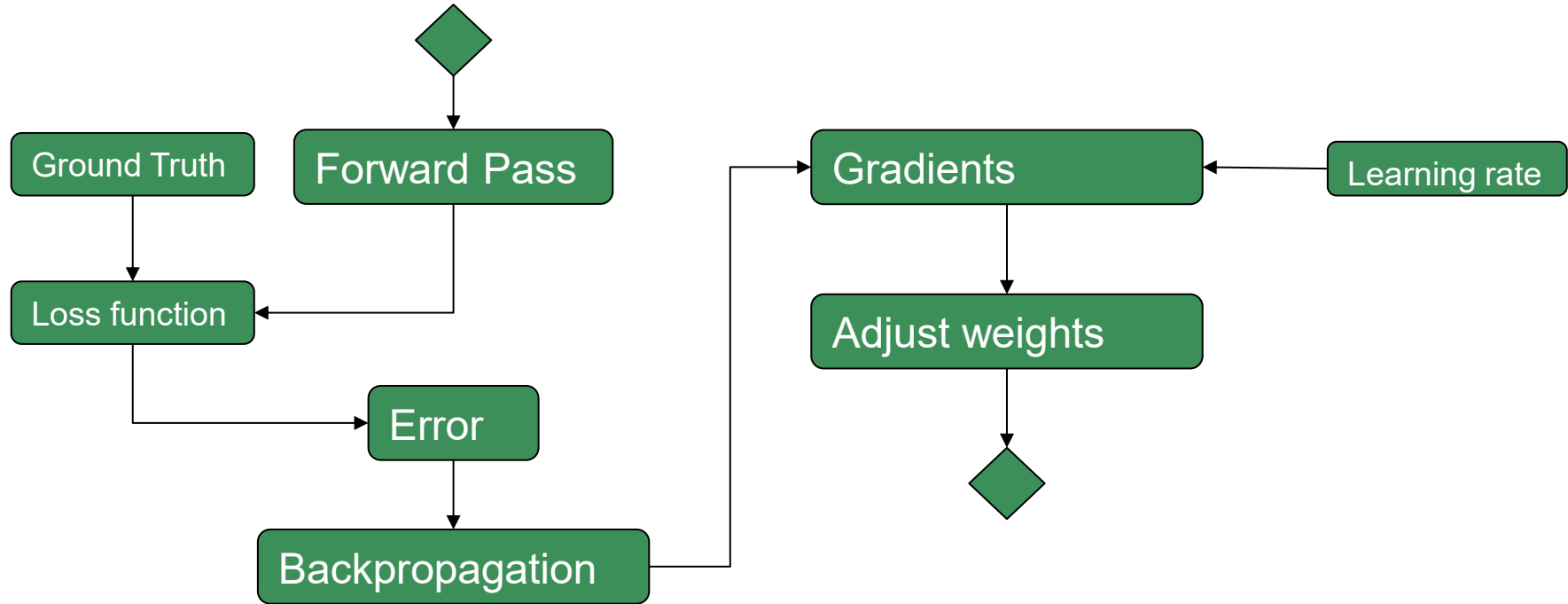
- Input
- Decision
- Output



Gradient boosting Decision trees



TRAINING



DATASETS

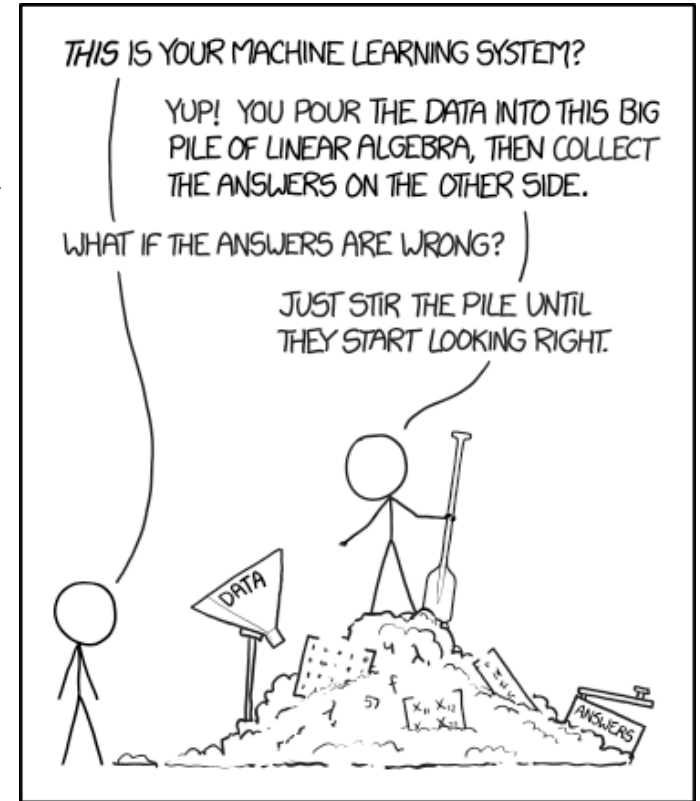
- Source
- Accuracy
- Authenticity
- Scope
- Usability



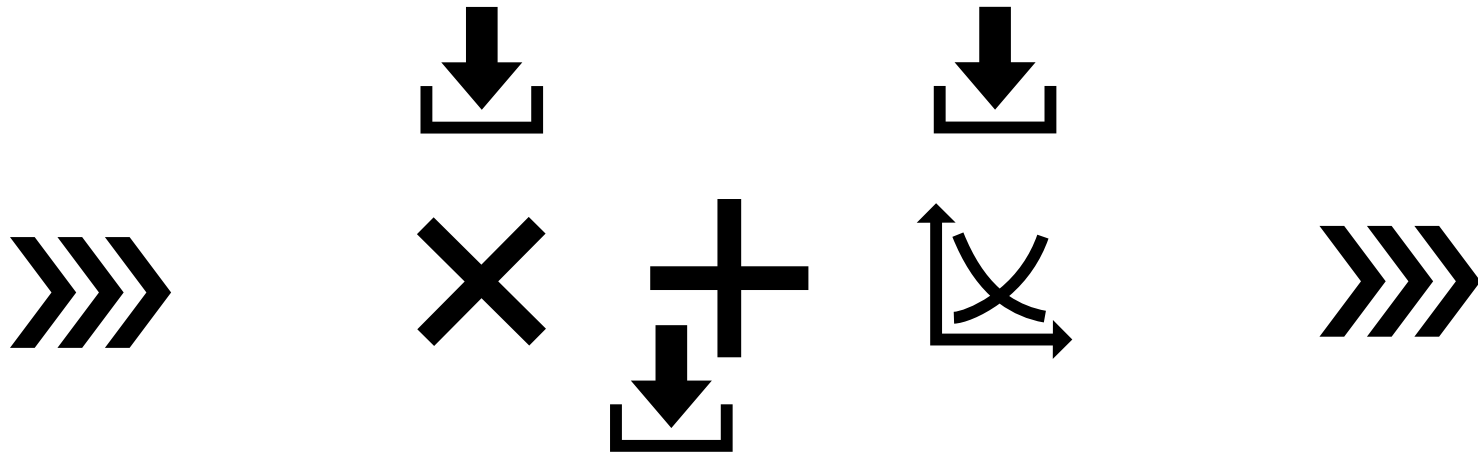
MODELS

- A model is only as good as its data
 - (at best...)

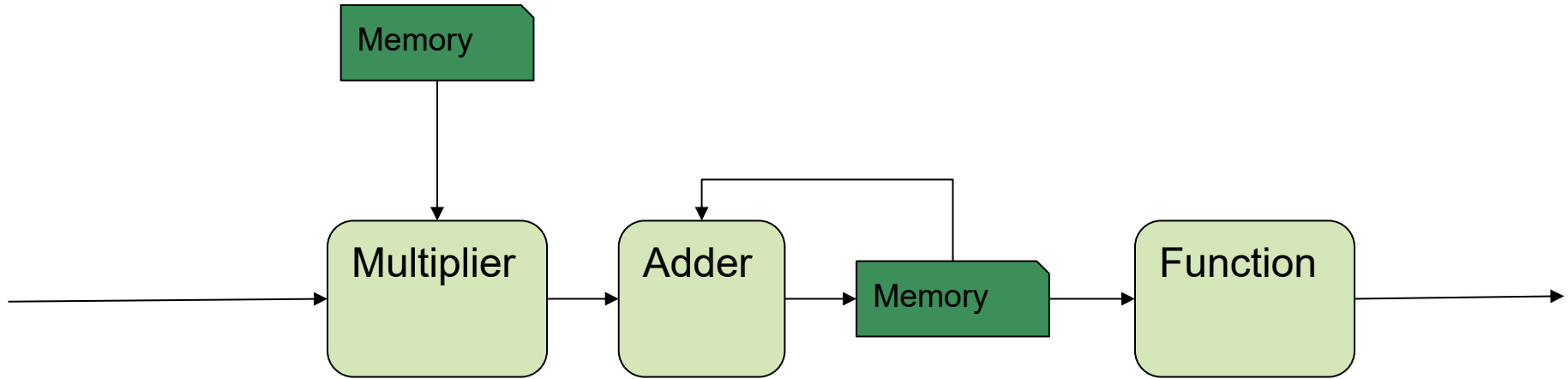
- “All models are wrong, but some are useful.”
 - George Box



