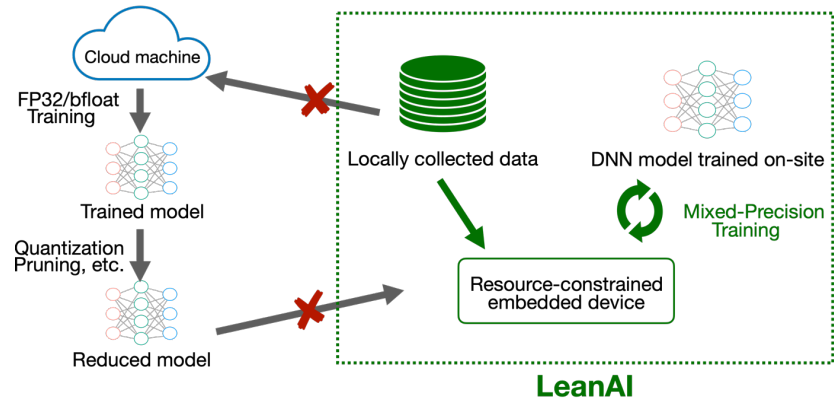
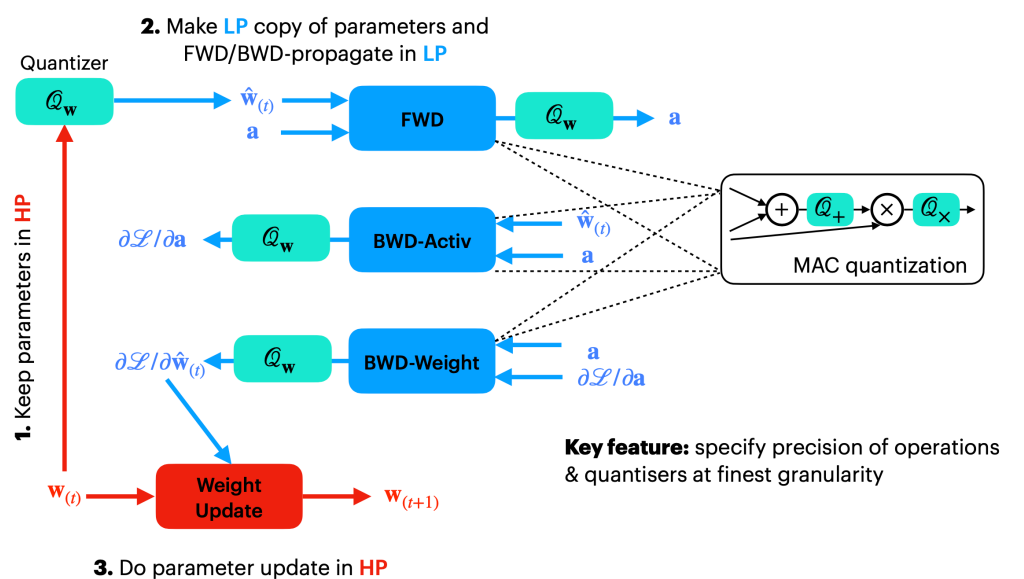


## Context and Objectives



- Need for **learning acceleration** mechanism in both **cloud** (for large-scale models) and **on-site** settings (e.g. autonomous driving, privacy)
- Working on both arithmetic and algorithmic levels
- Design of dedicated HW operators

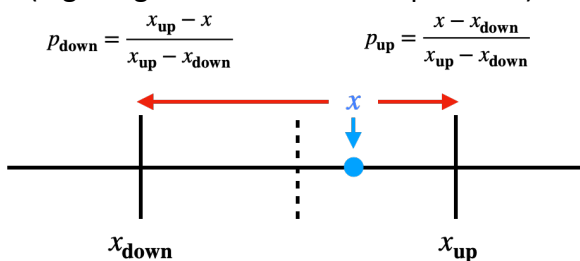
## MPTorch: Mixed-Precision DNN Compute Simulator



Support for CPU & GPU simulation + FPGA-based accelerator prototyping

## Stochastic Rounding for DNN Training Acceleration

Stochastic Rounding (SR) can recapture information that is discarded when bits are rounded off in long computation chains (e.g. long summations or dot products)

