# Deep Learning Inference with FPGA

**Corso INFN "Tecniche Di Machine Learning Con Dispositivi FPGA per Gli Esperimenti Di Fisica Delle Particelle"**
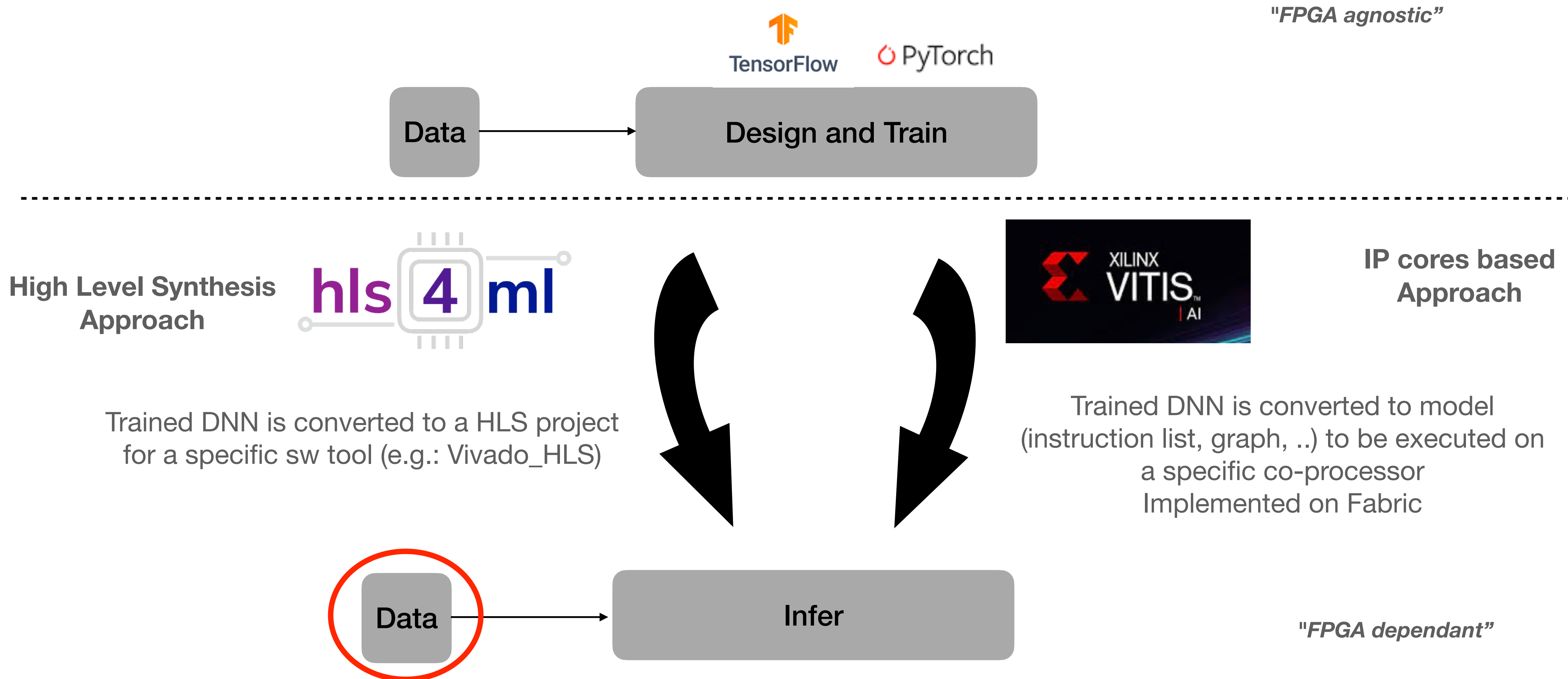
**Riccardo Travaglini - 02/11/2022**

# The ML/AI big picture

- Training and Inference
- Several techniques:
  - BDT, Random forest, DNN, CNN, RNN, …
- Need to accelerate processing and get a lower power consumption

- GPU are best suited for training
- FPGA only used for inference acceleration - best for:
  - Low latency (can be also fixed and reproducible)
  - Low power (not true for floating point arithmetic)
  - Supported architectures (depends on the tool you are using)
    - mainly based on NN: DNN, CNN, RNN, Autoencoders, …  (I can be outdated!)

# From algorithm to FPGA
## Two typical alternative flows

# High Level Synthesis approach

- Given a chosen ML framework: <100% can be supported for HLS translation

- Synthesis and P&R must be performed

  - Usually takes longer on big FPGA needed for AI/ML

- The model (f.i. the CNN) must fit into the FPGA resources (Multi-FPGA requires tools still in development)

- Requires (mostly custom) interface with the FPGA I/O for data

- Weights can be hardcoded into the FPGA

- Model performance can be estimated; they are measurable and reproducible

# IP core based approach

- Given a chosen ML framework: <100% can be supported for HLS translation

- Synthesis and P&R not needed

  - Platform with IP cores already available or can be done once ("overlay")

- Replica of the model can be executed in parallel just like "threads" (depends on the IP core "occupancy")

- Model computed like executing instruction on a co-processor

- Data and weights dynamically loaded to the IP-core (typically IO interface through software), limited by supported formats (e.g. batch limitations for CNN inputs)

- Can be driven by CPU: platform/overlay available

- Performance can be difficult to estimate - profiling tools can be available
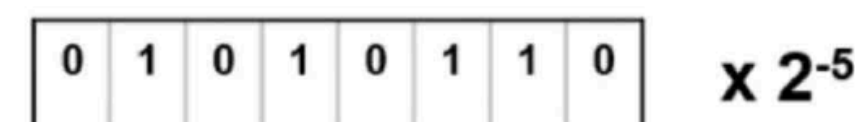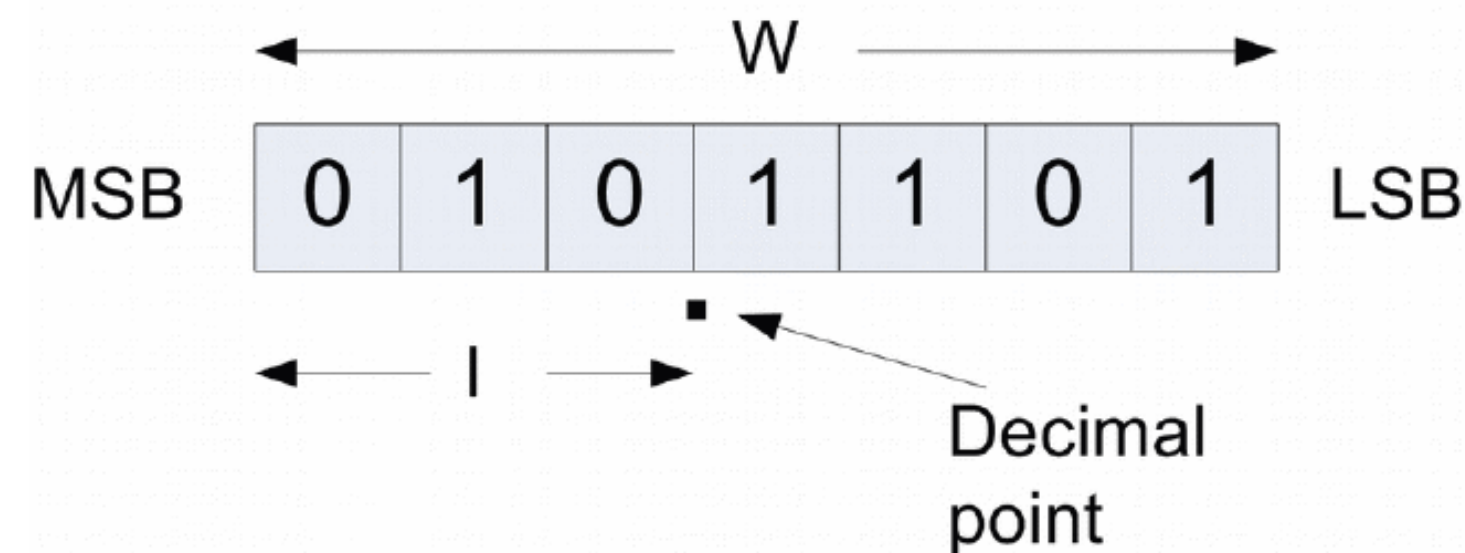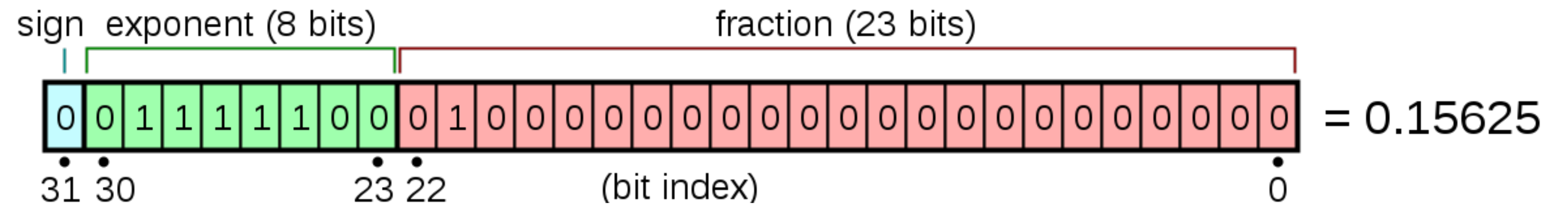
# FPGA Constraints impact
## FPGA have limited resources

- Floating-point (single precision too) adder and multiplier are resource hungry

- Nets with high fan-out have bad impact on timing (slow max. Clock frequency)

# Number representation
## Floating - fixed point - integer - int8

- Positional notation: bit i contributes $2^{\pm I}$

- Signed/unsigned

- Most used in FPGA: arbitrary length fixed point, int8

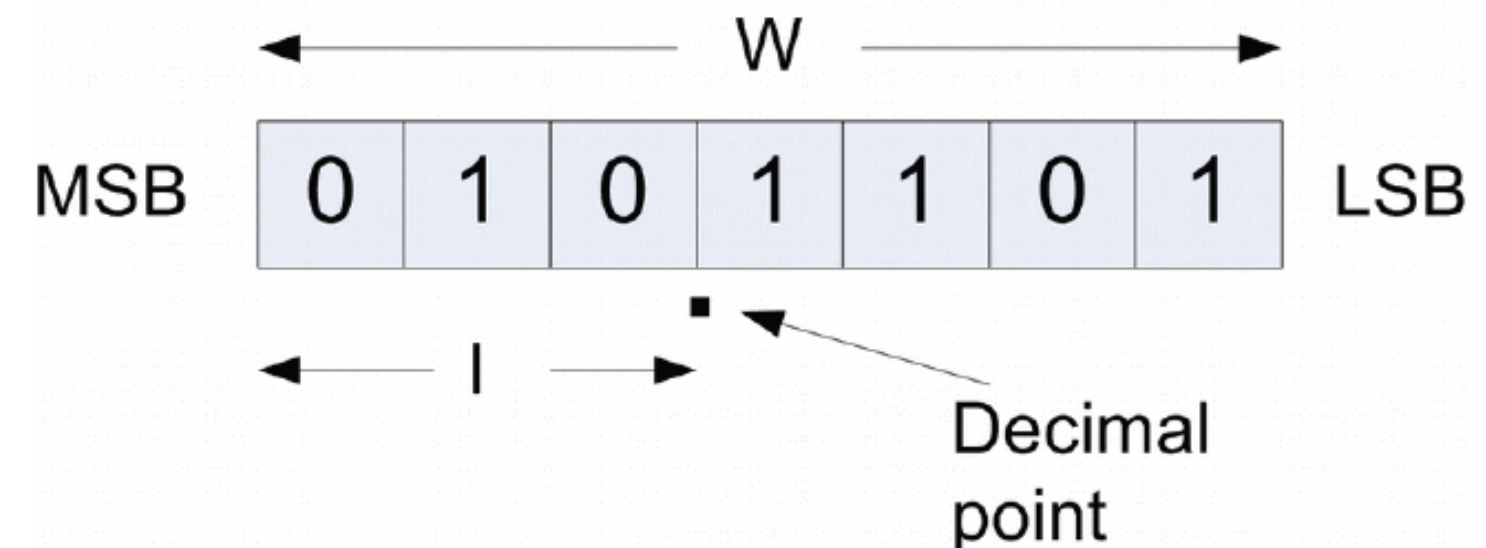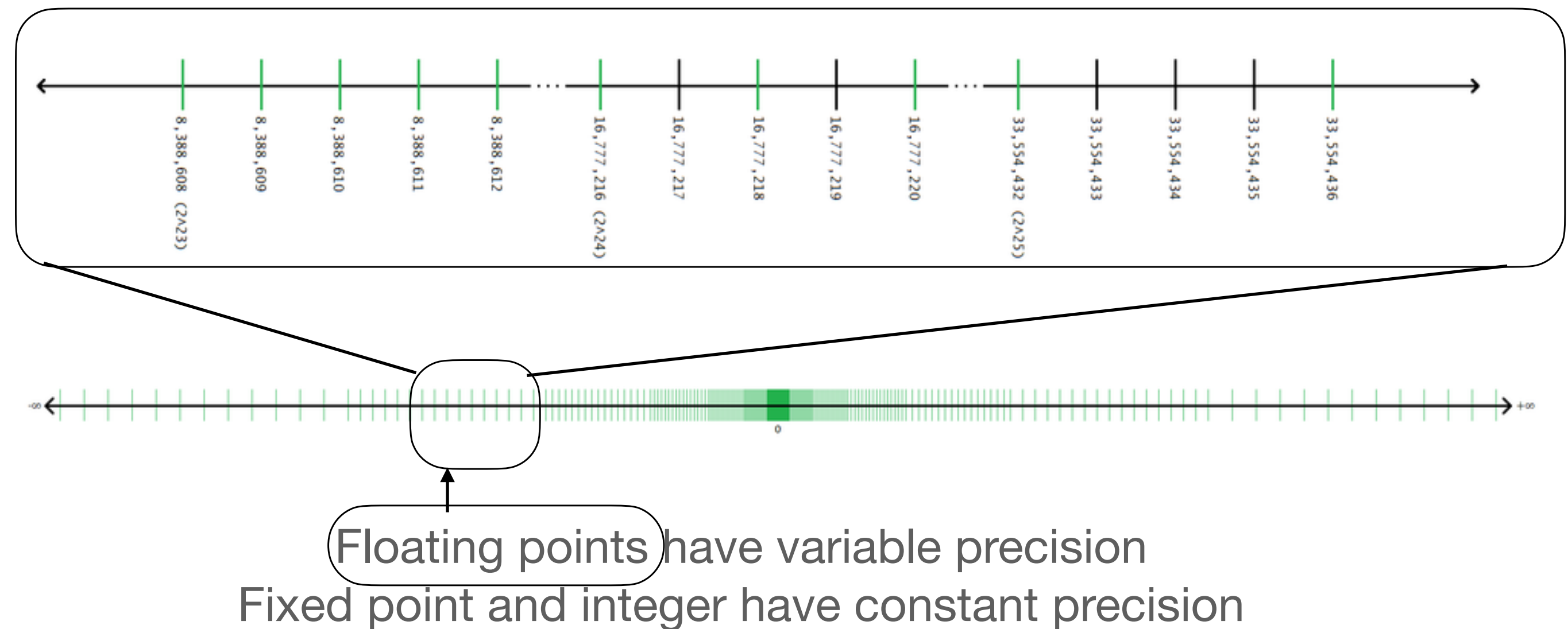- TensorFlow Lite Int8 with exponent (i.e. scale) ; a common offset can be provided



sign  exponent (8 bits)        fraction (23 bits)

| 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |   = 0.15625

31 30                23 22            (bit index)            0

W

MSB | 0 | 1 | 0 | 1 | 1 | 0 | 1 | LSB

I

Decimal point

| a[7] | . | . | . | . | . | . | a[0] |

| Integer | x $2^{Exponent}$

| 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 |  x $2^{-5}$

Exponent is shared (fixed) for a set of variables

# Saturation and precision
## Quantization

- Fixex point with integer I bits is limited to ~ $2^I$

- I-bits unsigned integer limited to $2^I - 1$

- Signed range from ~$-2^{I-1}$ to ~ $2^{I-1}$



8,388,608 (2^23)
8,388,609
8,388,610
8,388,611
8,388,612
16,777,216 (2^24)
16,777,217
16,777,218
16,777,219
16,777,220
33,554,432 (2^25)
33,554,433
33,554,434
33,554,435
33,554,436

Floating points have variable precision
Fixed point and integer have constant precision

W

MSB | 0 | 1 | 0 | 1 | 1 | 0 | 1 | LSB

I

Decimal point

# Quantization

- Design in full precision (typically FP32) e quantize in inference (possibly scaling)

- Quantization-aware training

- Quantization can be global, per layer (NN), per channel (CNN)


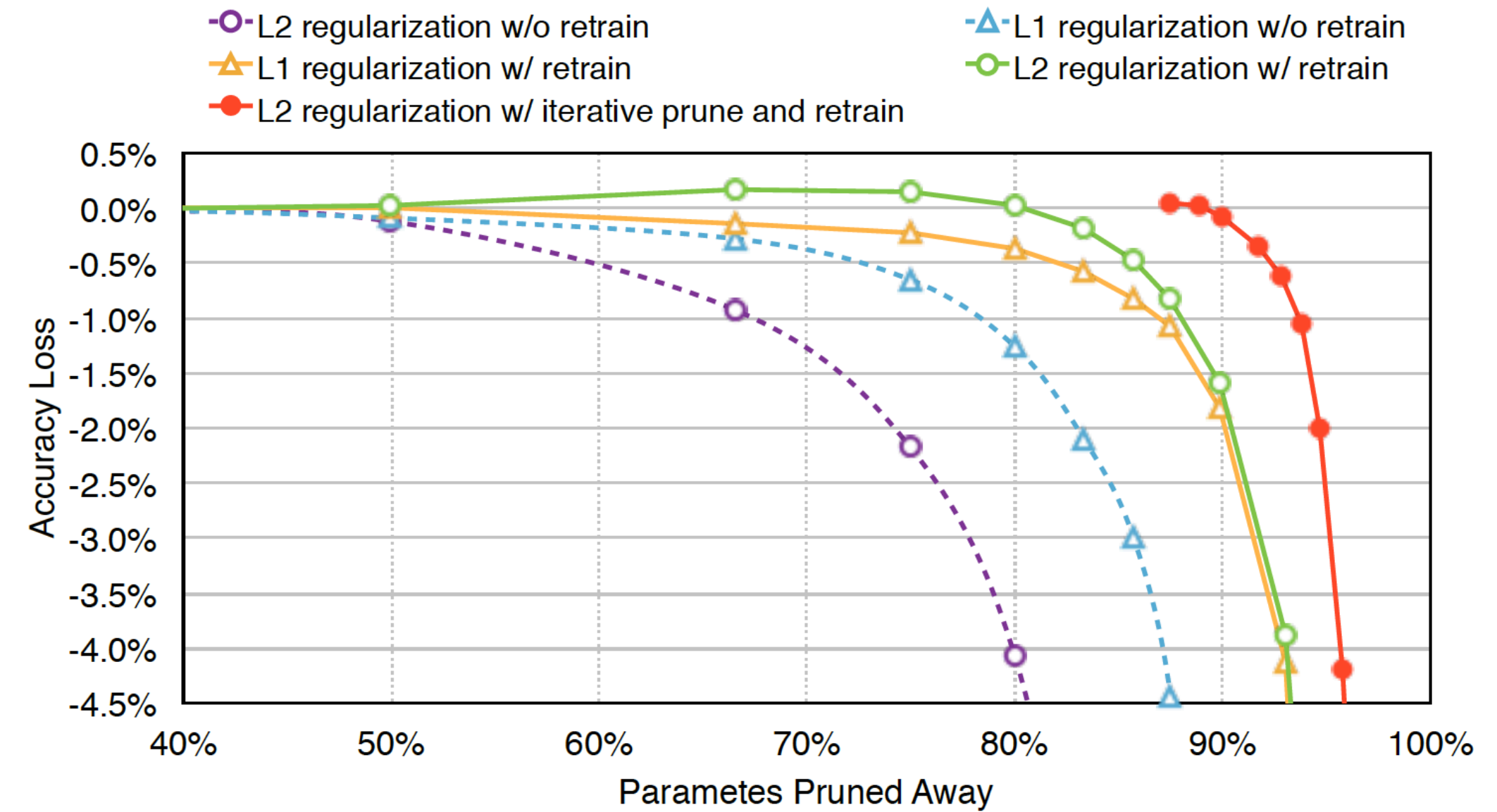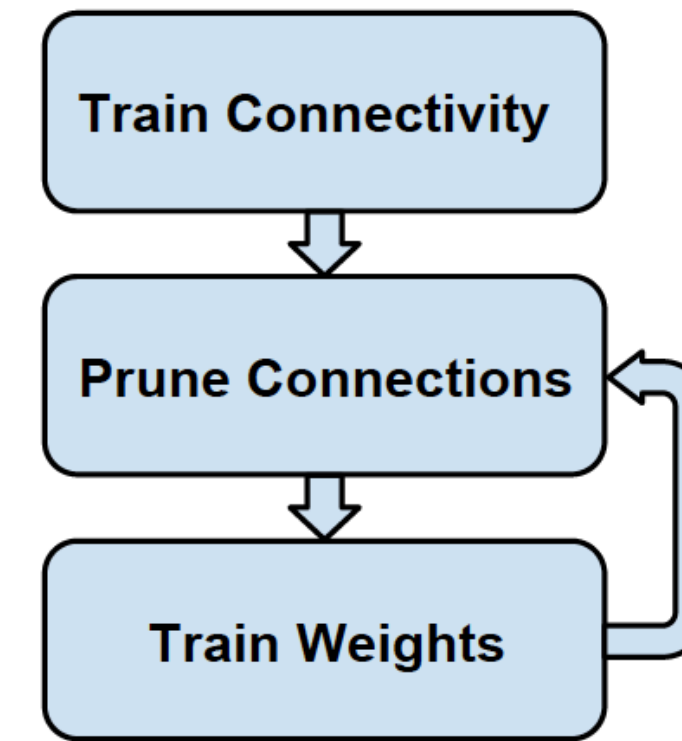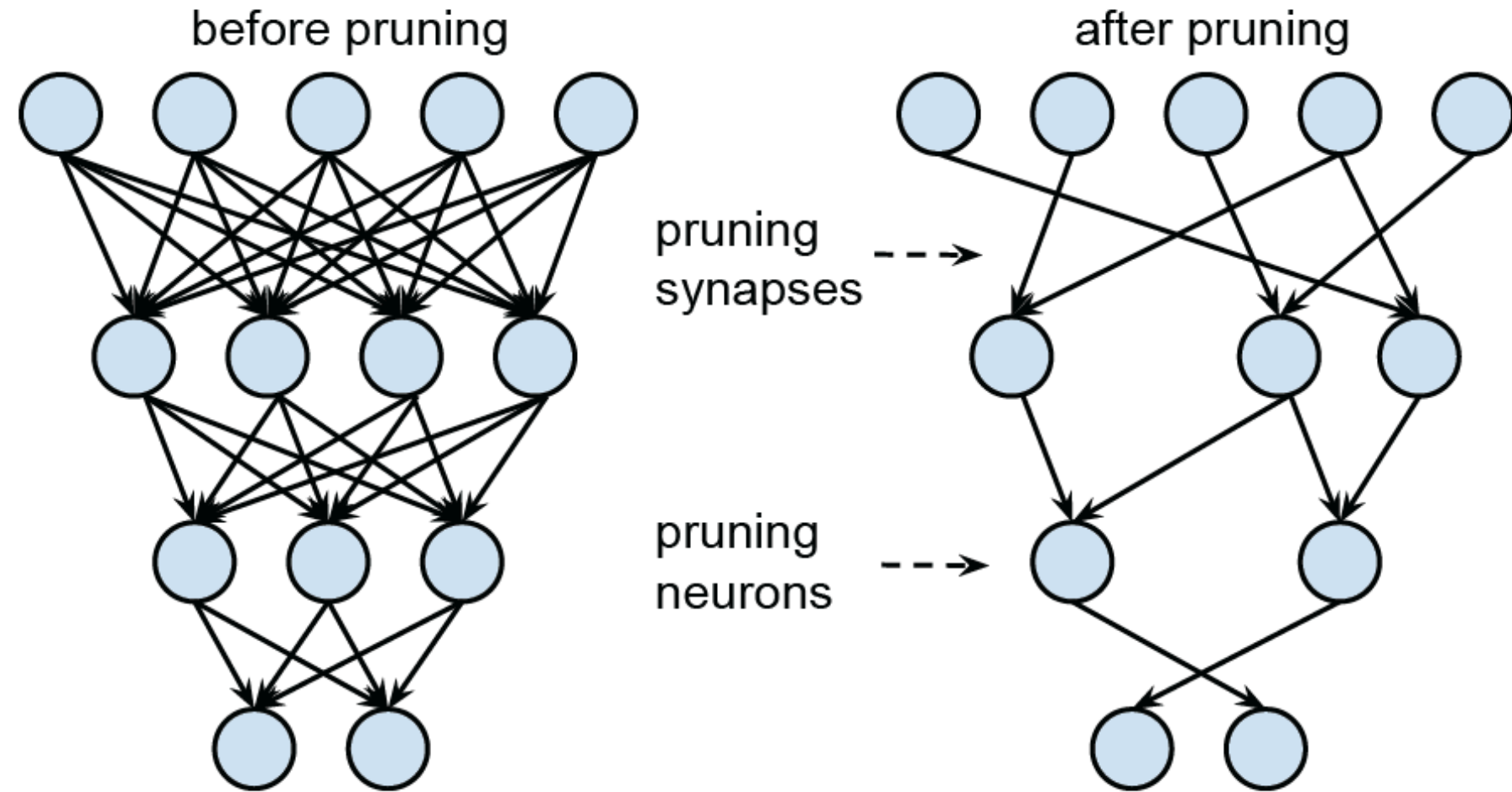- Several tools available

- See hands-on for some examples

# FPGA Constraints impact
## FPGA have limited resources

- Floating-point (single precision too) adder and multiplier are resource hungry

- Nets with high fan-out have bad impact on timing (slow max. Clock frequency)

# Pruning
## Resource optimization



before pruning          after pruning

pruning synapses

pruning neurons

Train Connectivity → Prune Connections → Train Weights

L2 regularization w/o retrain
L1 regularization w/ retrain
L2 regularization w/ iterative prune and retrain
L1 regularization w/o retrain
L2 regularization w/ retrain

Accuracy Loss
Parametes Pruned Away

Learning both Weights and Connections for Efficient Neural Networks
Song Han, Jeff Pool, John Tran, William J. Dally
arXiv:1506.02626 [cs.NE]

# Resource and performance estimation

- HLS-based



- IP based

Resource utilization and clock frequency are fixed by design

Performance difficult to estimate

# Implementation

Data strategy:

- Model standalone

  - How to send/receive data to/from the board?

  - How to interface NN block with the FPGA I/O?

  - I think FPGA expert is mandatory

- Model as a co-processors

  - Bitfile available (either as a dynamically loaded "overlay")

# Running inference for co-processors

- Data can be managed in C/C++, Open-CL or Python with the CPU

- API provided to:

  - Load the bitfile (overlay)

  - Get tensors (size and "reference")

  - Move data through DMA (connecting also to high speed memories into the FPGA boards)

  - Load "n" threads

  - Control execution

  - Access output data