OPERATIONS OF A NEURAL NETWORK



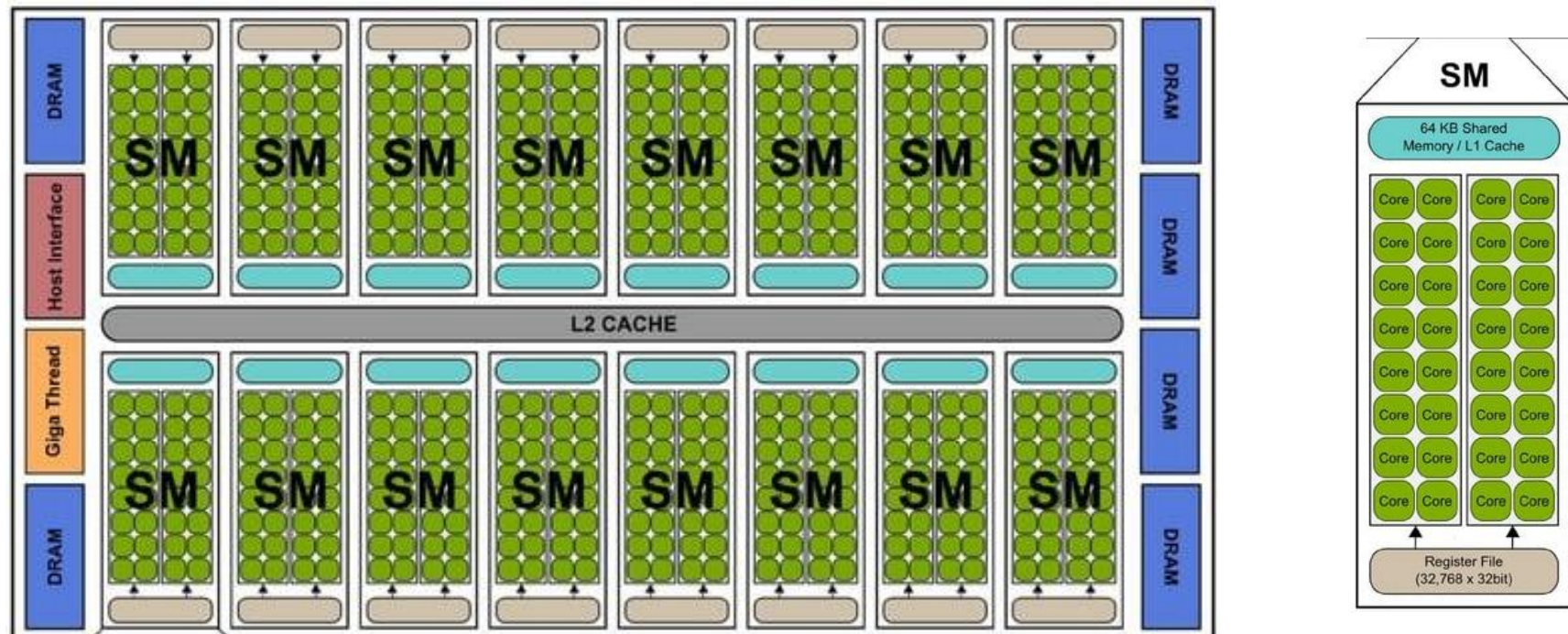
OPERATIONS OF A NEURAL NETWORK





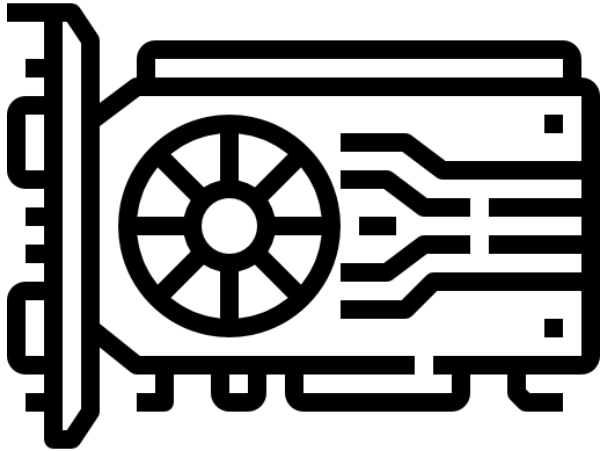
HARDWARE

GRAPHICS PROCESSING UNIT



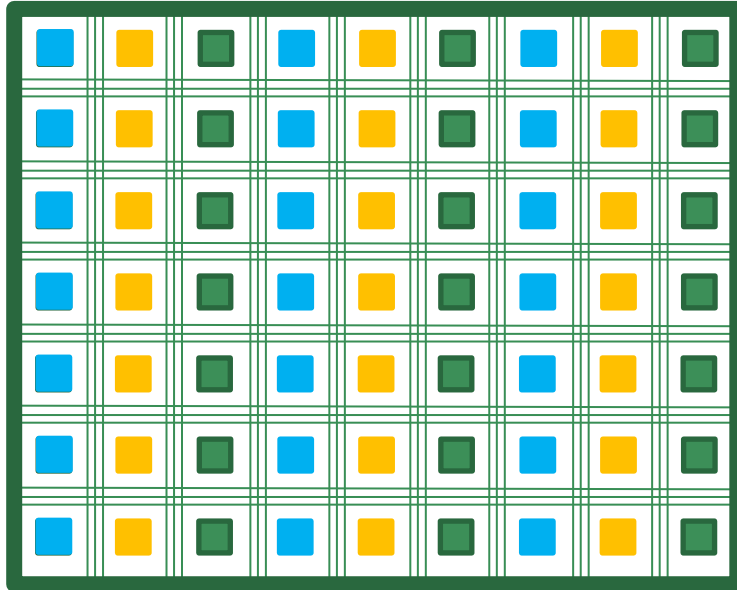
Source : Modified from
https://www.researchgate.net/publication/236666656_Accelerating_Fibre_Orientation_Estimation_from_Diffusion_Weighted_Magnetic_Resonance_Imaging_Using_GPUs/figures?lo=1

GRAPHICS PROCESSING UNIT



- Highly parallel
- Flexible
- Well supported
- Limited I/O
- Large footprint
- High power usage

FPGA



- Logic
- Memory
- Digital signal processing slices

Reconfigurable

Efficient

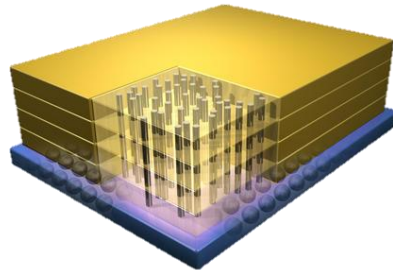
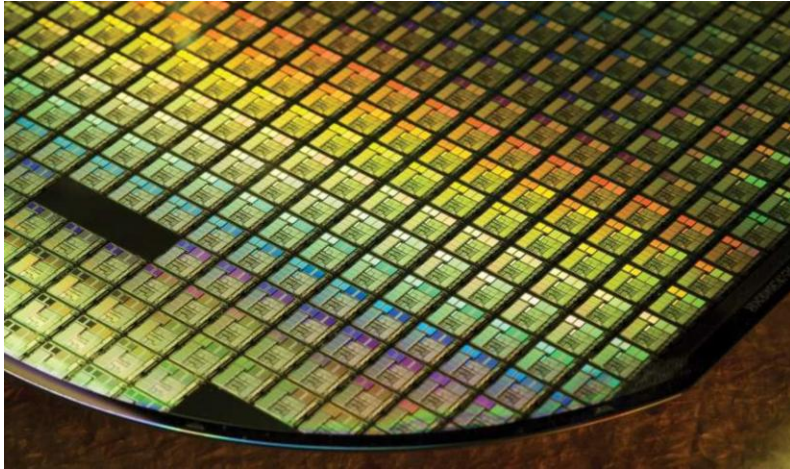
I/O capacity

Programming

Limited clock

Limited resources

ASIC



Efficient++

Custom I/O

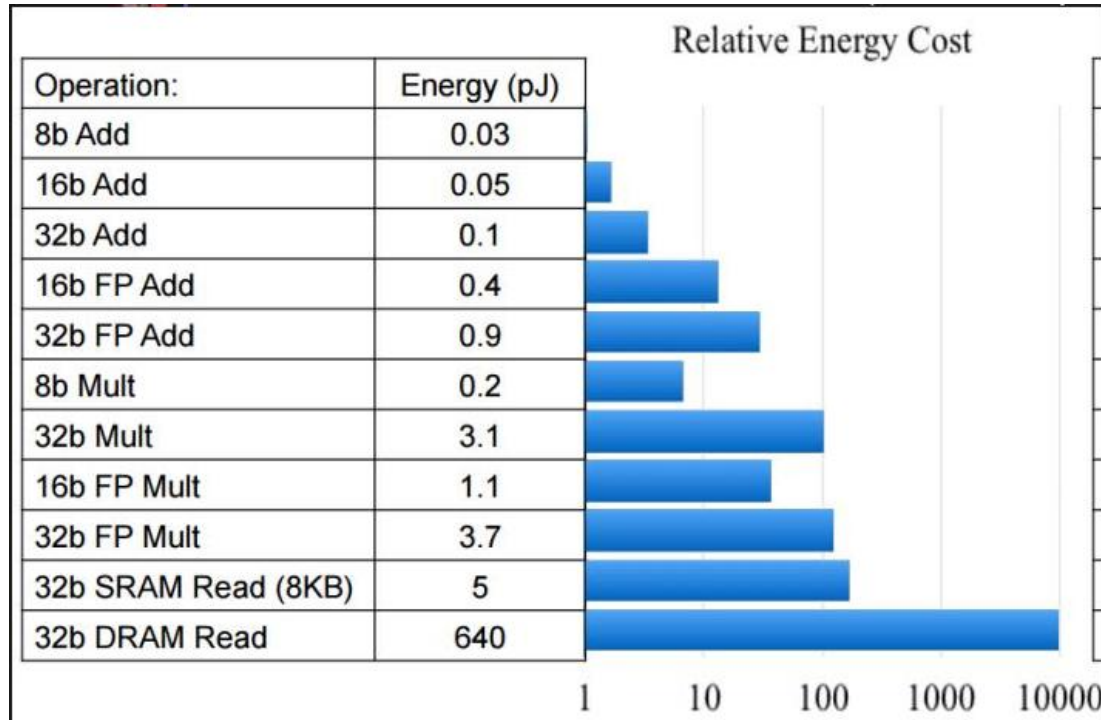
3DIC

Reconfigurable

Expensive

Long design cycle

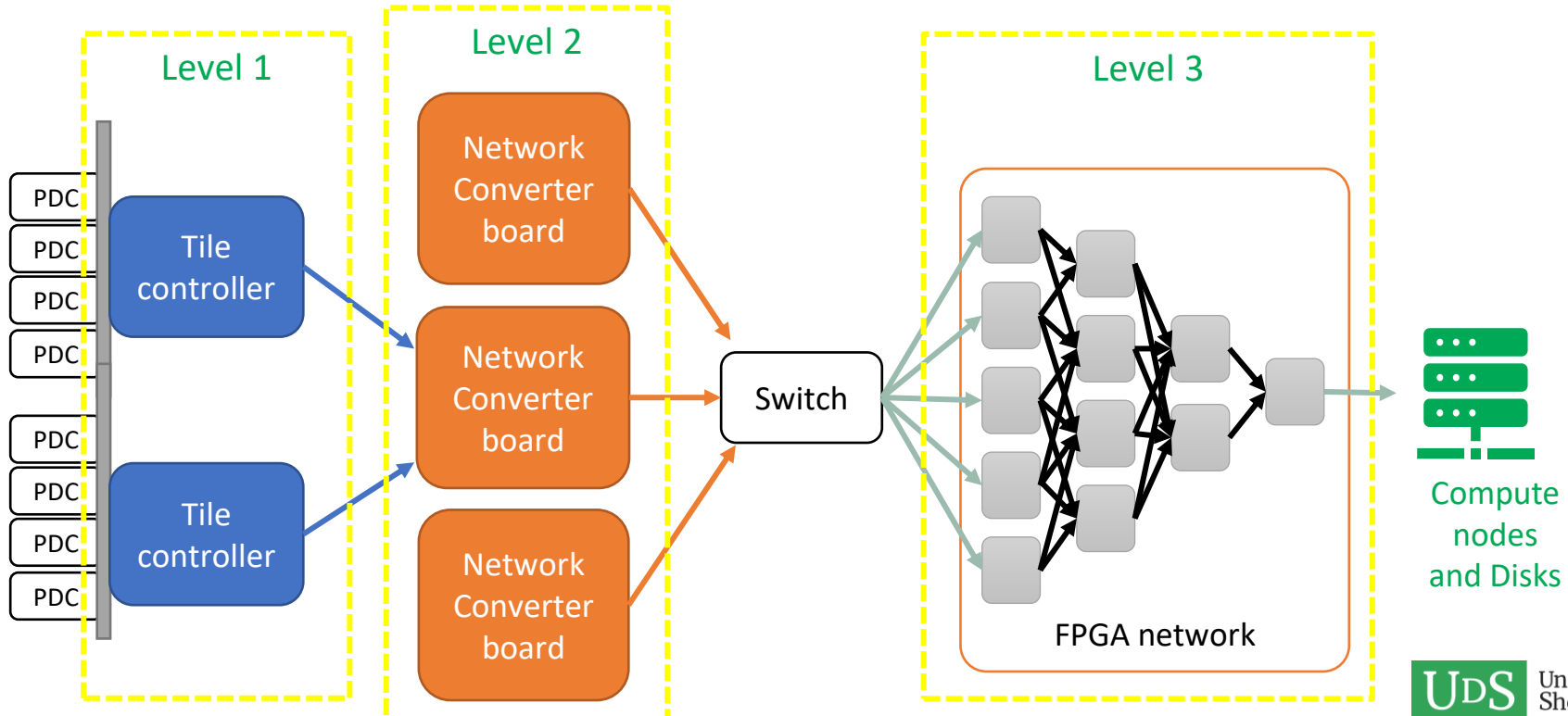
ENERGY



EFPGA

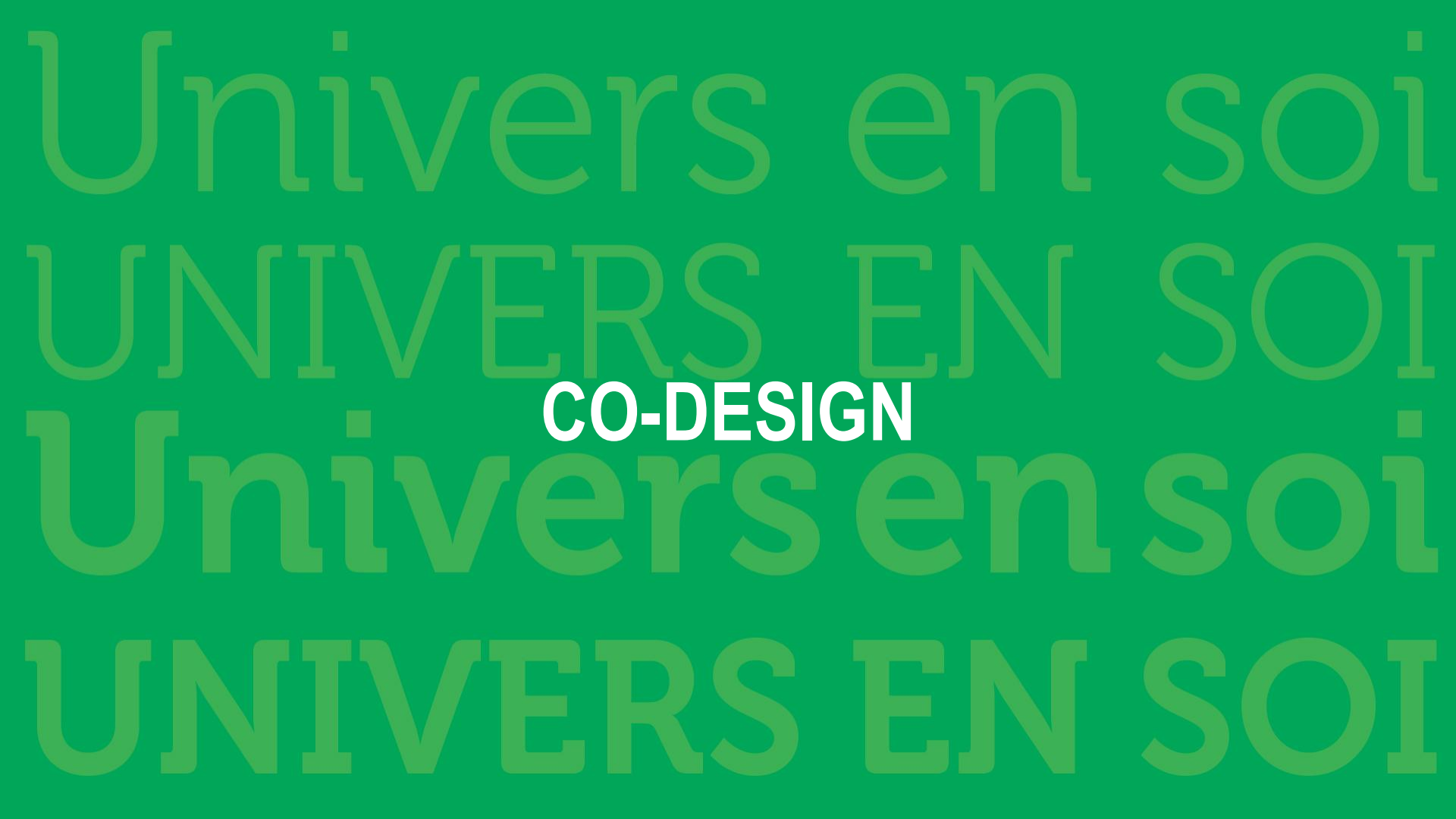
- Embedding FPGA logic in an ASIC designs
 - Combine the low power and interconnectivity of ASIC with the reconfigurability of FPGA
 - More options to update the ML architecture
 - Pre-determined dimensions
 - Can be expensive depending on technology node

DISTRIBUTED COMPUTE

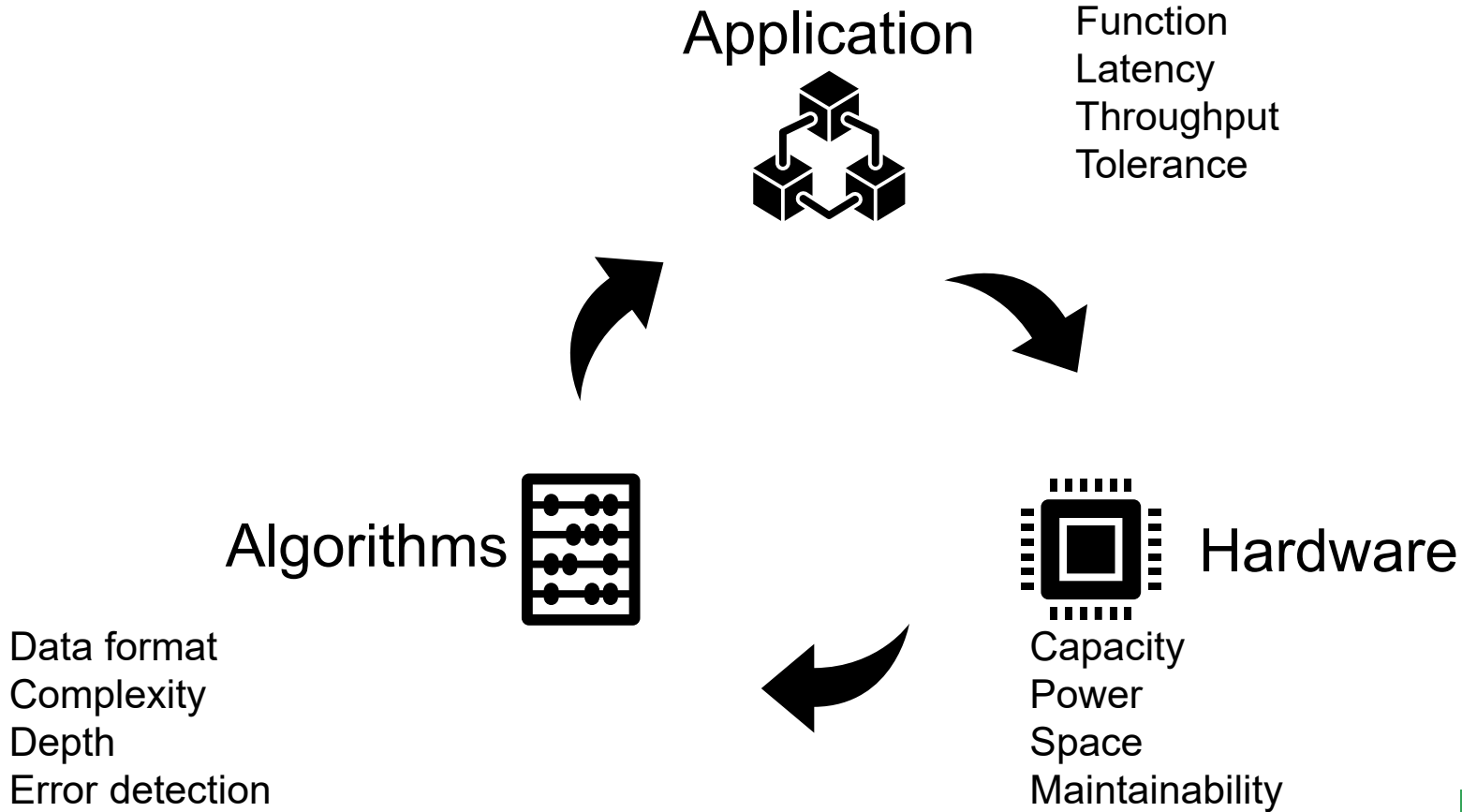


DISTRIBUTED COMPUTE

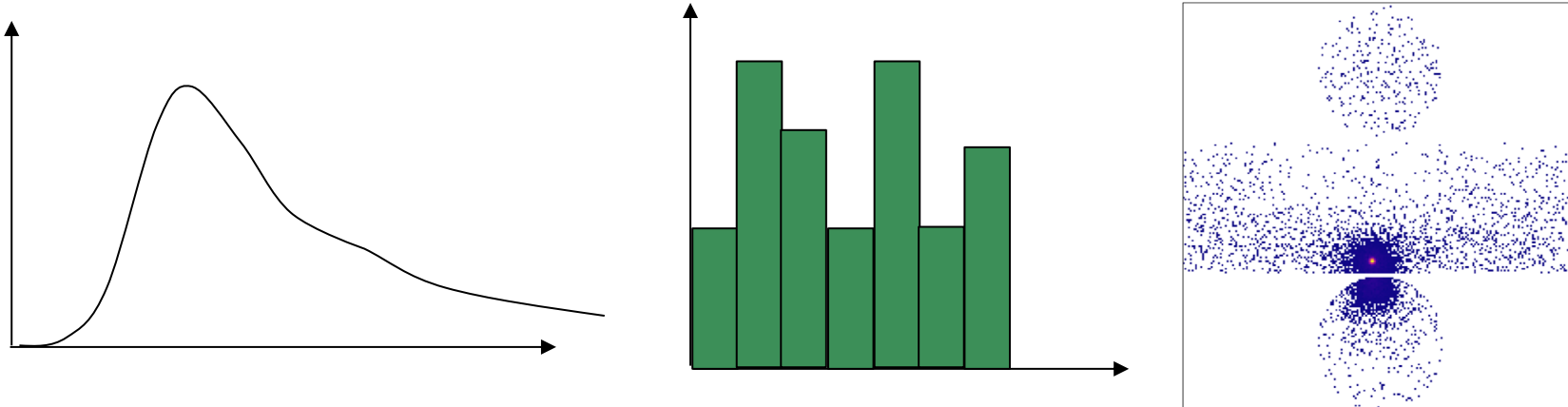
- Look for unused compute potential
 - ASIC
 - Network Switches
 - Data movers
- Some model architectures reduce data
 - Lower requirements on data links
- Tools are under developement



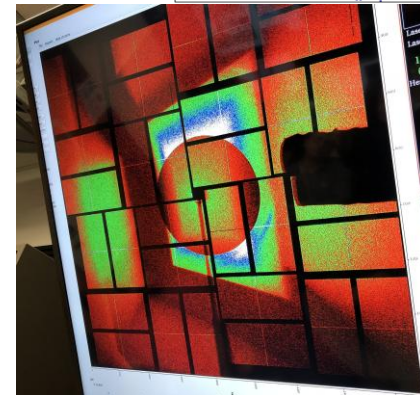
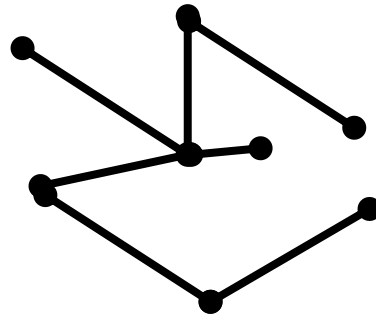
CO-DESIGN



INPUT-OUTPUT



Feature	Value
Energy	412
Position	(8,6,4)
Fprompt	0,56
Direction	(1,1,6)
Color	Blue

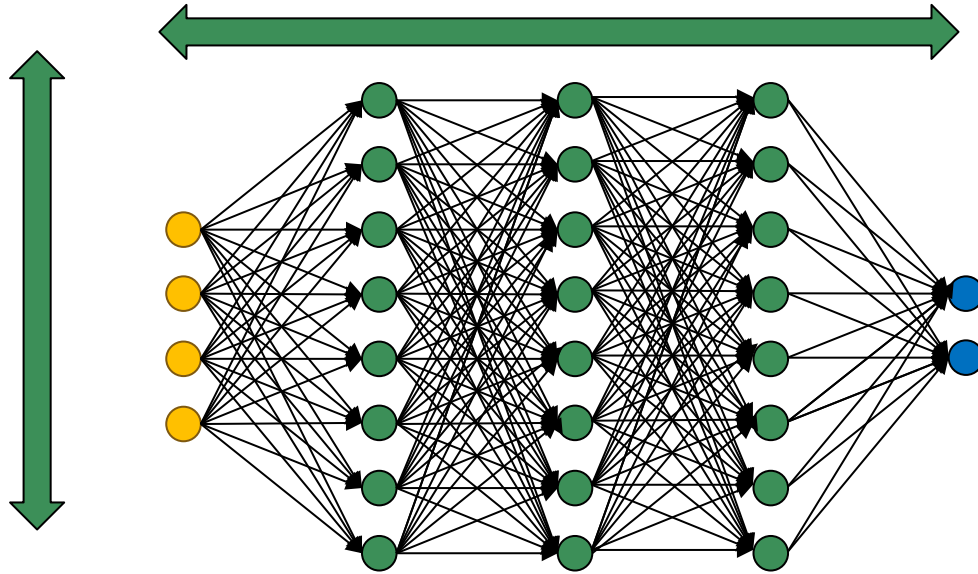


INPUT/OUTPUT

- Type
- Dynamic range
 - Normalisation
- Representation
- Transformation?

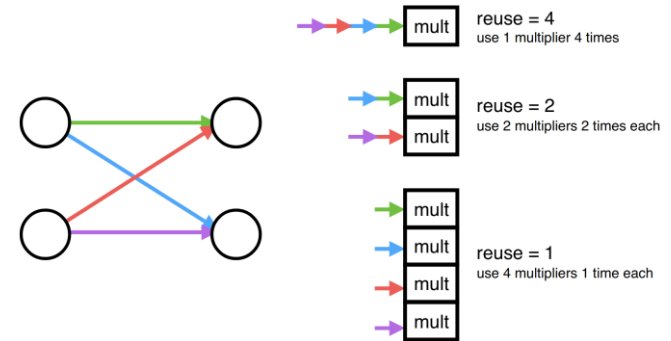
Type	Range
FP64/double	1.7E +/- 308 (fifteen digits)
FP32/float	3.4E +/- 38 (seven digits)
INT32/ long int	-2,147,483,648 to 2,147,483,647
INT16/ Short int	-32,768 to 32,767
Unsigned INT8	0 to 255

STRUCTURE/DIMENSIONS



2,4,8,16,32,64,128.

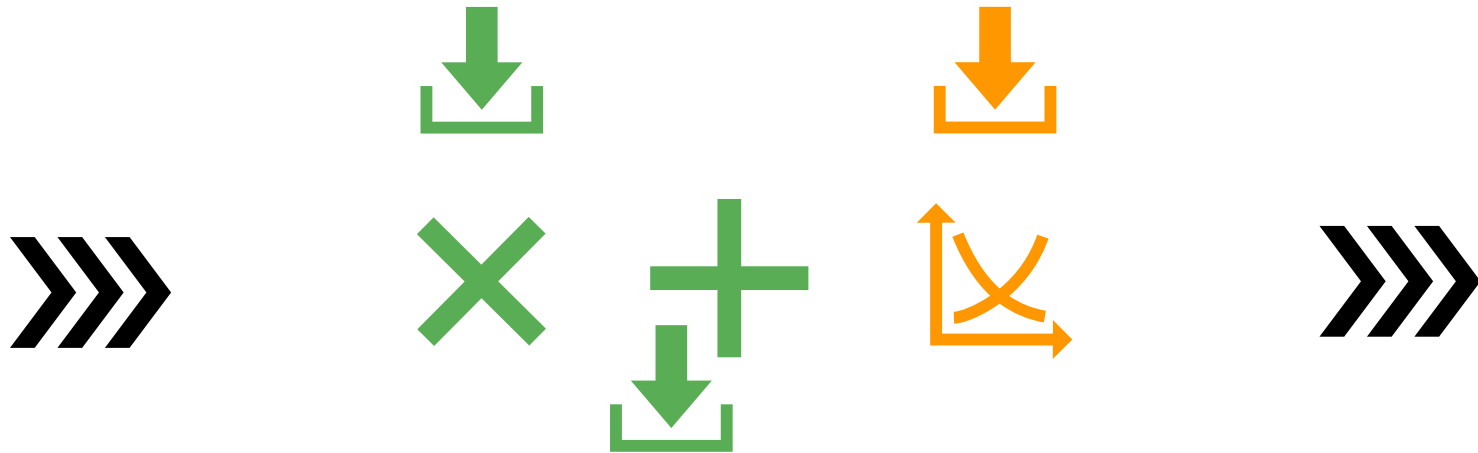
Unrolled Vs Reuse



Hls4ml

<https://fastmachinelearning.org/hls4ml/concepts.html>

OPERATIONS



ACTIVATION FUNCTION

- Calculation is often expensive
- Look-up table may be memory intensive
- Not all activations are implemented

- Use the same activation if you can
- Retrain with a simpler activation function

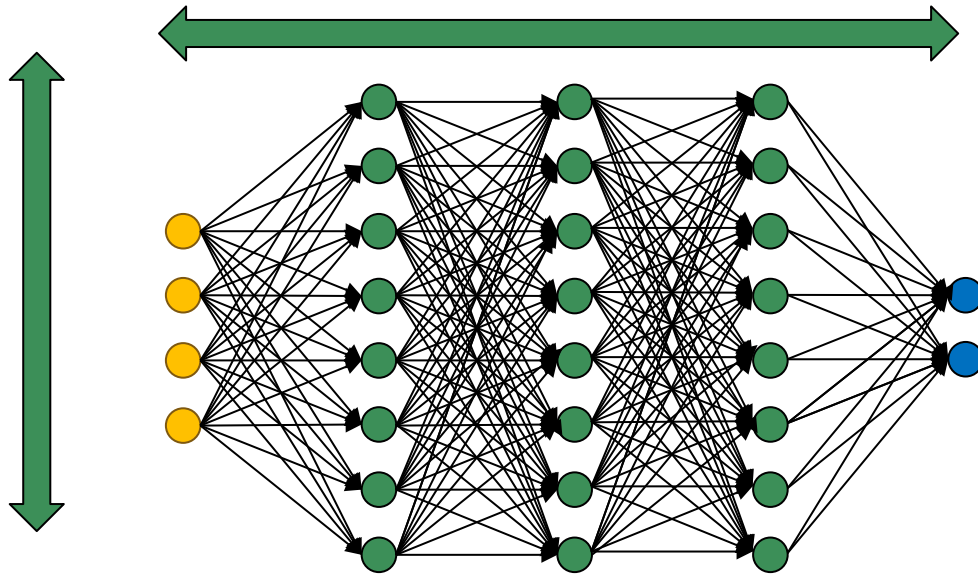


MODEL STRUCTURE

- Not all operations are implemented by tools
- Non-standard architectures can fail to synthesize, or have unexpected behavior

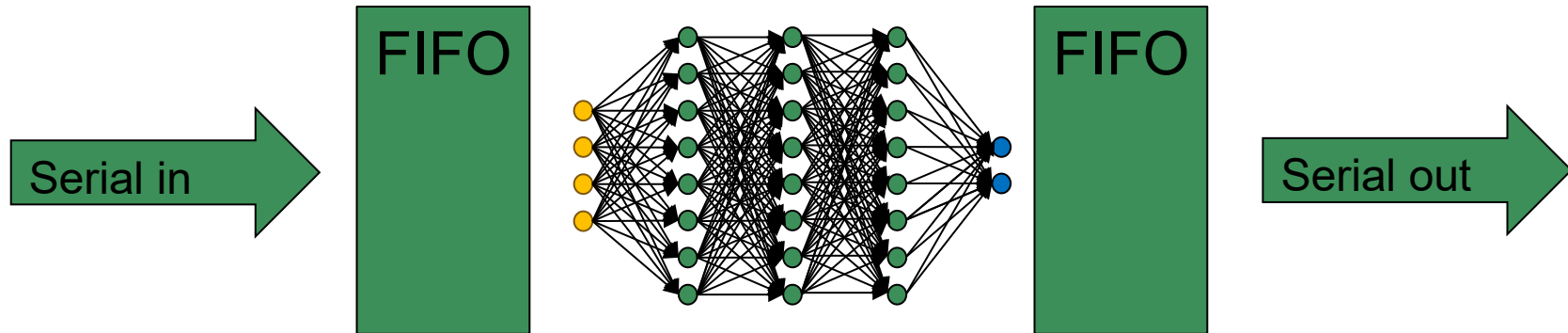
DATA FLOW

- Latency vs throughput



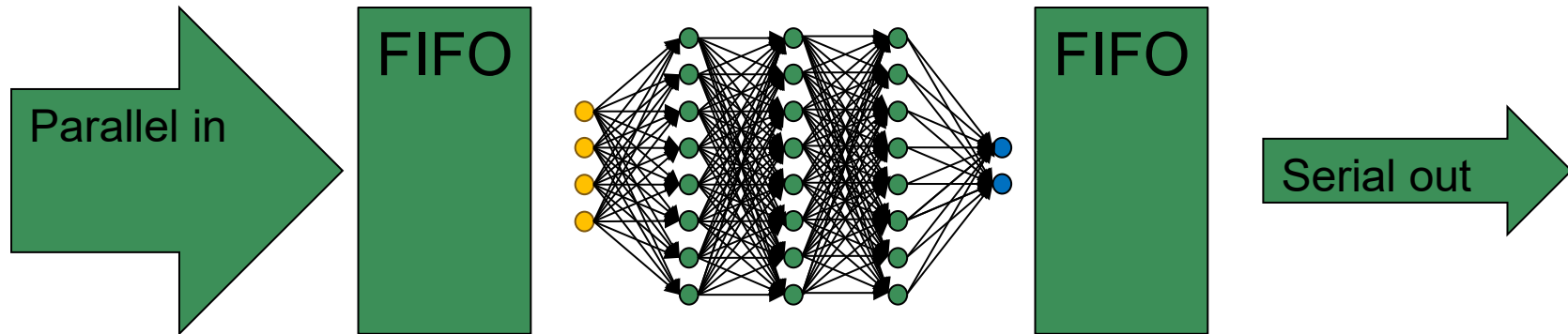
DATA FLOW

- Buffers
- Communication



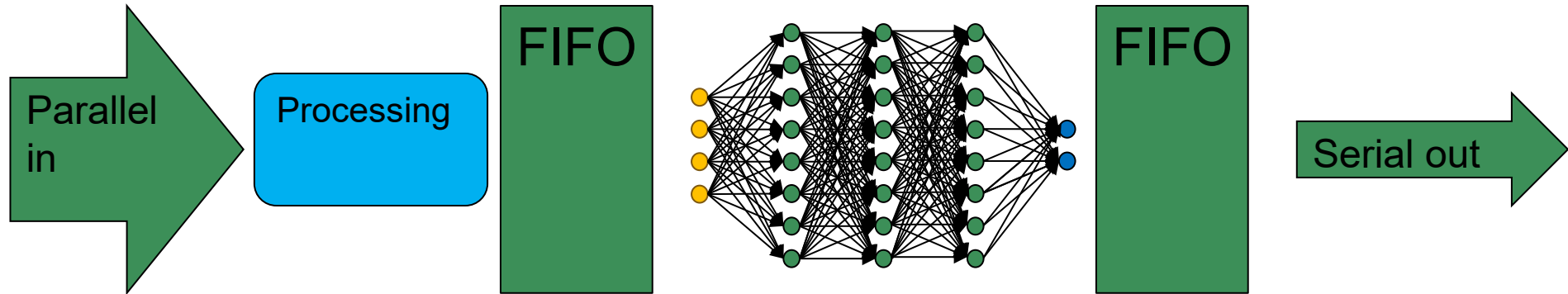
DATA FLOW

- Buffers
- Communication



DATA FLOW

- Buffers
- Communication



SUMMARY

- Input-Output
- Dynamic range
- Structure/dimensions
- Operations
- Data flow requirements

- How will this work in hardware?

Univers en soi

UNIVERS EN SOI

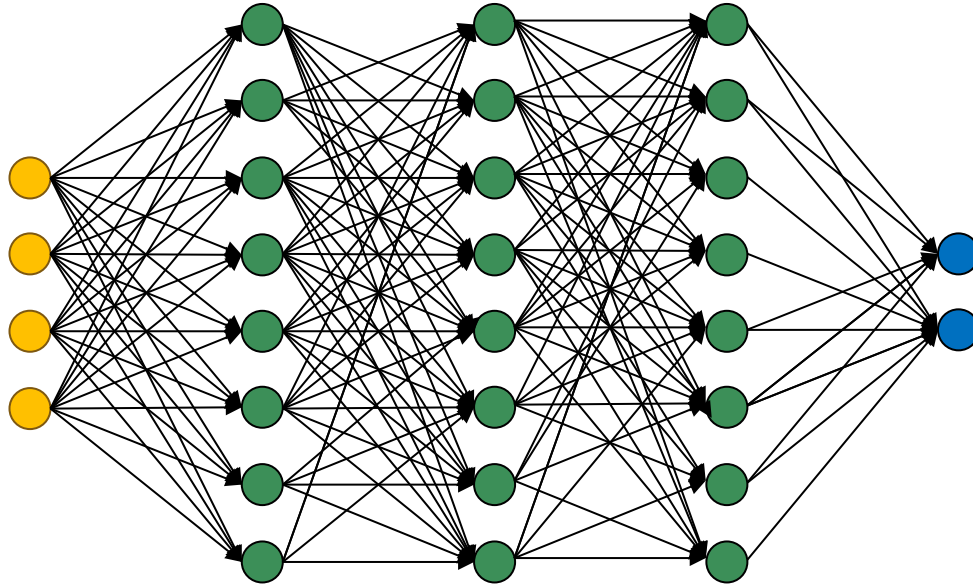
SQUISH

Univers en soi

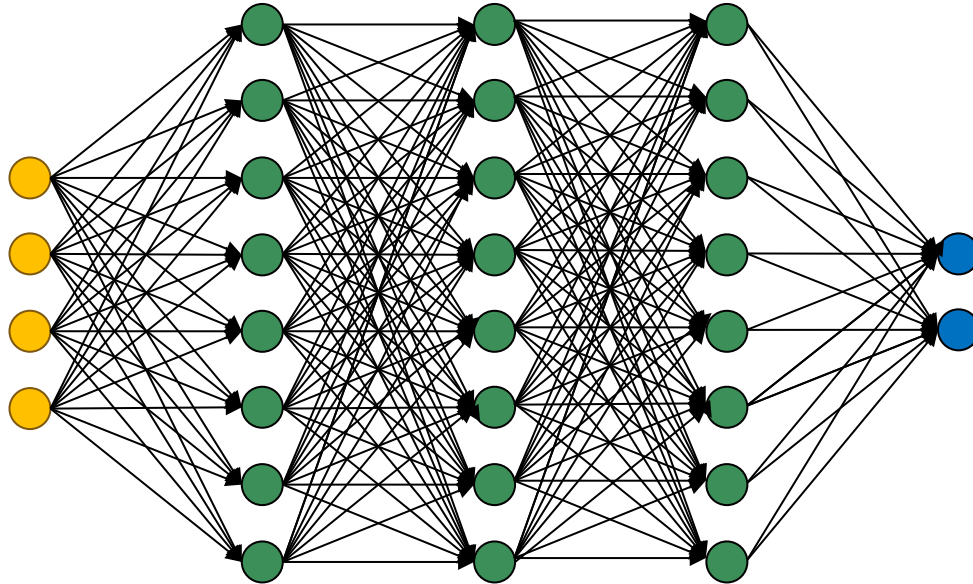
UNIVERS EN SOI

PRUNING

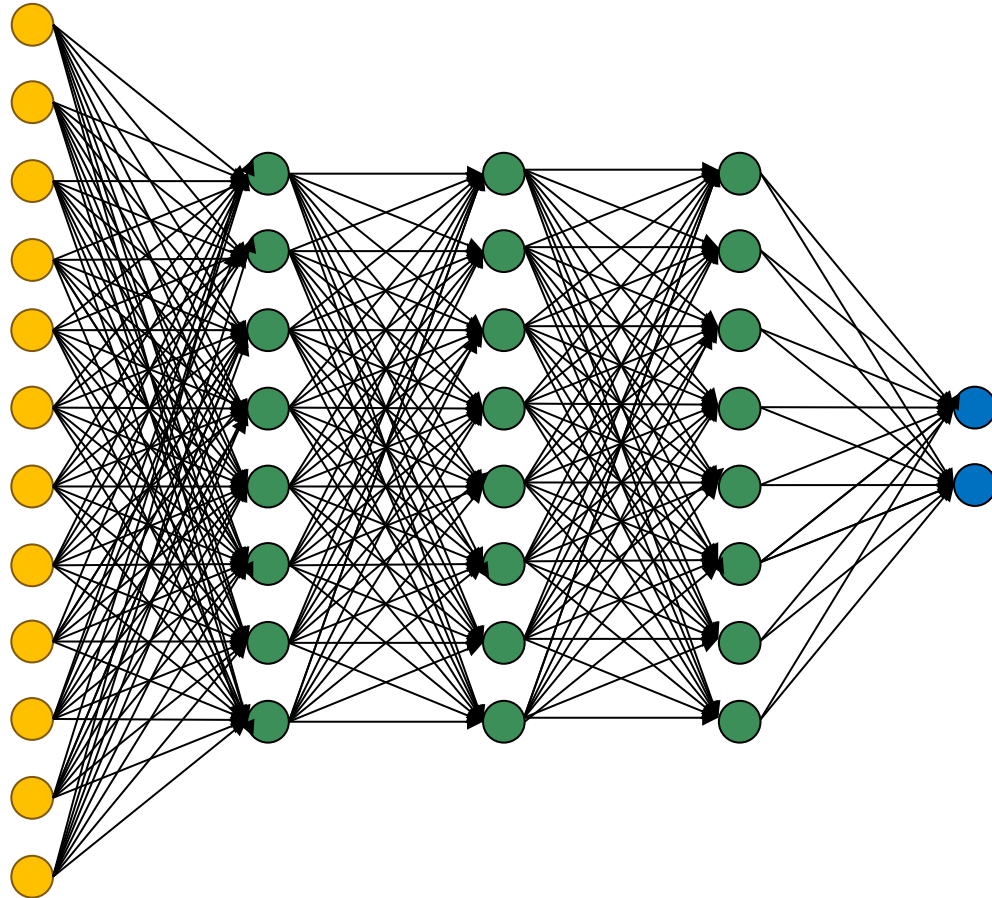
!! Logic overhead



NODES AND LAYERS

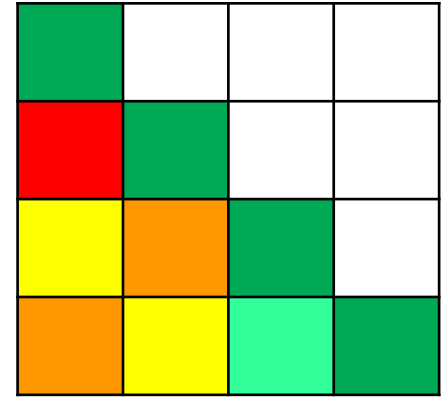
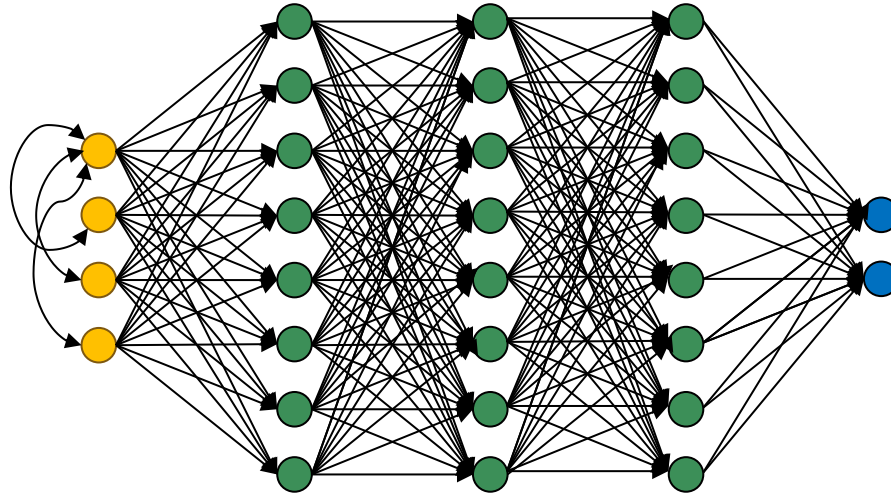


INPUTS



INPUTS – FEATURE IMPORTANCE

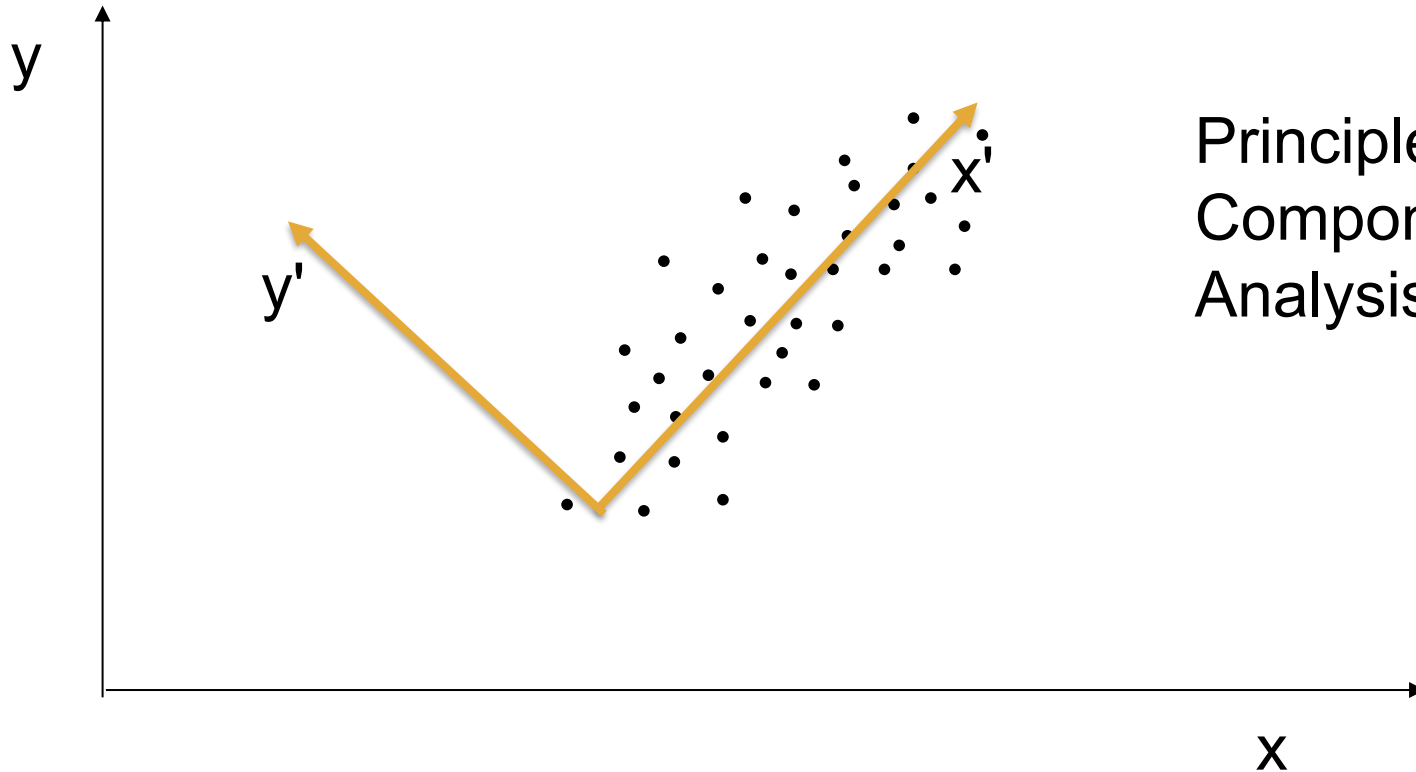
- Permutation
 - Swapping inputs



INPUTS – FEATURE IMPORTANCE

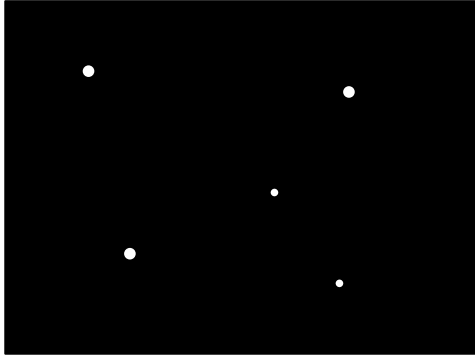
- Feature Ablation
 - Remove, retrain, compare
- SHAP (SHapley Additive exPlanations)
 - <https://shap.readthedocs.io/en/latest/>

INPUT DECORRELATION

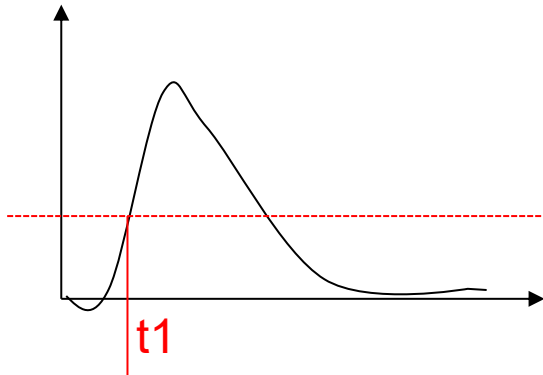


Principle
Component
Analysis

FEATURE EXTRACTION



(15,10,3)
(10,45,3)
(30,20,1)
(40,8,1)
(42,42,3)



If $f(t) > \text{threshold}$
then t

QUANTIZATION - WEIGHTS

- Reduces memory footprint
- Allows simpler multipliers/adders
- Reduces latency and power

- Careful of dynamic range!
 - Overflow
 - Representation

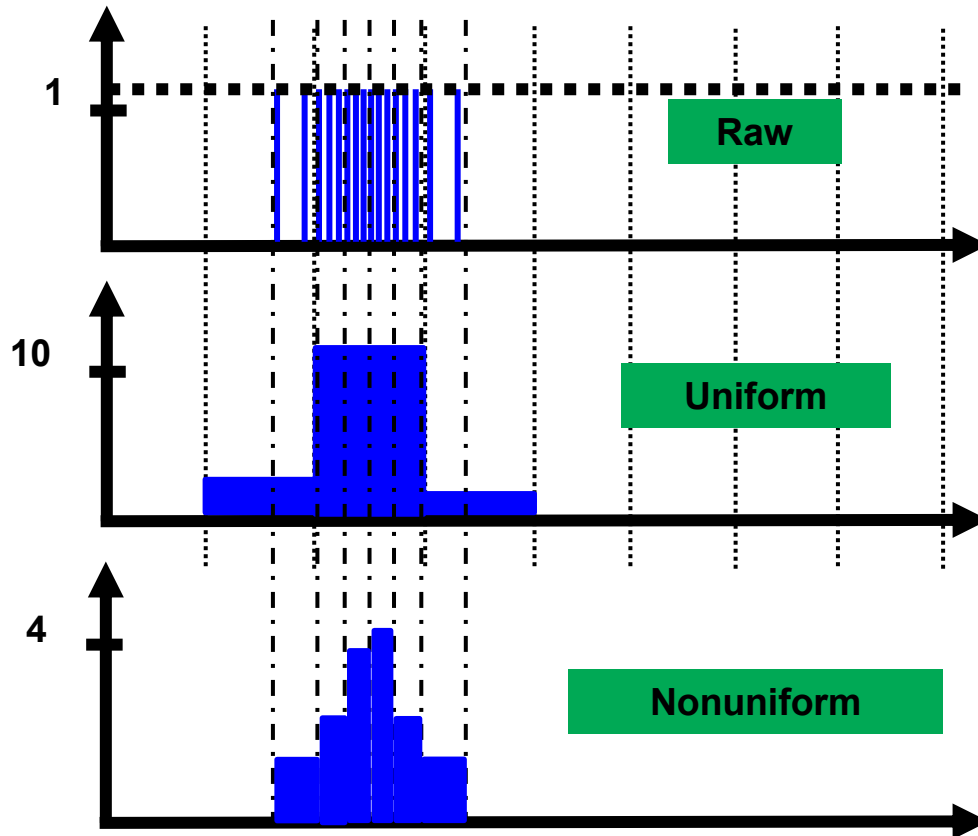
FP32 to INT8

4x reduction in memory

20x reduction in power

Reduction in clock cycles

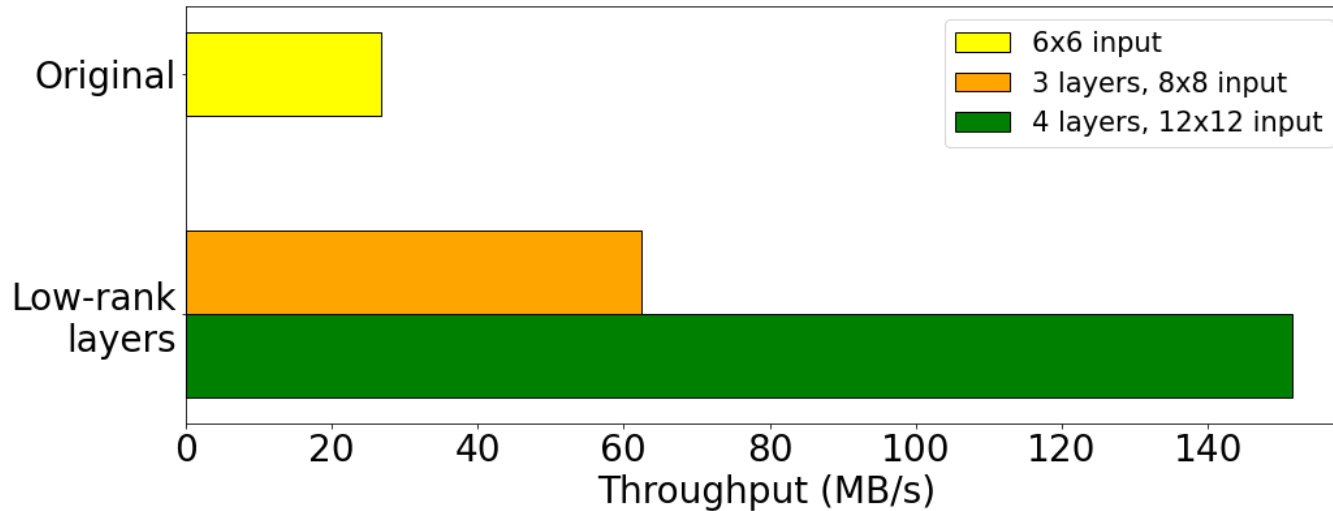
QUANTIZATION - INPUTS



- ✓ Optimized for the data probability distribution function (PDF)
- ✓ Minimizes the quantization error
- ✓ Reduces our input shape dimensionality

MATRIX FACTORISATION

- Reduce the size of a layer by decomposing into smaller matrices.



MATRIX FACTORISATION

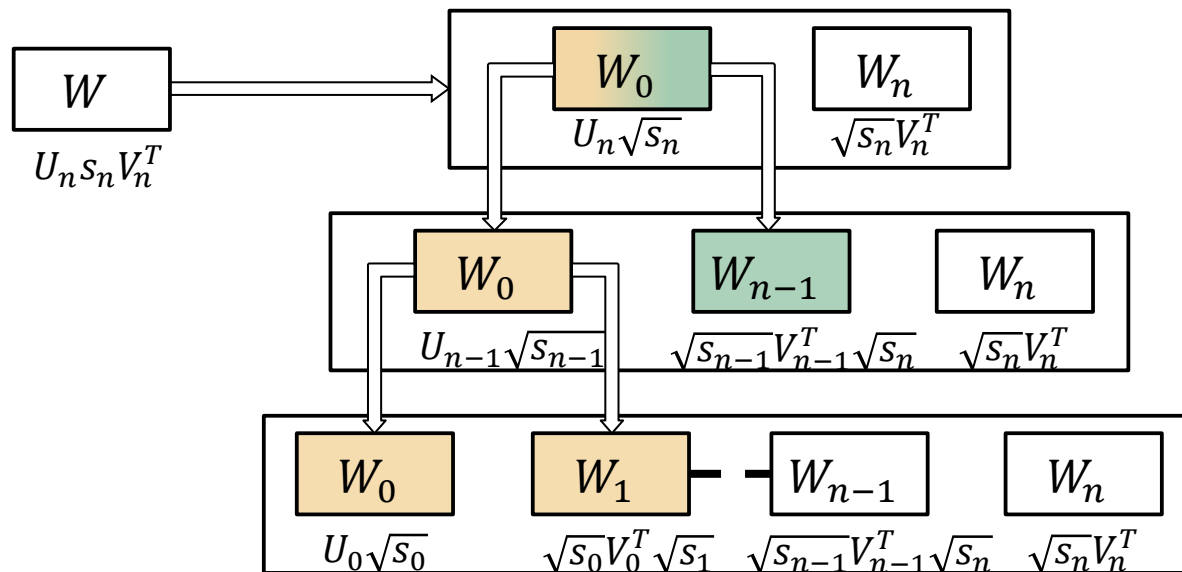
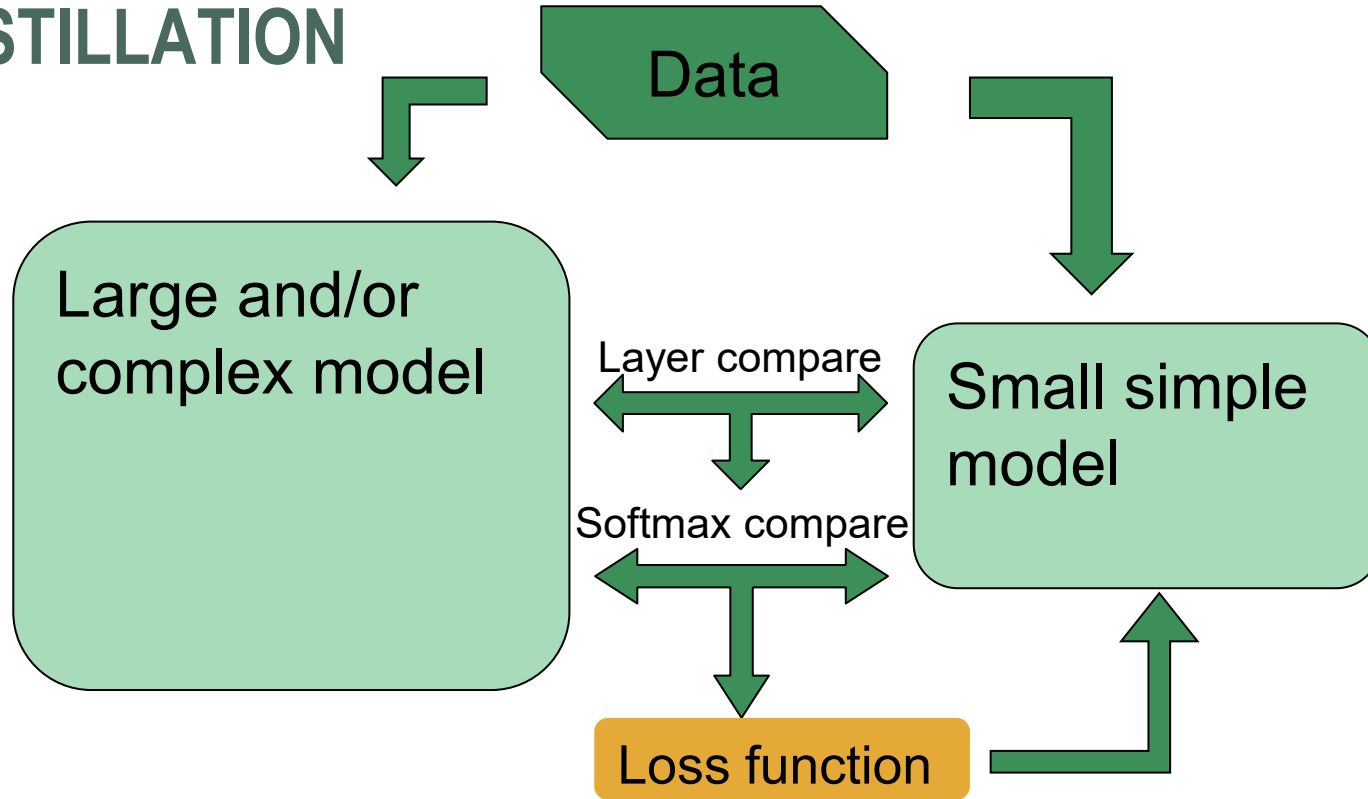
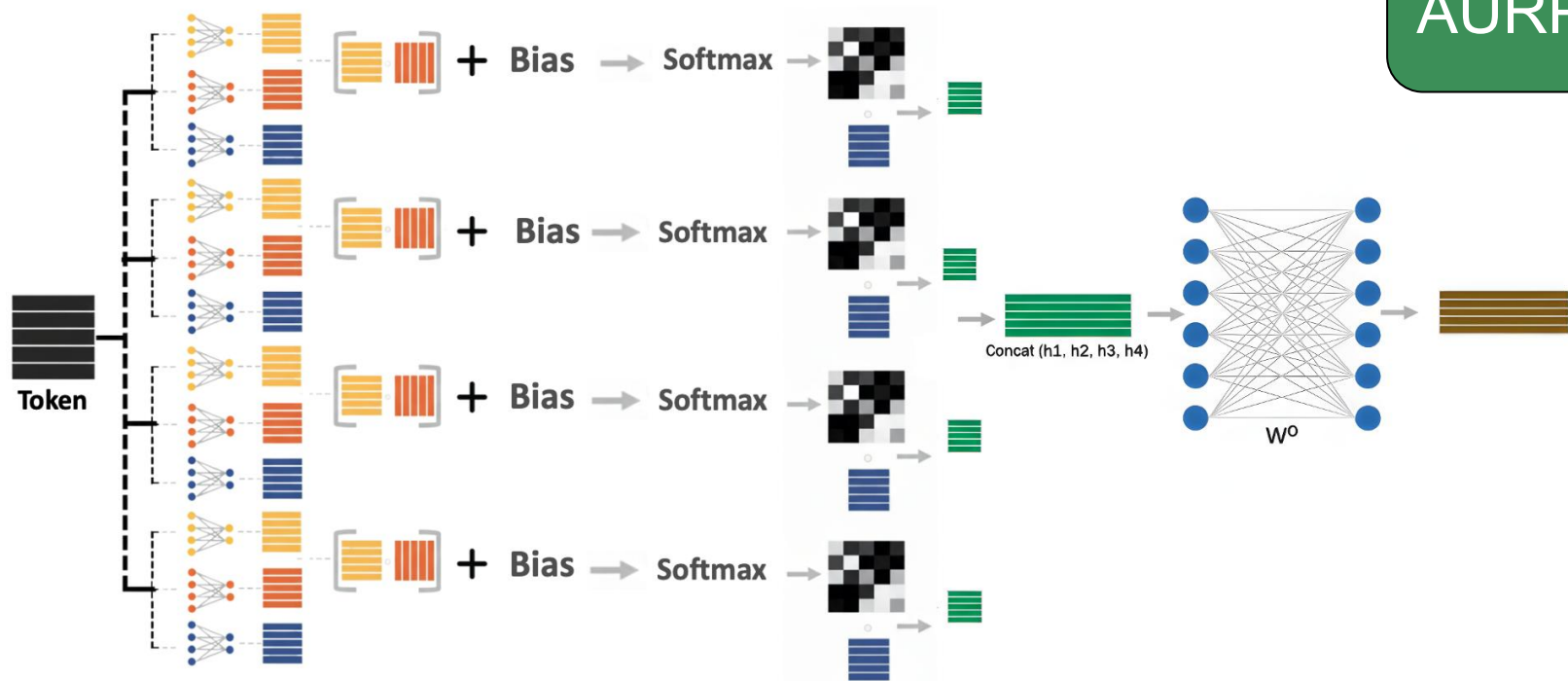


Figure 4: Decomposition of a single network layer into smaller multiple layers.

DISTILLATION



DISTILLATION



0,99
AURPC

DISTILLATION

- Physics informed distillation + Quantization aware training

Performance	0,99 AURPC	No change
Parameters	2300	157x smaller
Memory footprint	2,3 kB	640x smaller
Latency	151 cycles	
Throughput	1,32 Mevents/s	
Usage	<40 %	

SQUISH

- Prune
 - Connections, Nodes, Layers
 - Inputs – Importance analysis
- Input decorrelation
- Feature extraction/signal processing
- Quantization
 - Weights
 - Inputs
- Matrix factorisation
- Distillation

Underlined is
model compression

Univers en soi

UNIVERS EN SOI

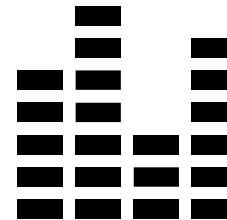
VALIDATION

Univers en soi

UNIVERS EN SOI

DESIGN FOR FAILURE

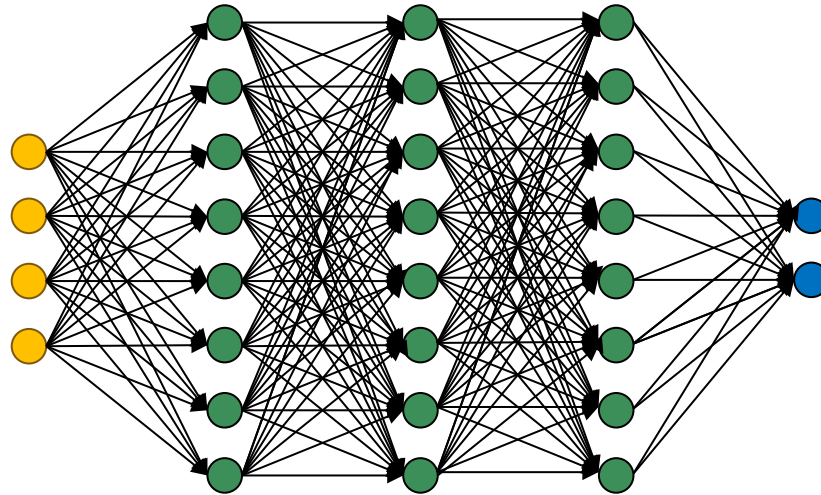
- How should the system fail?
- Prediction uncertainty
- Train failure/unknown case
- Sanity check collection



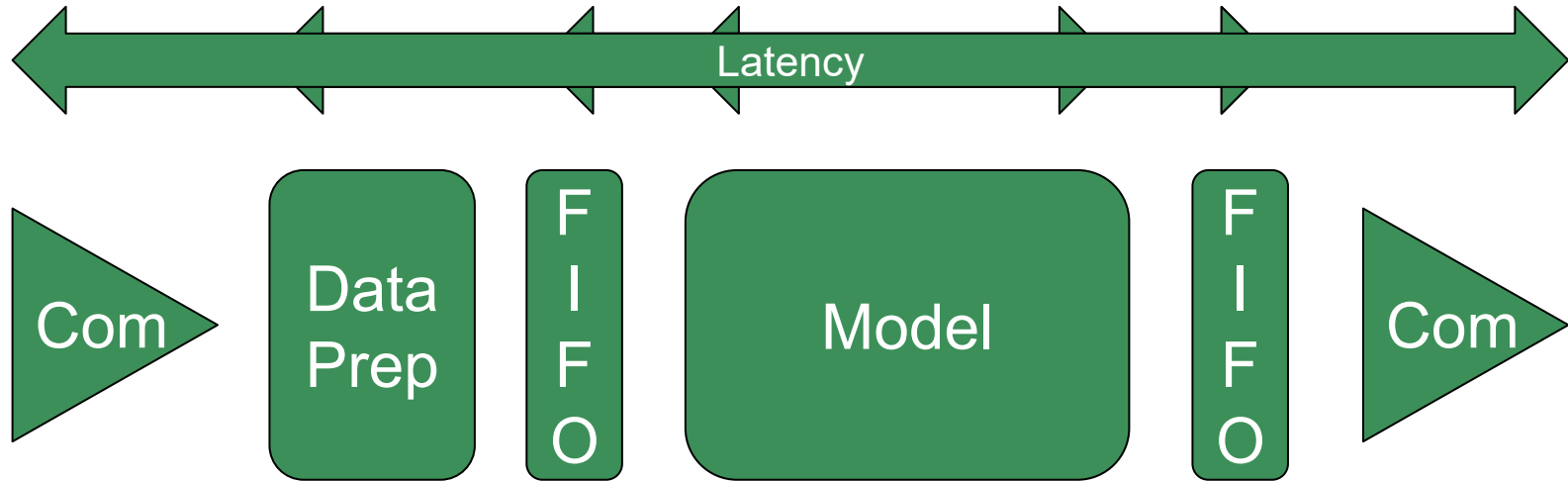
INPUT/OUTPUT VALIDATION

- Design your test set carefully
 - Good inputs, covering the entire application space
 - Extrapolation, outside the scope, but plausible
 - Bad inputs, implausible
 - Noise
 - Detector failure

MODEL STABILITY/SENSITIVITY

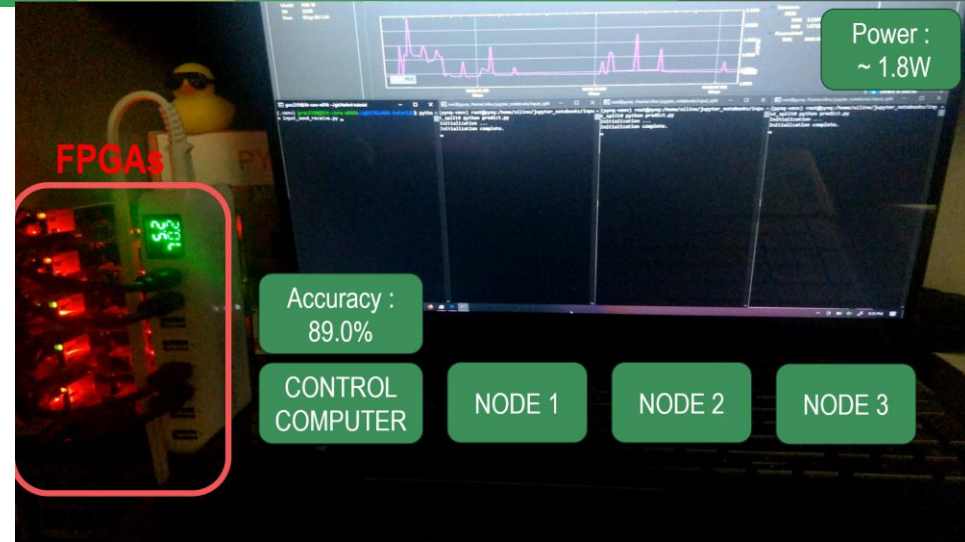


LATENCY



POWER

- Challenging to measure

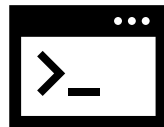
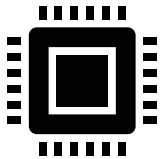
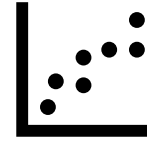
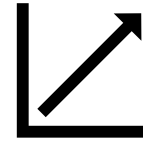
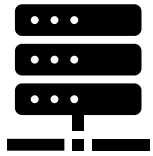


Idle – Full Load = Avg power

$$\frac{\text{Avg power}}{\text{throughput}} = \text{Energy per inference}$$

COMPARISON

- What is the model replacing?



MAINTAINABILITY

- Code readability
- Documentation
- Training data archive
- Update process
 - Functionability update
 - Library compatibility updates
 - Hardware updates

Univers en soi

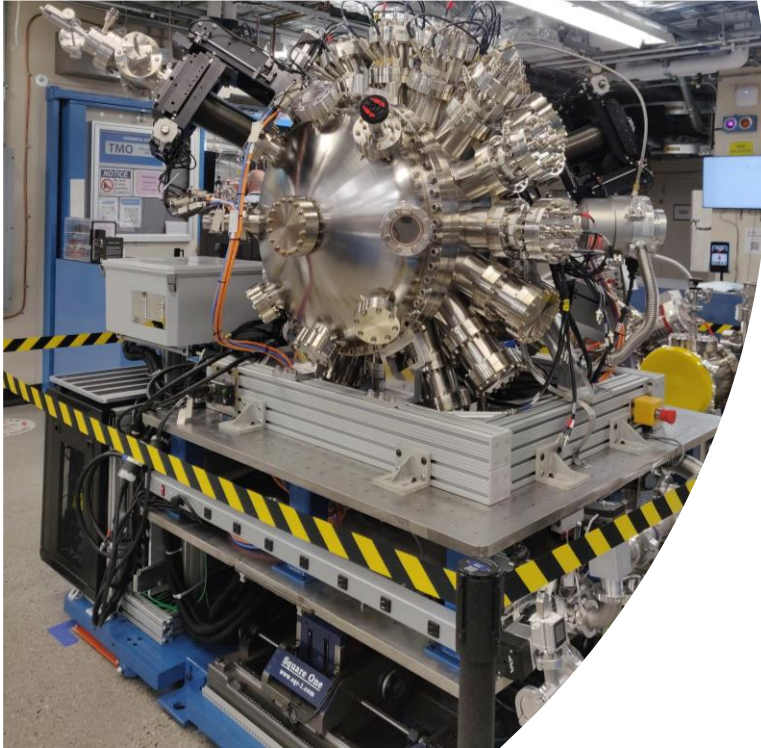
UNIVERS EN SOI

EXAMPLES

Univers en soi

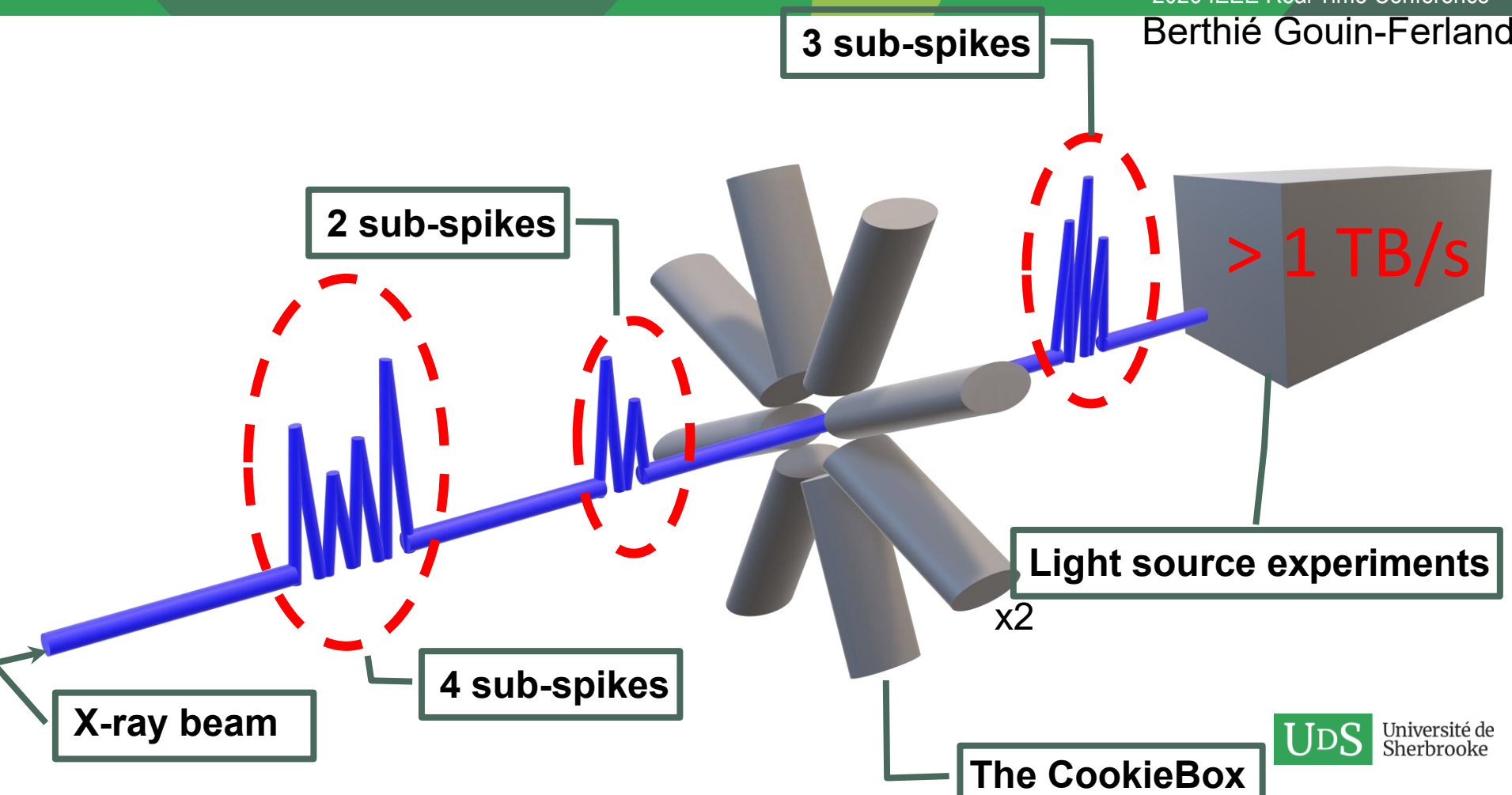
UNIVERS EN SOI

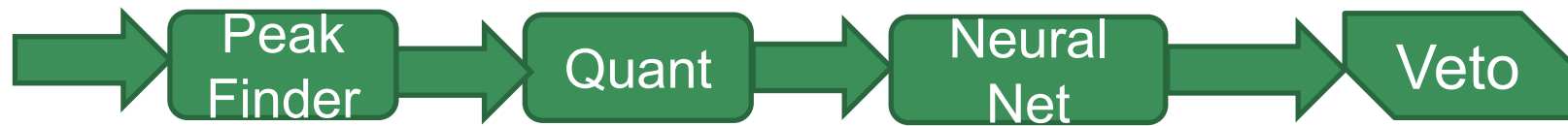
COOKIEBOX



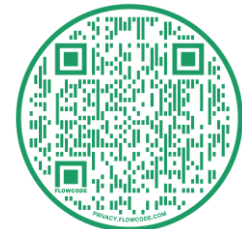
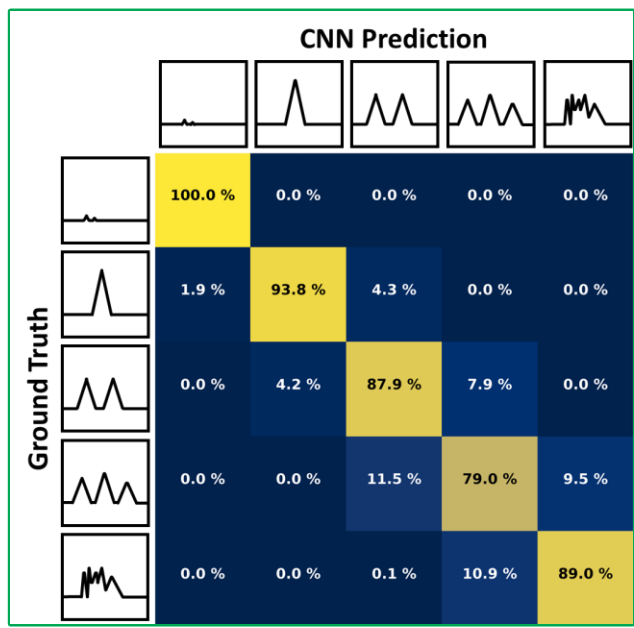
- 16 Time-of-Flight Spectrometers
- 800 GB/s
- 1 MHz rate
- Less than 1 μ s to complete analysis

COOKIEBOX - PROBLEM





Layers
Conv 10x3x3
MaxPool
Conv 20x3x3
MaxPool
Conv 30x3x3
MaxPool
Fully connected 5
Fully connected 5
Softmax classification

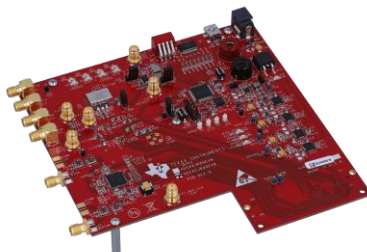


COOKIEBOX – LAB MEASUREMENTS

Berthié Gouin-Ferland Quentin Wingeiring Mehdi Rahimifar



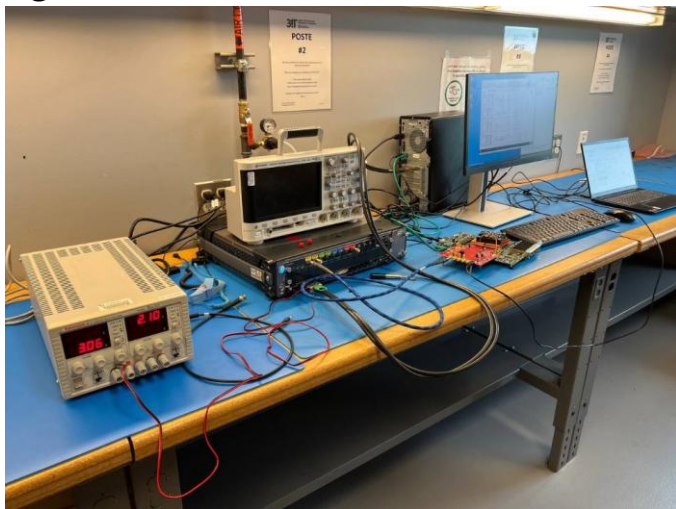
Signal generator



Digitizer



FPGA



Complete Analysis – 400 ns
AI Module – 200 ns



Univers en soi

UNIVERS EN SOI

CONCLUSION

Univers en soi

UNIVERS EN SOI



CONCLUSION

Neural Networks and Decision Trees are compatible with hardware

Co-design is necessary to optimise the system

Squishing the model helps its overall performance

Validation and maintainability are often overlooked!

Edge and distributed ML is a great tool to reduce resource usage

Univers en soi

UNIVERS EN SOI

TOOLS

Univers en soi

UNIVERS EN SOI

SOME TOOLS

- Hls4ml
- Coyote
- SNAC-Pack
- Conifer - <https://github.com/thesps/conifer>
- SNL – to be announced - <https://arxiv.org/abs/2508.21739>
- da4ml - <https://arxiv.org/abs/2507.04535>
- NeuraLUT - <https://arxiv.org/abs/2504.00592>

Univers en soi

UNIVERS EN SOI

THANK YOU!

Univers en soi

UNIVERS EN SOI