



**Problem:** Linear model training using low-precision SGD.

**Existing Approach:**

- 1, One hardware design for each precision.
- 2, One quantized dataset for each precision.

**Our Approach (MLWwaving):**

One hardware design and one memory layout for any precision.

$\vec{x}$ : model,  
 $\vec{g}$ : gradient,  
 $\Gamma$ : learning rate,  
 $df$ : derivative of loss function

```

epochs*/
For i = 1 to N do /*N
samples ( $\vec{a}_i, b_i$ )*/*
ax =  $Q_s(\vec{a}_i) * \vec{x}$ ; /*dot
product*/
scale =  $\gamma * df(ax, b_i)$ ; /*serial part*/
 $\vec{g} = scale * Q_s(\vec{a}_i)$ ; /*gradient
comp*/
 $\vec{x} = \vec{x} - \vec{g}$ ; /*model
update*/

```

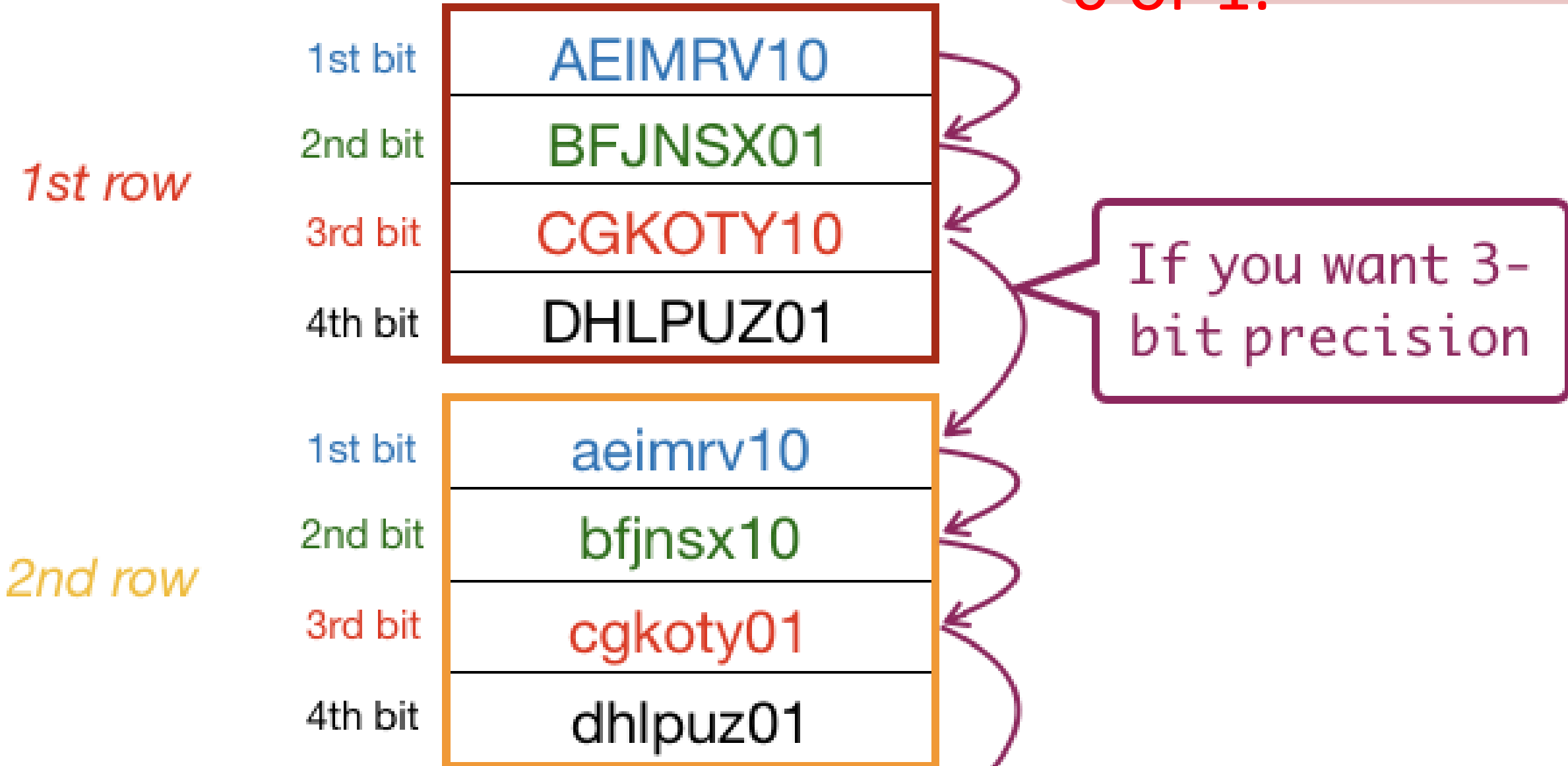
## MLWeaving: Software/Hardware Co-design

### MLWeaving memory layout (software)

1, Original full-precision fixed-point table for training dataset a

1st row	ABCD	EFGH	IJKL	MNOP	RSTU	VXYZ	1010	0101
2nd row	abcd	efgh	ijkl	mnop	rstu	vxyz	1100	0011

A~Z (a~z) is binary, 0 or 1.

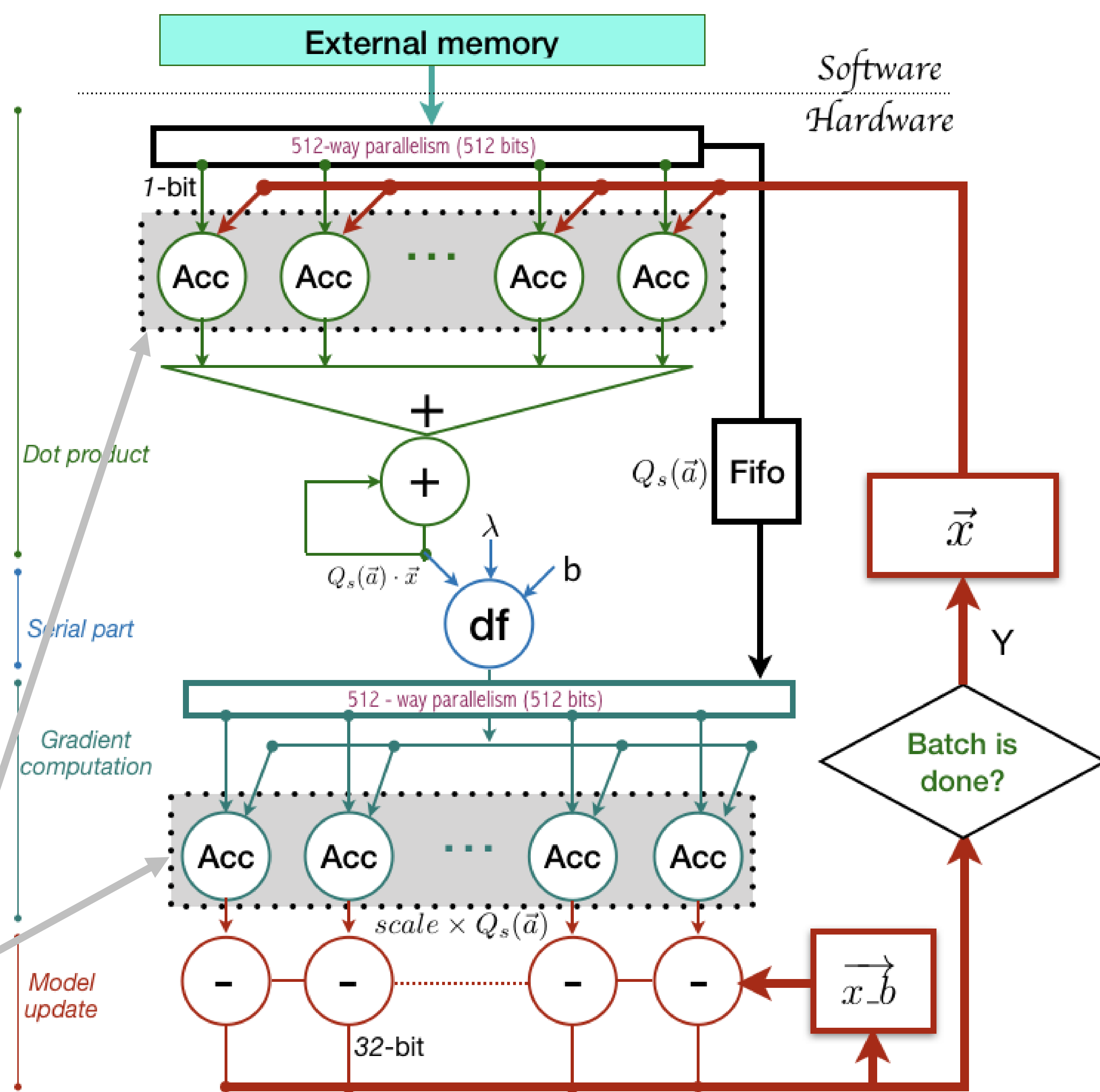


2, MLWeaving memory layout for training dataset a

Bit-level flexibility of precision from software side

Directly consume the data from MLWeaving memory layout with bit-serial multipliers.

### MLWeaving arithmetic (hardware)

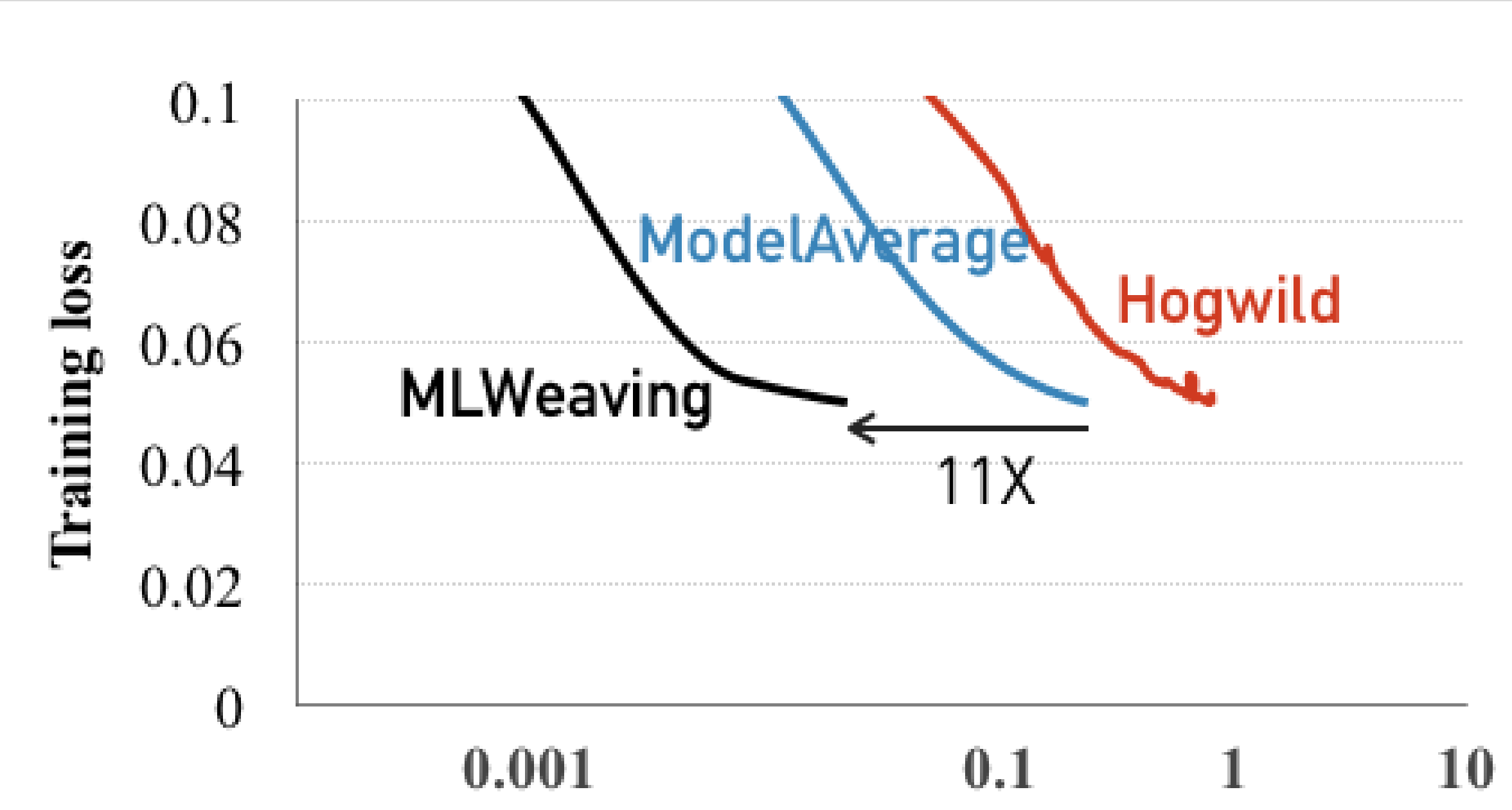
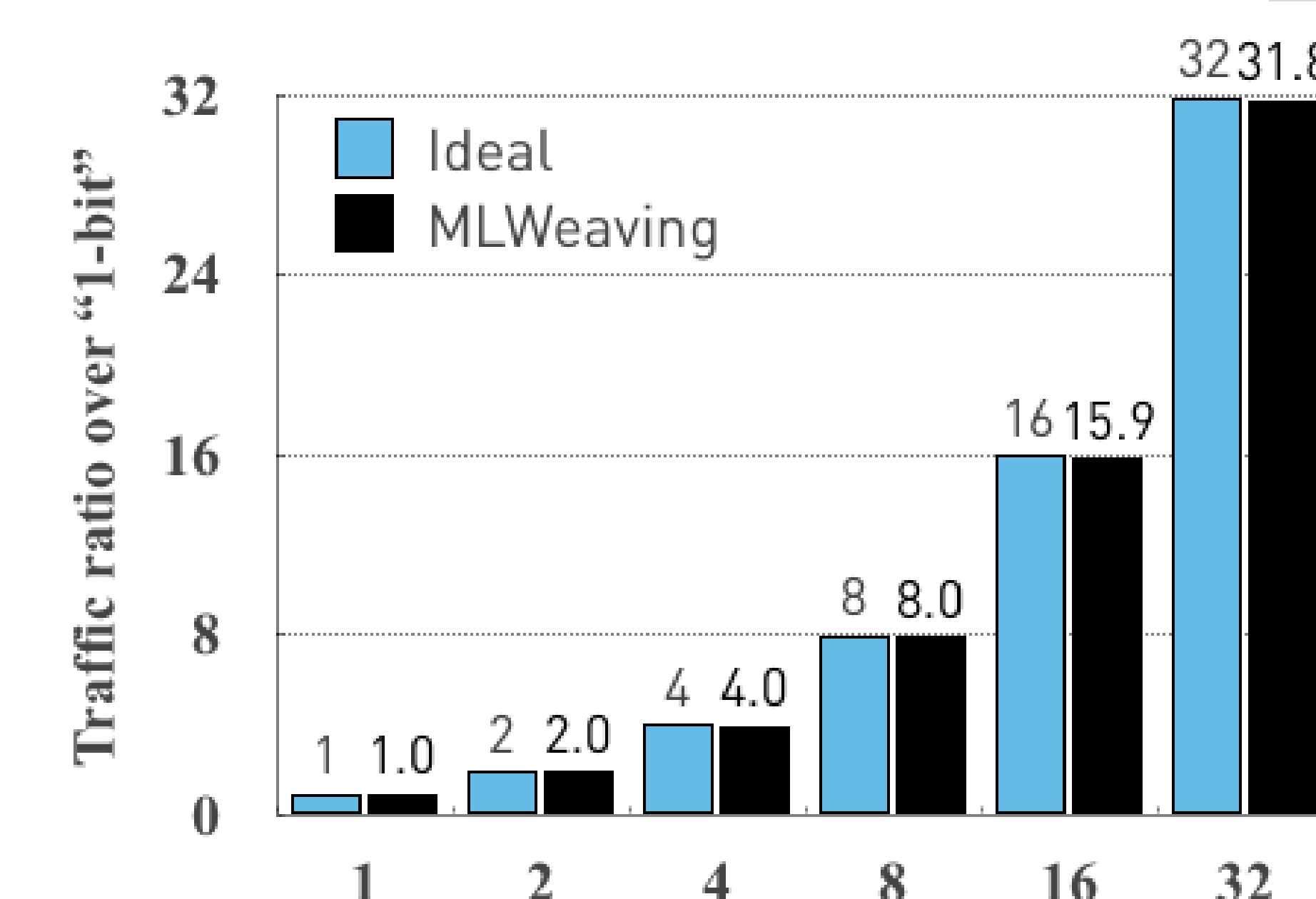
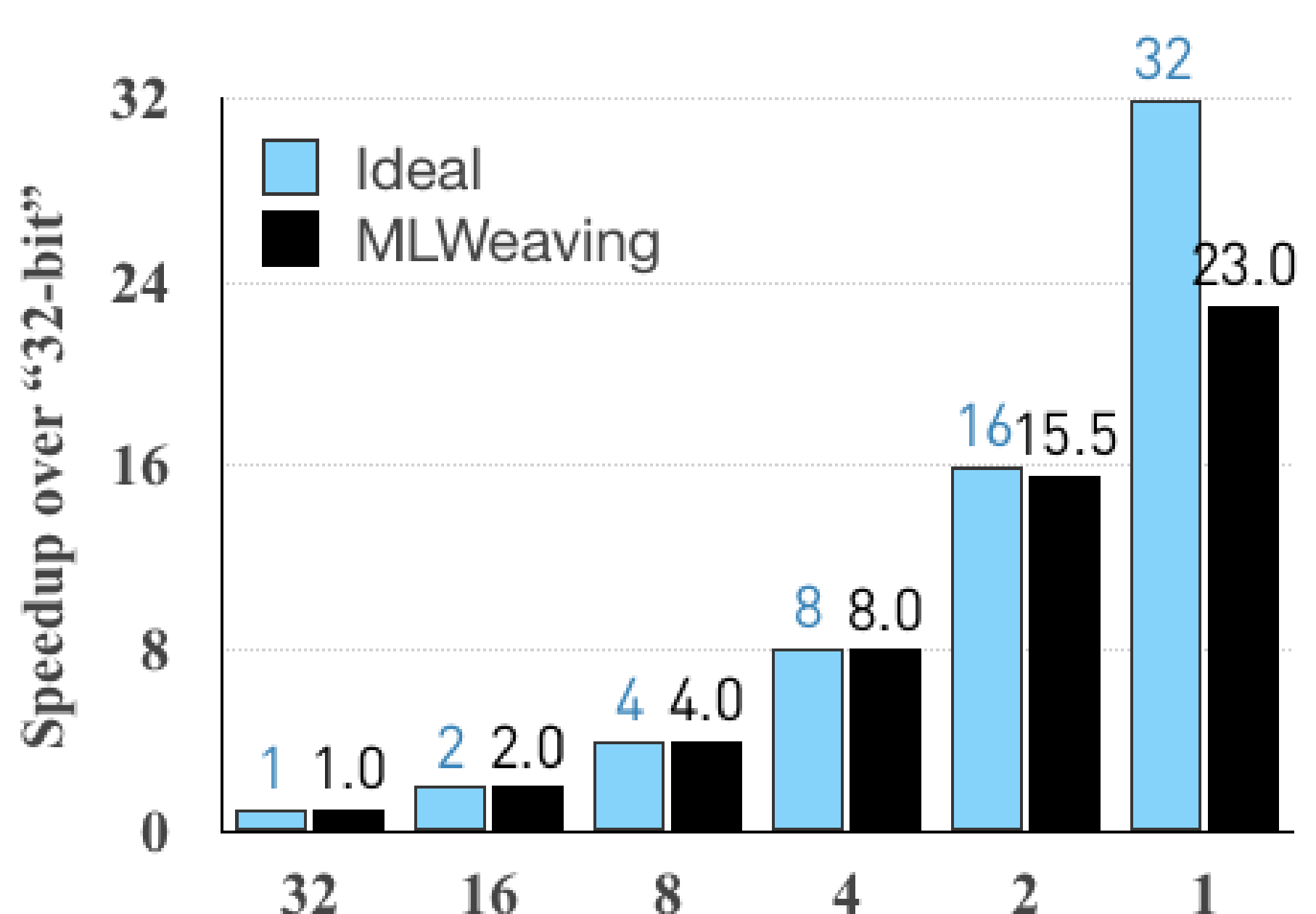


## Experiment

- Hardware: an Intel Broadwell CPU (14 cores, 35MB LLC, 60 GB/s memory bandwidth) and an Intel Arria 10 FPGA (directly access the CPU memory via one QPI and two PCIe, with memory bandwidth: 15GB/s).
- Dataset: Epsilon (40,000 samples, 2000 features).
- Hogwild (ModelAverage): state-of-the-art parallel implementations of SGD on CPUs, using 14 cores, AVX2 and 8-bit dataset.

## Findings:

- 1, MLWeaving can roughly achieve linear speedup (time or memory traffic), when a lower number of bits is used, as shown in Figures a, b.
- 2, MLWeaving on an FPGA can achieve 11X speedup over its CPU rivals in Figure c.



(a) Time vs. Precision

(b) Memory traffic vs. Precision

(c) MLWeaving vs. CPU rivals