# MSD: Mixing Signed Digit Representations for Hardware-efficient DNN Acceleration on FPGA with Heterogeneous Resources

**Jiajun Wu**, Jiajun Zhou, Yizhao Gao, Yuhao Ding,
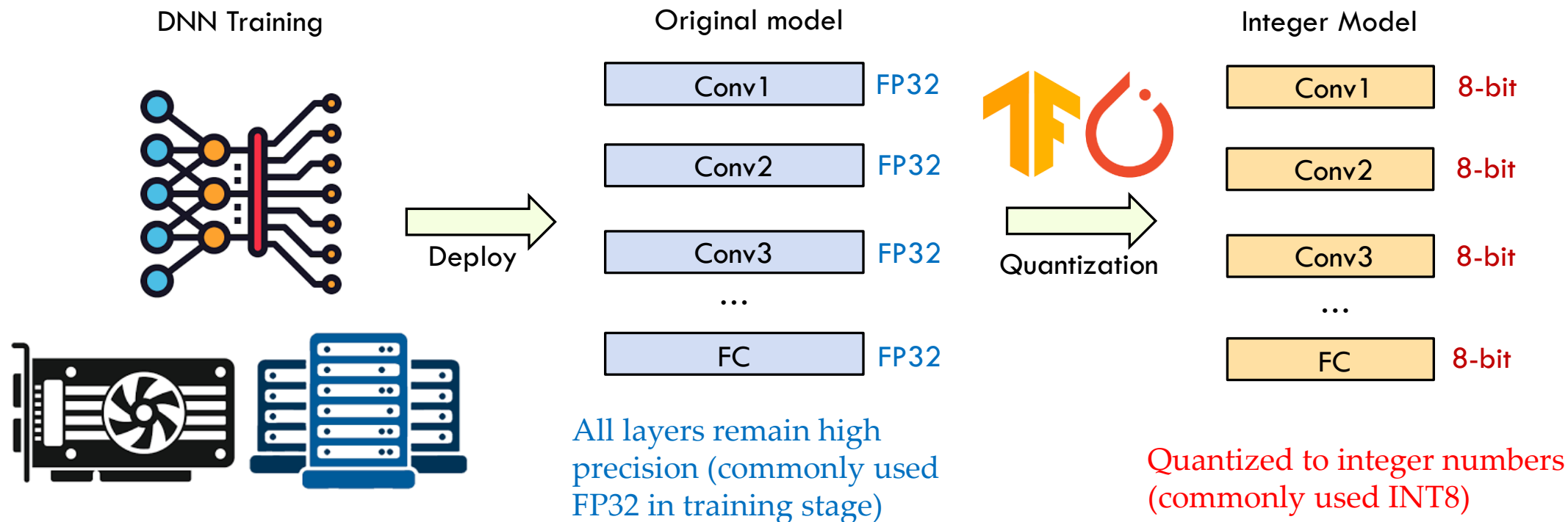Ngai Wong, Hayden Kwok-Hay So

University of Hong Kong

*FCCM 2023*

9 May 2023

# Contents

- **Motivation & Background**

- The Proposed Methodology

- Experiment Results

- Conclusion & Future Works

# DNN Quantization



DNN Training

Original model

Integer Model

Conv1 — FP32
Conv2 — FP32
Conv3 — FP32
...
FC — FP32

Deploy

Quantization

Conv1 — 8-bit
Conv2 — 8-bit
Conv3 — 8-bit
...
FC — 8-bit

All layers remain high precision (commonly used FP32 in training stage)

Quantized to integer numbers (commonly used INT8)

Can we further compress the computation of int8 to improve the inference performance without loss of accuracy?
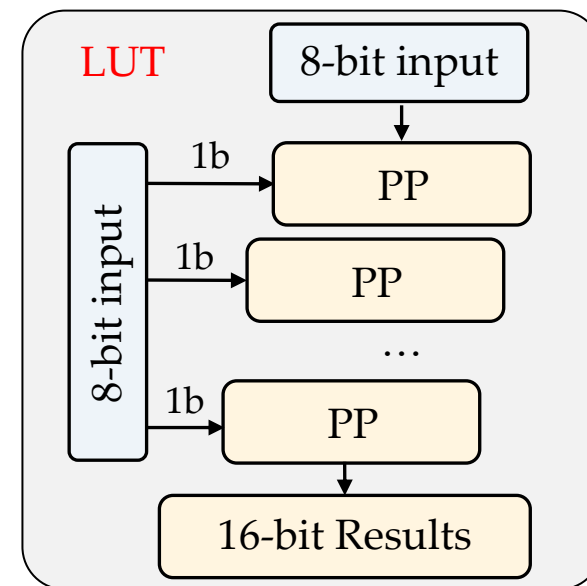
# Deploy Quantized Multiplication

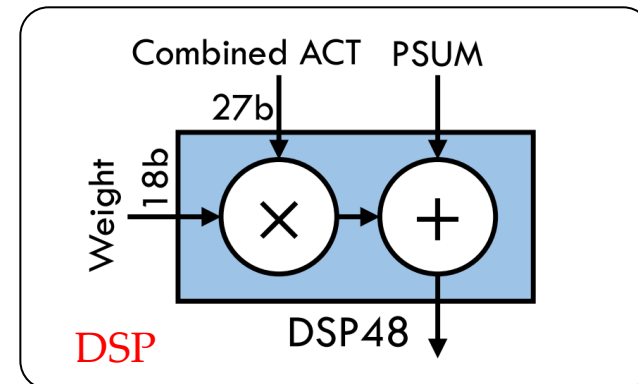Conventional Int8 multiplication

$$14 \times 27 = 00001110 \times 00011011$$

$$
\begin{array}{r}
00001110 \\
\times\ 00011011 \\
\hline
00001110 \\
00001110 \\
00000000 \\
00001110 \\
00001110 \\
00000000 \\
00000000 \\
+\ 00000000 \\
\hline
\end{array}
$$

Partial Products (pp)

16-bit results

**FPGA Implementation**

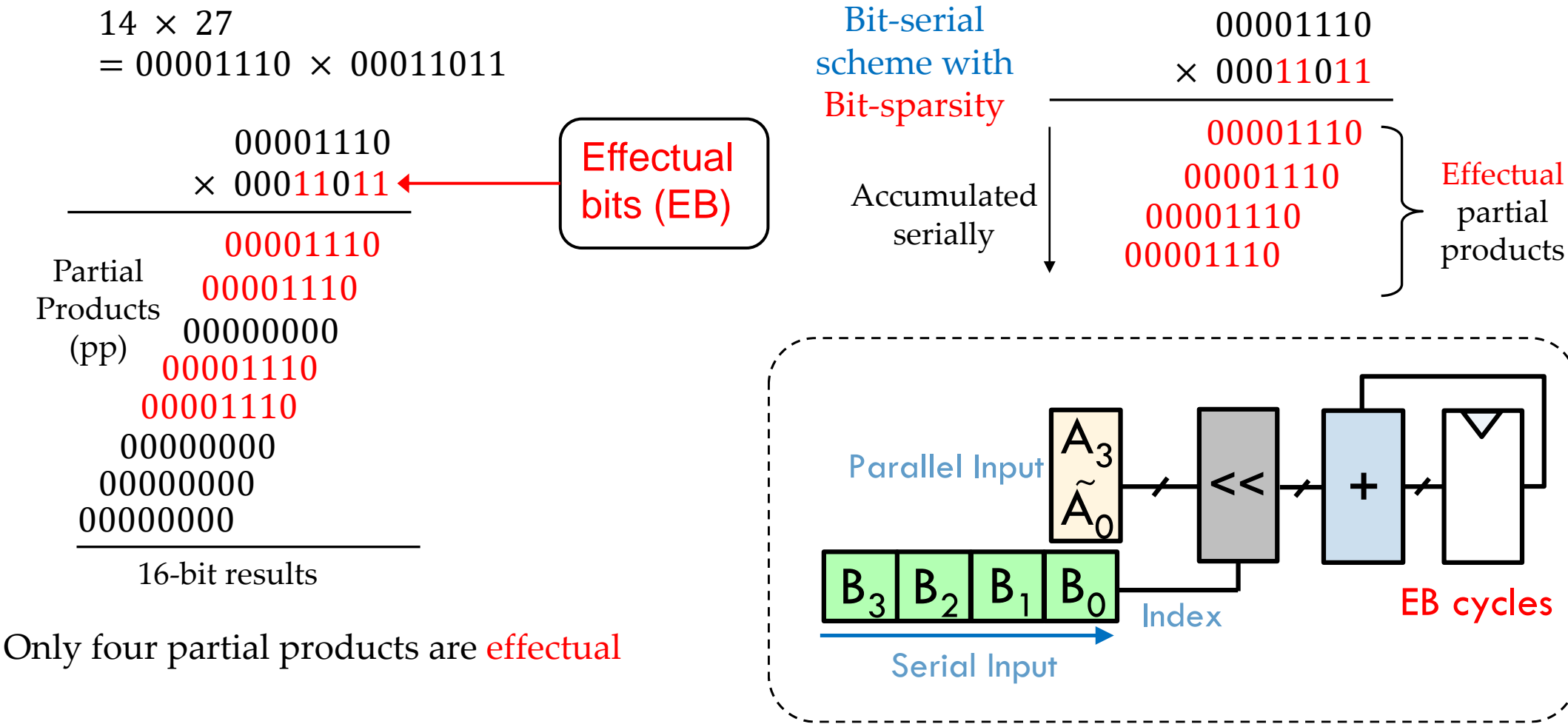**How do we further improve the performance?**

1. To use both DSPs and LUTs for multiplication.
2. To make our LUT implementation smaller.

**Parallel Multiplier**
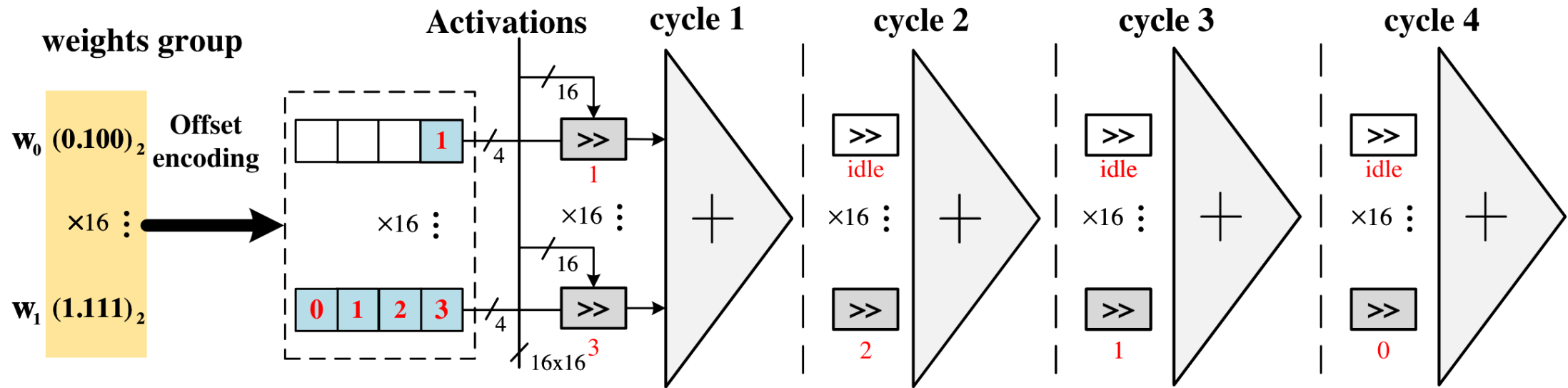


Combined ACT    PSUM
27b
Weight  18b
DSP    DSP48



LUT
8-bit input
8-bit input
1b    PP
1b    PP
...
1b    PP
16-bit Results

Latency: 1 cycle

# Bit-Serial Scheme with Bit-Sparsity

**An alternative approach for the multiplier on LUT**

$14 \times 27$
$= 00001110 \times 00011011$

$$
\begin{array}{r}
00001110 \\
\times \; 000\textcolor{red}{11}0\textcolor{red}{11} \\
\hline
\textcolor{red}{00001110} \\
\textcolor{red}{00001110} \\
00000000 \\
\textcolor{red}{00001110} \\
\textcolor{red}{00001110} \\
00000000 \\
00000000 \\
00000000 \\
\hline
\end{array}
$$

Partial Products (pp)

16-bit results

Only four partial products are effectual

Effectual bits (EB)

Bit-serial scheme with Bit-sparsity

$$
\begin{array}{r}
00001110 \\
\times \; 000\textcolor{red}{11}0\textcolor{red}{11} \\
\hline
\textcolor{red}{00001110} \\
\textcolor{red}{00001110} \\
\textcolor{red}{00001110} \\
\textcolor{red}{00001110} \\
\end{array}
$$

Accumulated serially

Effectual partial products



Parallel Input $A_3 \; \tilde{A}_0$

$B_3 \; B_2 \; B_1 \; B_0$

Index

Serial Input

EB cycles

# Workload Imbalance in Bit-Sparsity Scheme

The workload imbalance issue in bit-sparsity:



Ref: BitCluster (TCAD, Volume: 41, Issue: 11, November 2022)

Need a method to use a restricted and small number of EB without causing major loss of accuracy.

# Contents

- Motivation & Background

- **The Proposed Methodology**
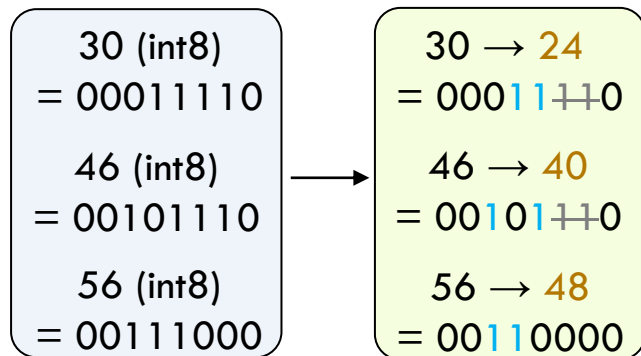
- Experiment Results

- Conclusion & Future Works

# Mixing Signed-digit Representations (MSD)

- A framework that automatically <span style="color:red">partitions</span> the DNN workloads to run on <span style="color:red">DSPs and LUTs</span>

- We proposed a new representation called <span style="color:red">restricted signed-digit (RSD)</span> to help restrict the number of EB

# Signed-Digit Representation

To solve the imbalance issue

30 (int8)
= 00011110

46 (int8)
= 00101110

56 (int8)
= 00111000

$\longrightarrow$

$30 \rightarrow 24$
= 00011$\overline{11}$0

$46 \rightarrow 40$
= 00101$\overline{11}$0

$56 \rightarrow 48$
= 00110000

Restriction: EB = 2

Large quantization error ☹

Find another number format with **higher representation capability**

Our Approach: Signed-digit representation

Let $X[i]$ (i-th bit in the number with n-bit) expand to 0, 1 and -1 ($\overline{1}$):

$$X = \sum_{i=0}^{n-1}\left(X[i] \times 2^i\right), \qquad X[i] = \{0, 1, -1\}$$

Effectual bits (EB)

Note that 2's complement is actually a special case of signed-digit:

$$X = (-1)^{X[n-1]} \times 2^{n-1} + \sum_{i=0}^{n-2}\left(X[i] \times 2^i\right), \qquad X[i] = \{0, 1\}$$

$$= \sum_{i=0}^{n-1}\left(X[i] \times 2^i\right), \qquad X[i] = \begin{cases} \{0, 1\} & if\ i \neq n-1 \\ \{0, -1\} & if\ i = n-1 \end{cases}$$
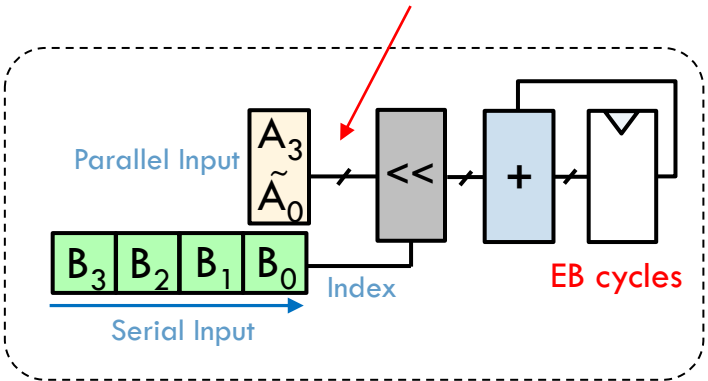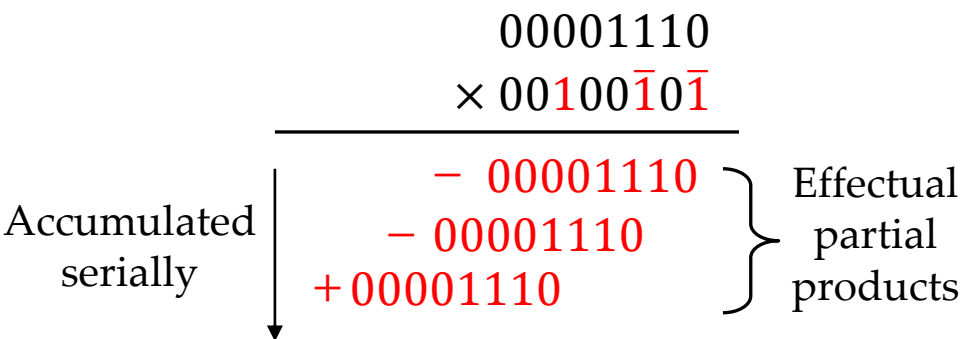
# Hardware Does Not Change a lot!

$$14 \times 27$$
$$= 00001110 \times 00011011$$

$$\begin{array}{r} 00001110 \\ \times\ 00011011 \end{array}$$

Accumulated serially

$$\begin{array}{r} 00001110 \\ 00001110 \\ 00001110 \\ 00001110 \end{array}$$

Effectual partial products

**Signed-digit scheme**

$$14 \times 27$$
$$= 00001110 \times 00100\bar{1}0\bar{1}$$

$$\begin{array}{r} 00001110 \\ \times\ 00100\bar{1}0\bar{1} \end{array}$$

Accumulated serially

$$\begin{array}{r} -\ 00001110 \\ -\ 00001110 \\ +\ 00001110 \end{array}$$

Effectual partial products

Parallel Input $\begin{array}{c} A_3 \\ \tilde{A}_0 \end{array}$  <<  +

$B_3$ $B_2$ $B_1$ $B_0$    Index    EB cycles

Serial Input

Subtraction is naturally the same with addition. We only need to add a simple circuit for negative values.

# Restricted Signed-Digit

Restrict EB: Restricted signed-digit (RSD)

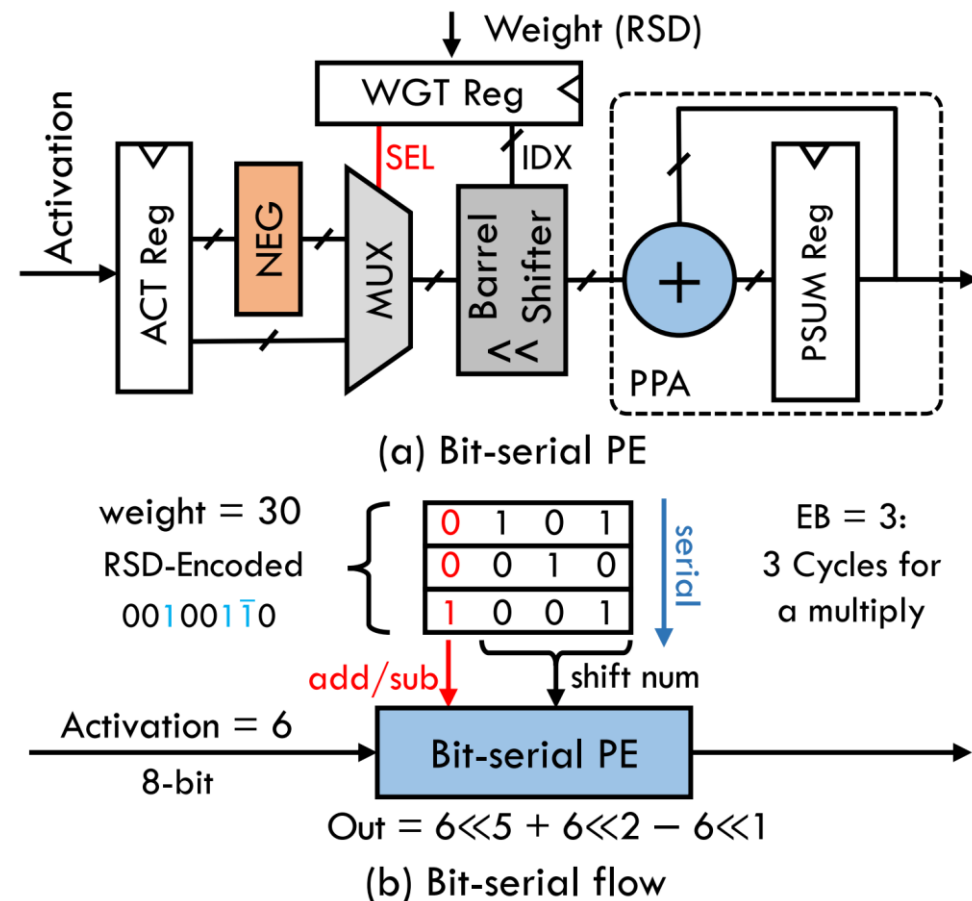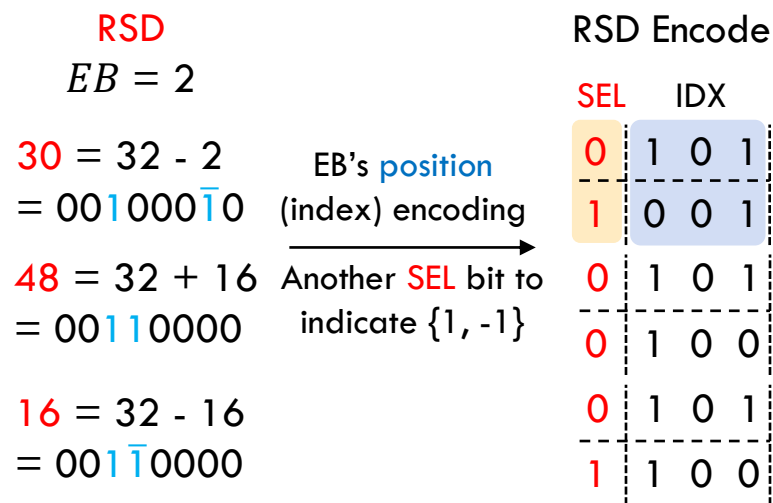| Original Numbers | 2's complement $EB = 2,$ | Error | RSD $EB = 2$ | Error |
|---|---|---|---|---|
| 30 (int8) = 00011110 | 30 → 24 = 00011110 | 6 | 30 = 32 - 2 = 00100010 | 0 |
| 46 (int8) = 00101110 | 46 → 40 = 00101110 | 6 | 48 = 32 + 16 = 00110000 | 2 |
| 56 (int8) = 00111000 | 56 → 48 = 00010001 | 8 | 56 = 64 - 8 = 01001000 | 0 |



RSD suits original weight distributions better

Algorithm: set up 2's integer power (0, 2, 4, 8,…) as the bases, iteratively find the base and sign bits according to the restriction.  30 -> get 32 and -2 in two iterations
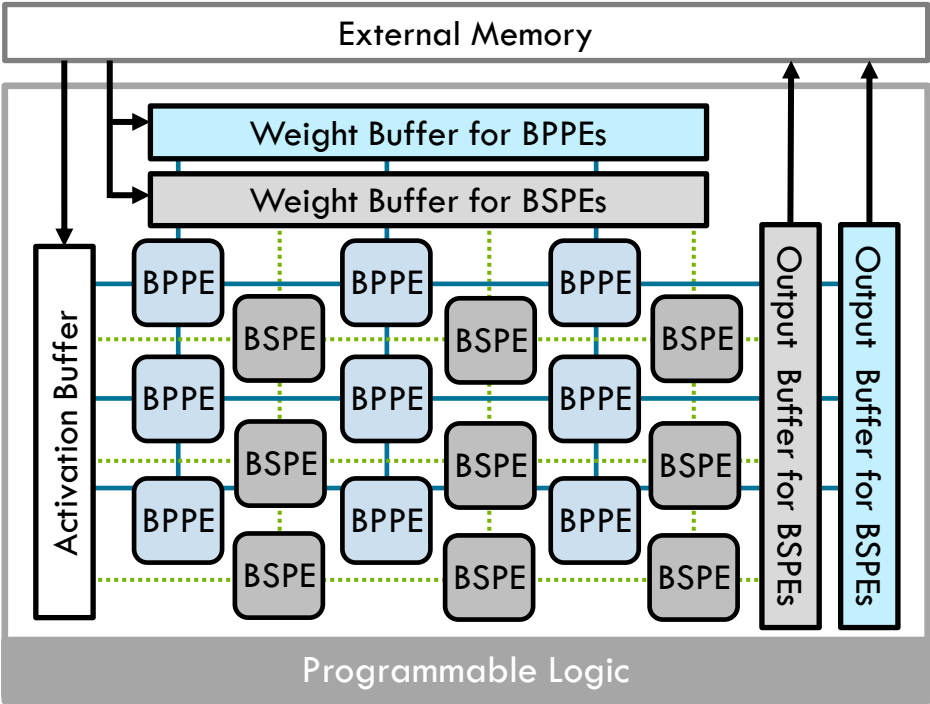
# Weights Encoding & Bit-serial PE Design
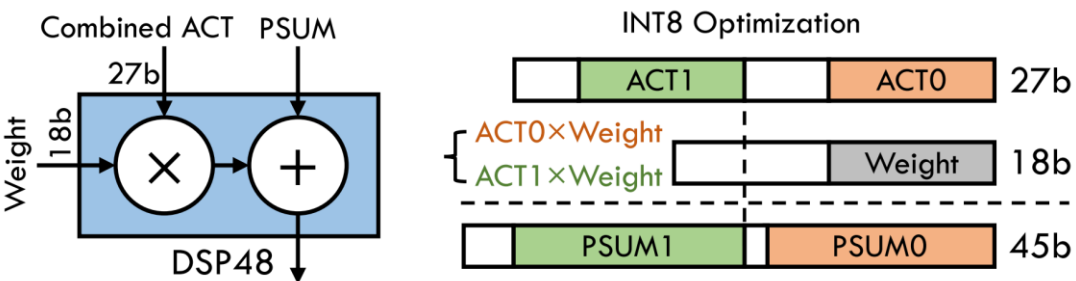
- We focused on layer-wise weights quantization in this work, and using the de facto 8-bit (int8) model as the starting point.



(a) Bit-serial PE

(b) Bit-serial flow

# Heterogeneous Architecture



INT8 multiplication on DSPs (standard representation)

|  | Weights | Activations |
|---|---|---|
| BSPE | RSD representation | Standard |
| BPPE | Standard | Standard |
| Consistence representation by mixing signed-digit (MSD) | | |

- BSPE & BPPE: Bit-serial & Bit-parallel PE, set up as a systolic array
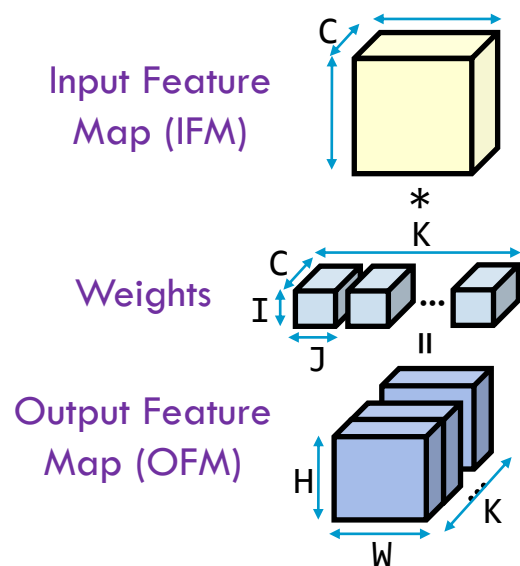- BSPEs are implemented on LUTs, while BPPEs are based on DSPs, running simultaneously

- Standard 2's complement can be regarded as a special case of signed-digit representation.

# Search Schedules

- Parameters that need to be searched by the scheduler (for each layer)

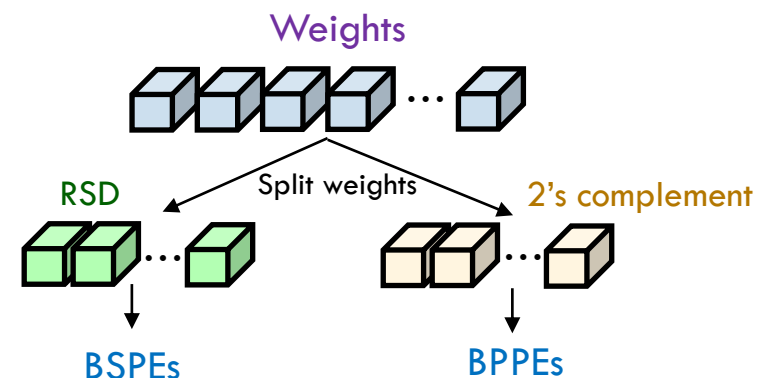Input Feature Map (IFM)

Weights

Output Feature Map (OFM)

Part I: 6-dimension for-loop for each layer
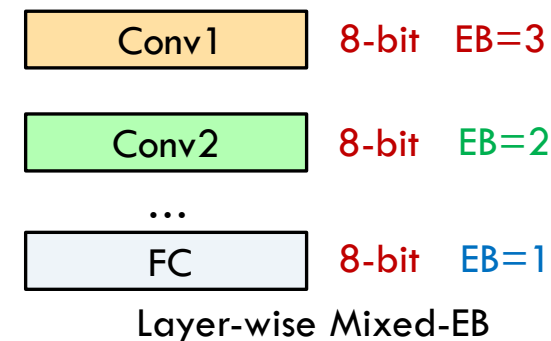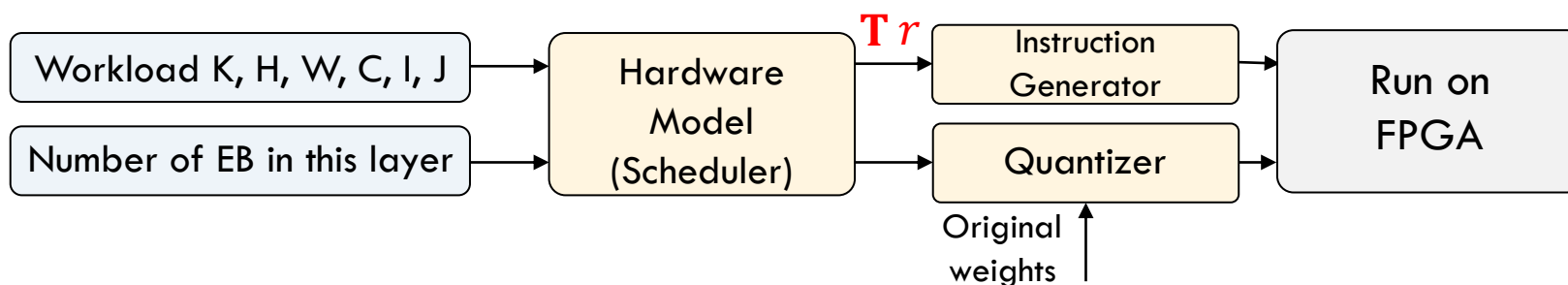
Tile size in each dimension:

$$\mathbf{T} = [t_K, t_H, t_W, t_C, t_I, t_J]$$

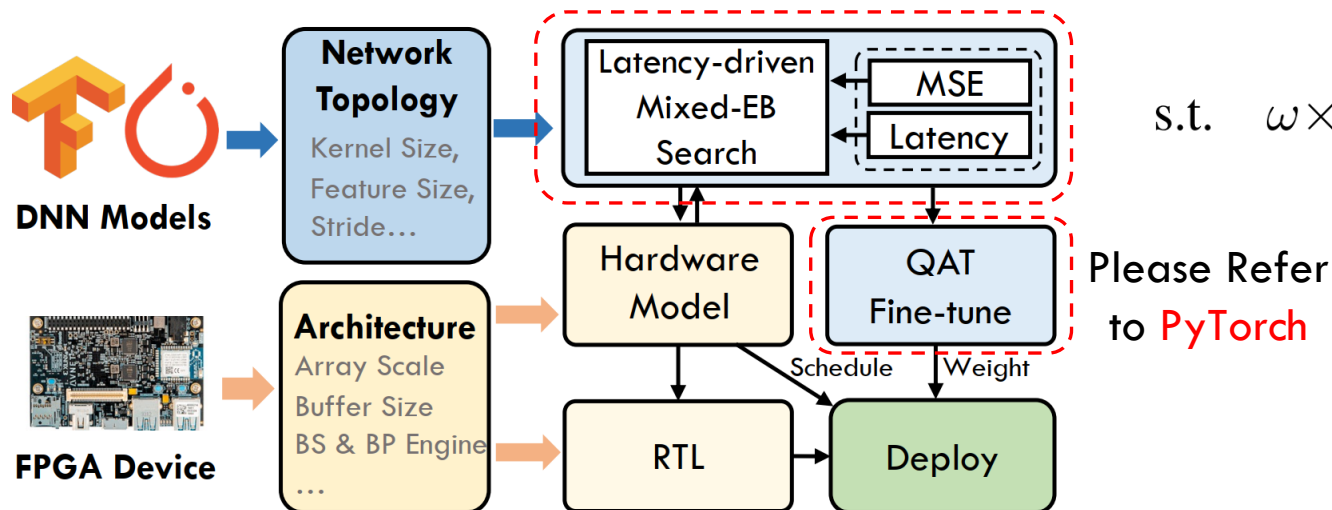Part II: Weight Split Ratio (workload partitioning)

Weights

RSD    Split weights    2's complement

BSPEs    BPPEs

Split ratio   $r = \dfrac{W_{BSPE}}{W_{Total}}$

| Workload K, H, W, C, I, J |
|---|
| Number of EB in this layer |

→ Hardware Model (Scheduler) → **T** $r$ → Instruction Generator → Run on FPGA

Quantizer

Original weights

| Conv1 | 8-bit | EB=3 |
| Conv2 | 8-bit | EB=2 |
| ... | | |
| FC | 8-bit | EB=1 |

Layer-wise Mixed-EB

# Hardware-aware Mixed-EB Quantization

- End-to-end framework



$$\min_{\mathbf{A}} \quad \sum_{j=1}^{m} MSE(j, \mathbf{A}[j])$$
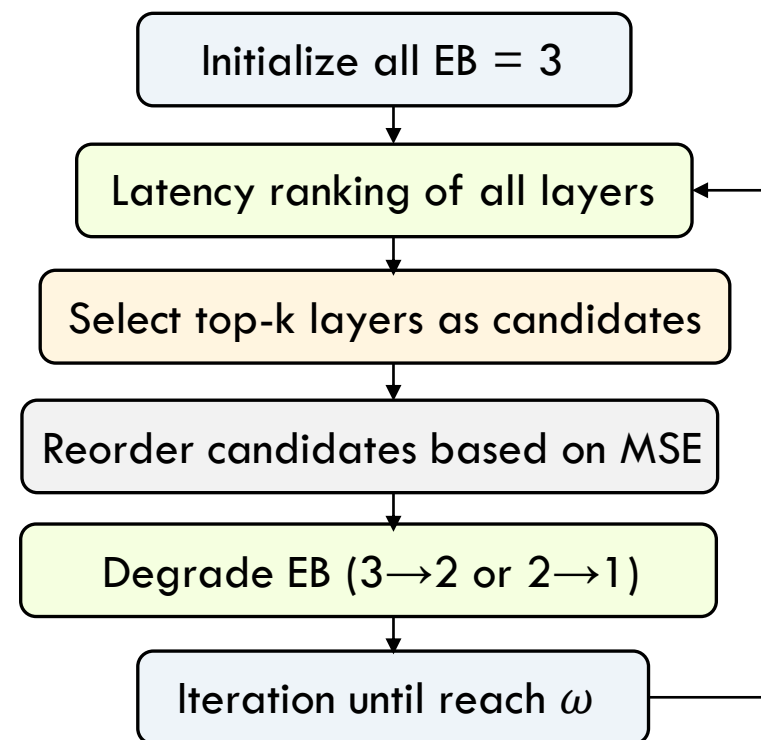
Objective: Minimize quantization error

$$\text{s.t.} \quad \omega \times \sum_{j=1}^{m} Lat_{\mathrm{L}}(j, \mathbf{A}[j]) \leq Lat_{\mathrm{base}}$$

Constraint: speedup ratio, $\omega$

Search Flow



- Inputs: DNN models and FPGA device constraints

- Mixed-EB search for each layer, under the constraint of MSE and latency (calculated based on the hardware model)

$$MSE = \sqrt{\Sigma_{i=1}^{n} \left( \frac{x - \hat{x}}{\sigma_i} \right)^2}$$

# Contents

- Motivation & Background

- The Proposed Methodology

- **Experiment Results**

- Conclusion & Future Works

# Mixed-EB Quantization

TABLE I: Mixed-EB speedup results with different constraints $\omega$. A larger $\omega$ means a higher speedup and a more aggressive quantization strategy.

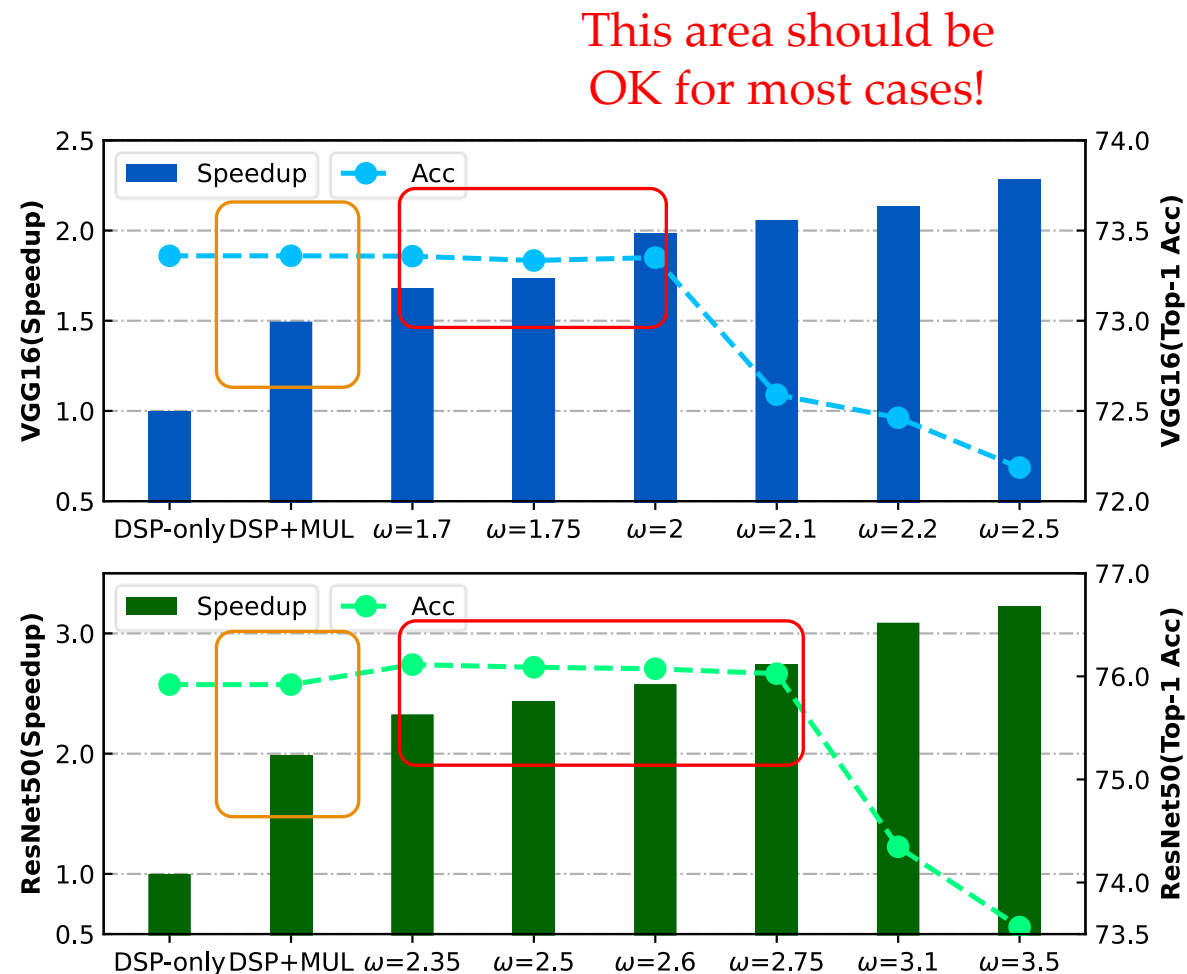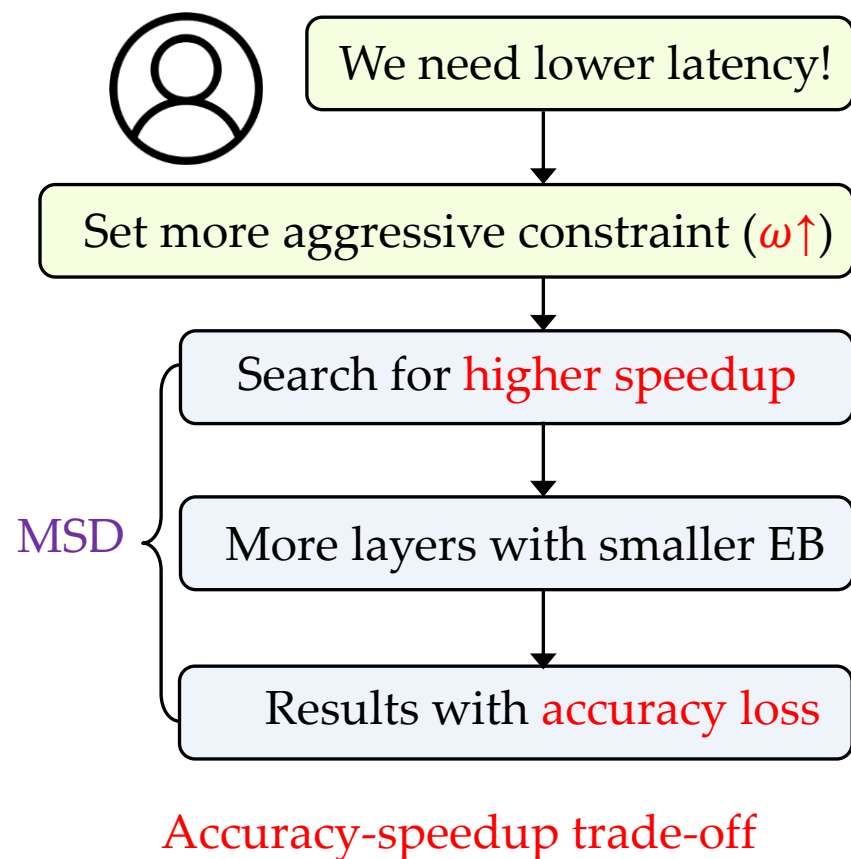| Model | $\omega$ | Layer-wise Mixed-EB Result **A** | Speedup |
|---|---|---|---|
| VGG-16 | 1.5 | [3 3 2 3 3 3 3 3 3 3 3 3 3 3 3 3] | 1.52 |
| | 1.6 | [2 3 2 2 2 3 2 2 3 3 3 2 3 3 3 2] | 1.60 |
| | 1.7 | [2 2 2 2 2 3 2 2 3 2 3 2 2 2 2 2] | 1.70 |
| | 2.0 | [2 2 2 2 2 2 2 2 2 2 2 2 2 2 1 2 2] | 2.00 |
| | 2.1 | [1 1 2 1 1 1 1 2 2 1 2 1 2 1 2 1] | 2.14 |
| | 2.2 | [1 1 2 1 1 1 1 2 1 1 2 1 2 1 1 1] | 2.24 |
| ResNet-18 | 1.5 | [3 3 3 3 3 3 3 3 3 3 3 3 3 3 2 2 3 3 3 3 3] | 1.51 |
| | 1.6 | [2 3 2 3 2 3 3 3 2 2 2 3 3 2 2 2 2 3 2 3 2] | 1.62 |
| | 1.65 | [2 3 2 3 2 3 2 3 2 2 2 3 3 2 2 2 2 3 2 3 2] | 1.65 |
| | 1.7 | [2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 1 2 2 2] | 1.71 |
| | 1.8 | [2 2 2 2 2 2 2 1 1 1 1 1 2 2 1 1 2 1 2 1 1] | 1.84 |
| | 1.9 | [2 2 2 2 2 2 2 1 1 1 1 1 2 2 1 1 1 1 2 1 1] | 1.90 |

$$\min_{\mathbf{A}} \quad \sum_{j=1}^{m} MSE(j, \mathbf{A}[j])$$

Objective: Minimize quantization error

$$\text{s.t.} \quad \omega \times \sum_{j=1}^{m} Lat_{\mathrm{L}}(j, \mathbf{A}[j]) \leq Lat_{\mathrm{base}}$$

Constraint: speedup ratio, $\omega$

The baseline is we only use DSPs for computation

- Device: Ultra-96
- Layer-wise EB configuration
- The final speedup meets the constraint of $\omega$

- Higher speedup constraint -> larger $\omega$
- More aggressive search for higher speedup, to reach the constraint

# Accuracy-Speedup Trade-off

We need lower latency!

Set more aggressive constraint ($\omega\uparrow$)

MSD
- Search for higher speedup
- More layers with smaller EB
- Results with accuracy loss

Accuracy-speedup trade-off

This area should be OK for most cases!
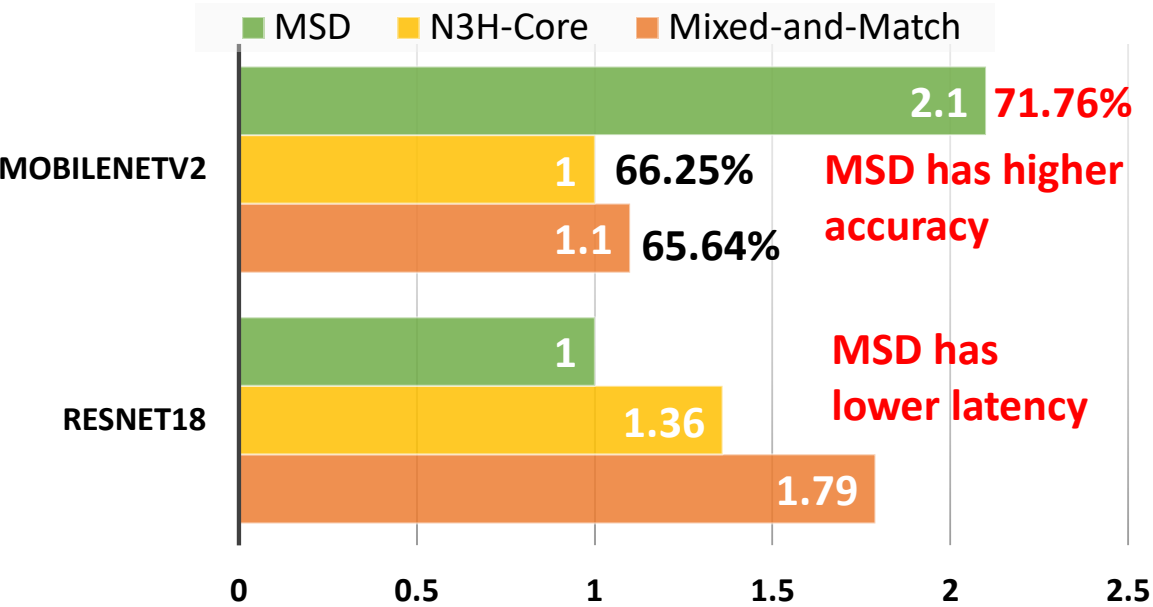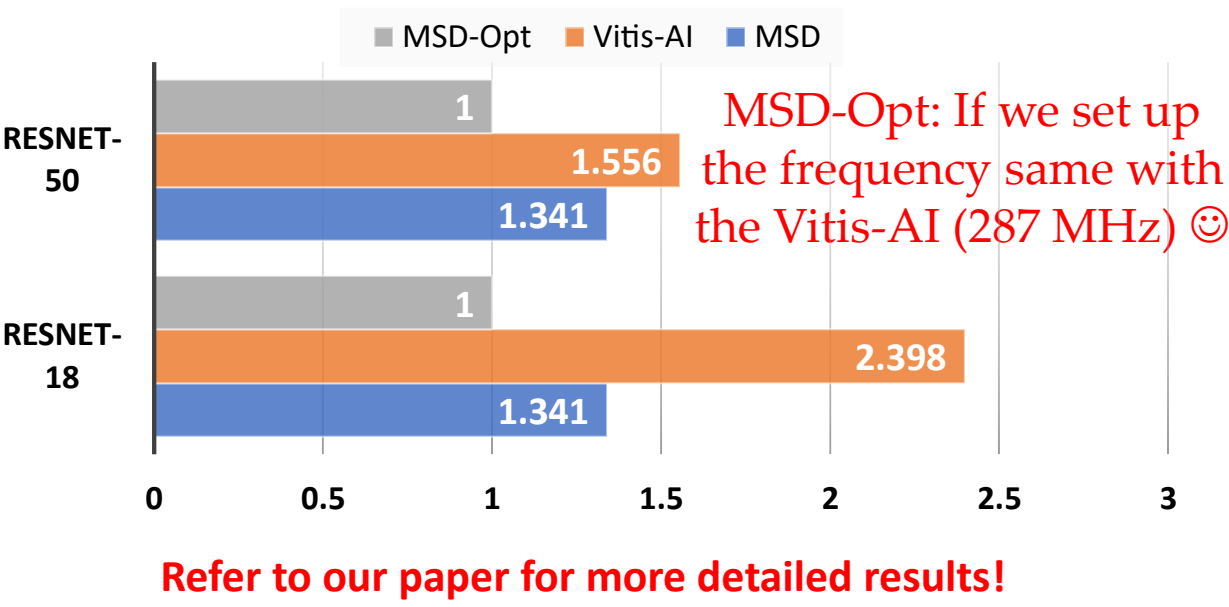


Also, RSD-based bit-serial scheme is more efficient than conventional parallel design on LUTs!

# Comparison with Previous Works



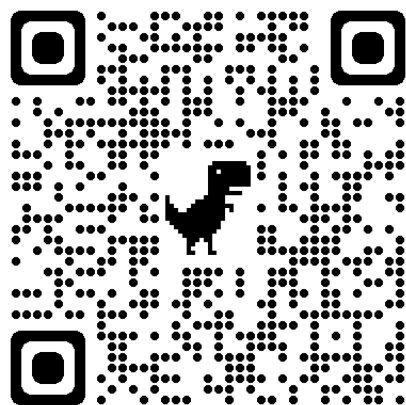**NORMALIZED LATENCY Performance with Accuracy on xc7z020**

Legend: ■ MSD ■ N3H-Core ■ Mixed-and-Match

MOBILENETV2:
- MSD: 2.1 — **71.76%**
- N3H-Core: 1 — **66.25%**
- Mixed-and-Match: 1.1 — **65.64%**

**MSD has higher accuracy**

RESNET18:
- MSD: 1
- N3H-Core: 1.36
- Mixed-and-Match: 1.79

**MSD has lower latency**

**Normalized Latency Performance on Ultra-96 Device**

Legend: ■ MSD-Opt ■ Vitis-AI ■ MSD

RESNET-50:
- MSD-Opt: 1
- Vitis-AI: 1.556
- MSD: 1.341

RESNET-18:
- MSD-Opt: 1
- Vitis-AI: 2.398
- MSD: 1.341

MSD-Opt: If we set up the frequency same with the Vitis-AI (287 MHz) ☺

**Refer to our paper for more detailed results!**

# Conclusion & Future Works

- MSD framework:
  - is a heterogeneous DNN acceleration framework that utilizes both LUTs and DSPs as computation resources and to exploit bit-sparsity.
  - fine-tunes and encodes the DNN weights into a bit-sparsity-aware format, making the bit-serial computation on LUTs more efficient.
  - uses a latency-driven algorithm to search for the optimal schedule and trade-off based on layer-wise mixed EB.
- We will explore more efficient scheduling methods and exploit FPGA-layout-tailored hardware design to enhance the hardware clock frequency further.
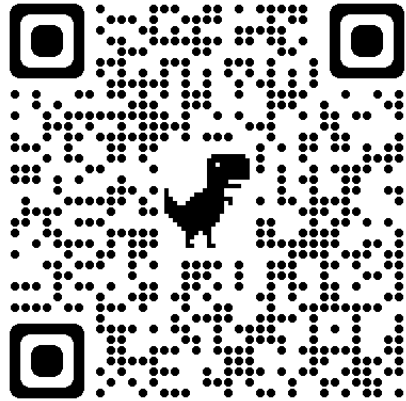


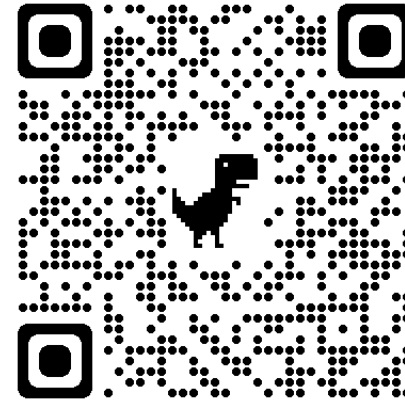More info of this paper in our group site!



MSD open-source in GitHub

# Thanks for Your Listening

## Find us in the poster session!

## Q & A



More info of this paper in our group site!



MSD open-source in GitHub