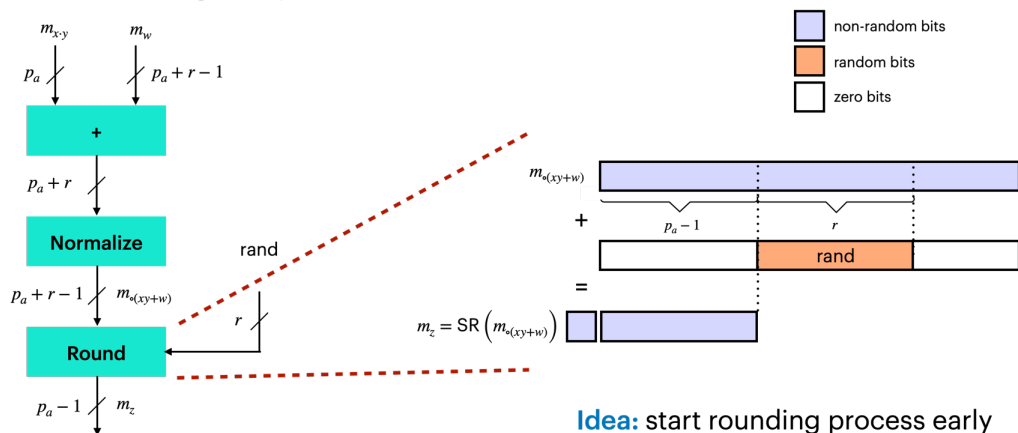


Shown to be beneficial for DNN training acceleration

**Challenge:** not obvious how to optimize in hardware

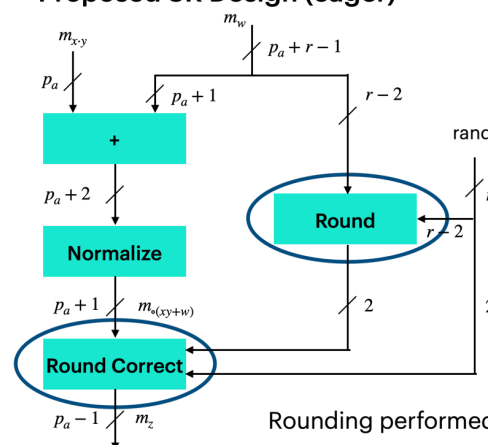
**Basic building block** is the multiply accumulate (MAC) unit:  $z = xy + z$

**Classic SR Design (lazy)**



Idea: start rounding process early

**Proposed SR Design (eager)**



Rounding performed in two stages

**Theoretical result:** probabilistic error analysis on the number of required random bits  $r$  need to implement SR wrt the length  $n$  of the compute chain

$$r \approx \lceil \log_2(n)/2 \rceil$$

## AdaQAT: Adaptive Quantization-Aware Training

Optimization-based method for **mixed-precision** (weights and activations) **DNN quantization**

**Idea:**

$$\mathcal{L}_{total} = \mathcal{L}([N_w], [N_a]) + \lambda \mathcal{L}_{HW}([N_w], [N_a])$$

$$:= [N_w] [N_a]$$

BitOps HW cost estimate

$$\frac{\partial \mathcal{L}}{\partial N_w} \approx \mathcal{L}([N_w], [N_a]) - \mathcal{L}([N_w], [N_a]) \quad \frac{\partial \mathcal{L}}{\partial N_a} \approx \mathcal{L}([N_w], [N_a]) - \mathcal{L}([N_w], [N_a])$$

$$\frac{\partial \mathcal{L}_{total}}{\partial N_w} \approx \frac{\partial \mathcal{L}}{\partial N_w} + \lambda \frac{\partial \mathcal{L}_{HW}}{\partial [N_w]}$$

$$\frac{\partial \mathcal{L}_{total}}{\partial N_a} \approx \frac{\partial \mathcal{L}}{\partial N_a} + \lambda \frac{\partial \mathcal{L}_{HW}}{\partial [N_a]}$$

Model	Method	# exploration epochs	# total epochs	Bit-width (W/A)	Accuracy (%) top-1	BitOPs (Gb)
ResNet-18	DQ [Uhl+19]	50	50	5.11/10.4	70.1	93.6
	FracBits [YJ21b]	120	50	4.00/4.00	70.6	34.7
	SDQ [Hua+22]	60	150	3.85/4	71.7	33.4
	<b>Our</b>	<1	100	<b>3.84/4.00</b>	71.4	<b>31.4</b>
MobileNet-V2	DQ [Uhl+19]	50	50	5.77/-	69.7	93.6
	FracBits [YJ21b]	120	150	4.00/4.00	71.3	5.35
	SDQ [Hua+22]	60	180	3.79/4	72.0	5.07
	<b>Our</b>	<1	150	<b>3.86/3.88</b>	71.3	<b>4.95</b>

Publications:

