

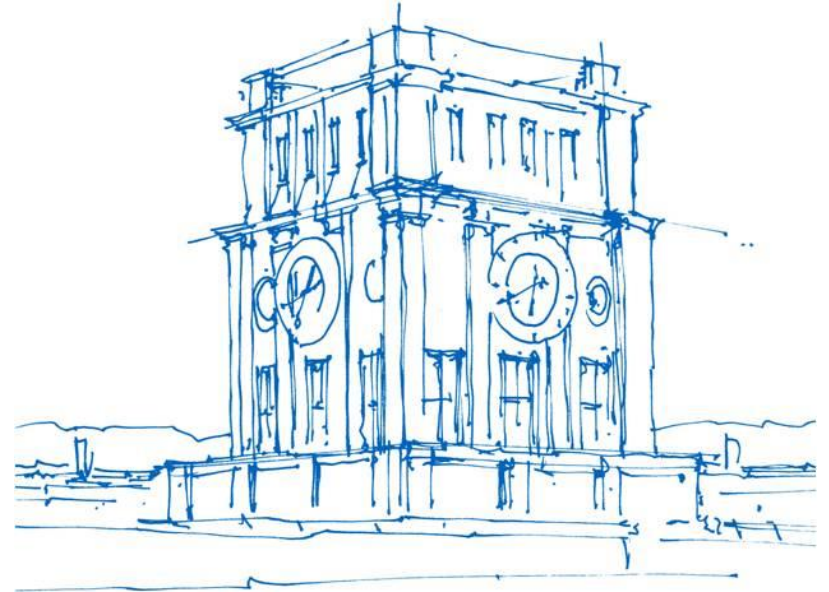
OrthrusPE: Runtime Reconfigurable Processing Elements for Binary Neural Networks

Nael Fafous

Technical University of Munich

Department of Electrical and Computer Engineering

Chair of Integrated Systems



Uhrenturm der TUM

Outline

- Optimization of Convolutional Neural Networks
- Challenges of Binary Neural Networks
- Motivation for Runtime Reconfigurable BNN Processing Elements
- OrthrusPE: Dual Modes
- Experimentation and Results

Optimization of Convolutional Neural Networks



Structural

Algorithmic

Hardware

Optimization of Convolutional Neural Networks

Structural

Low-Rank
Decomp.

Pruning

Quantization

Algorithmic

Winograd

FFT

GEMM

Hardware

Vectorization

Dataflow

Partitioning

Optimization of Convolutional Neural Networks

Structural

Low-Rank
Decomp.

Pruning

Quantization

Tucker
Decomp.

Channel-wise

Filter
Sharing

Weight
Sharing

CP
Decomp.

Filter-wise

Floating
Point

Fixed Point

Element-wise

Binarization

Winograd
Sparsity

Pow. 2

Spectral
Sparsity

Log Base

Algorithmic

Winograd

FFT

GEMM

Hardware

Vectorization

Dataflow

Partitioning

Array
Mapping

In

Out

OS

IS

IFmap

Filter

RS

WS

Bypass

IO

OW

IW

Overview of Binary Neural Networks



- Binary Weights
- Binary Activations
- Binary Weights AND Activations

Overview of Binary Neural Networks

- Binary Weights
- Binary Activations
- Binary Weights AND Activations → Replace Multiplications by XNOR ops
Replace Accumulations by Popcount ops



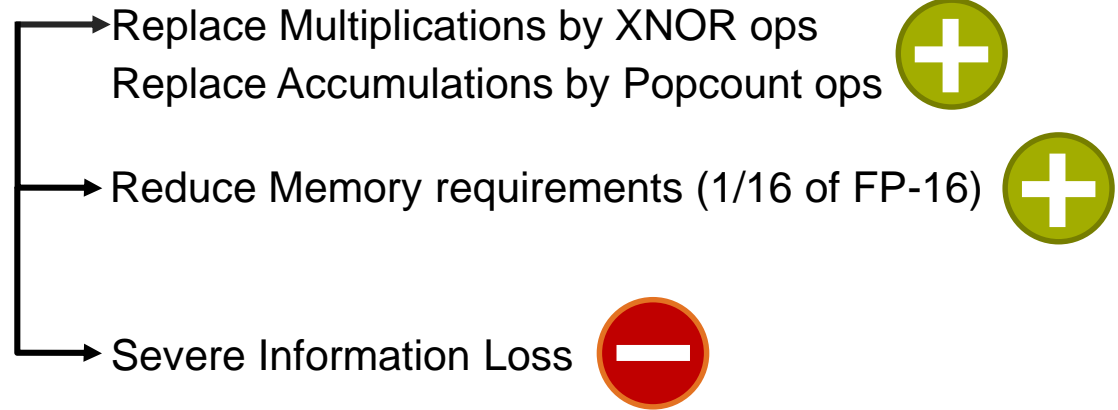
Overview of Binary Neural Networks

- Binary Weights
- Binary Activations
- Binary Weights AND Activations
 - Replace Multiplications by XNOR ops
 - Replace Accumulations by Popcount ops
 - Reduce Memory requirements (1/16 of FP-16)



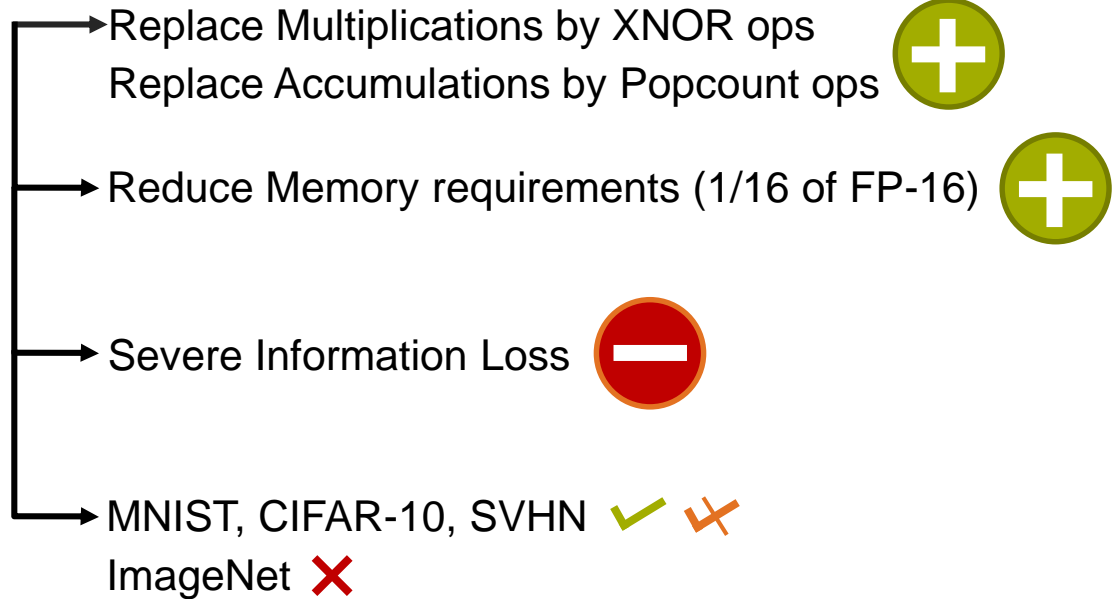
Overview of Binary Neural Networks

- Binary Weights
- Binary Activations
- Binary Weights AND Activations



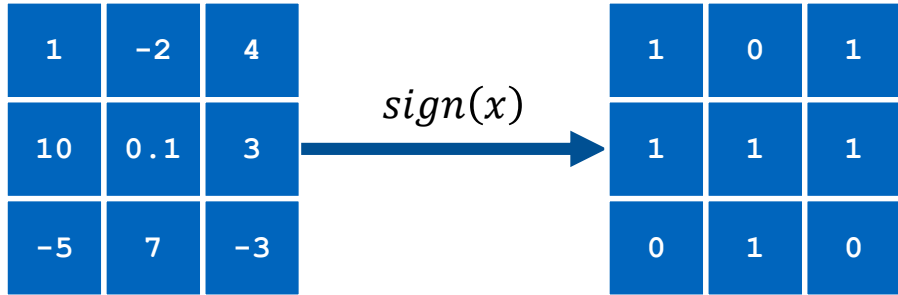
Overview of Binary Neural Networks

- Binary Weights
- Binary Activations
- Binary Weights AND Activations



Overview of Binary Neural Networks

Naïve Binarization

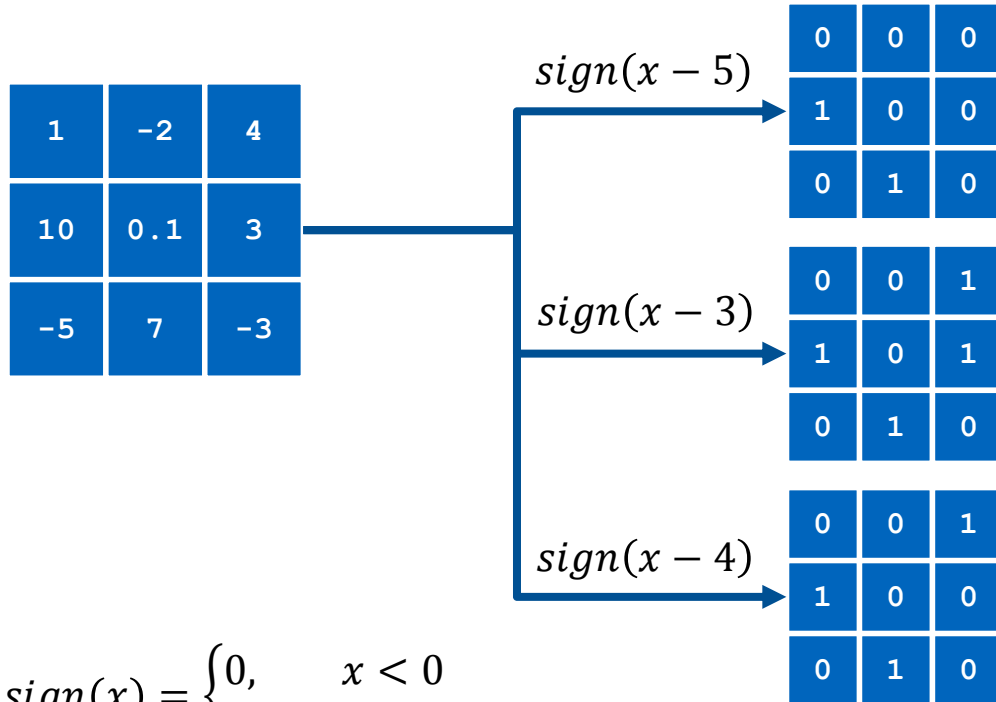


Severe information loss:
10 and 0.1 have the same effect on the network.

$$sign(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

Overview of Binary Neural Networks

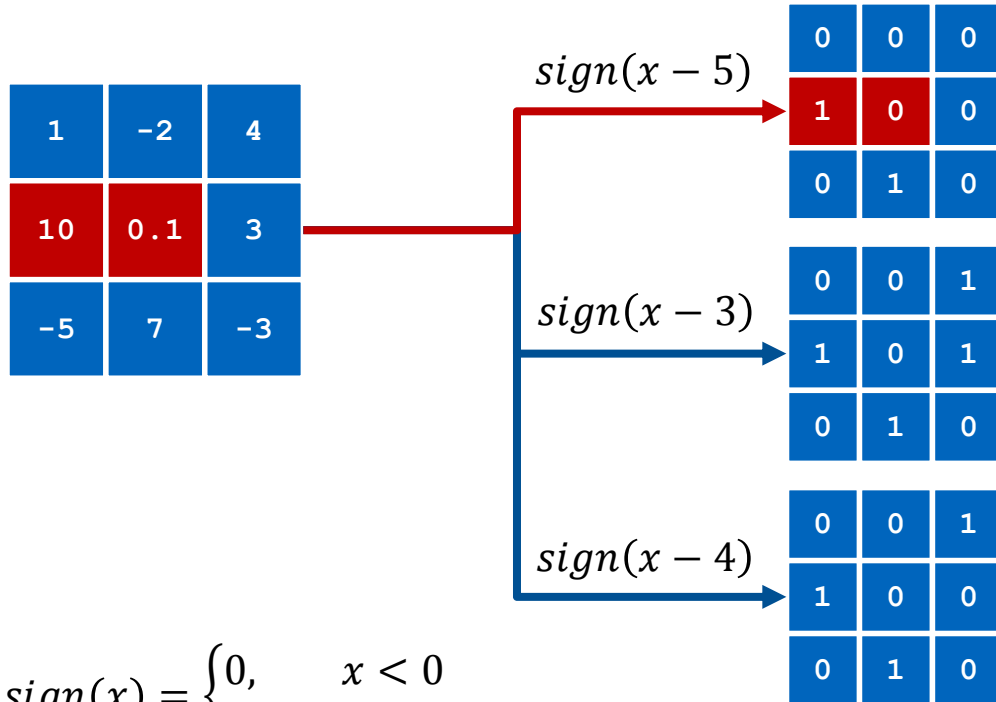
Approximation through binary bases



$$sign(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

Overview of Binary Neural Networks

Approximation through binary bases



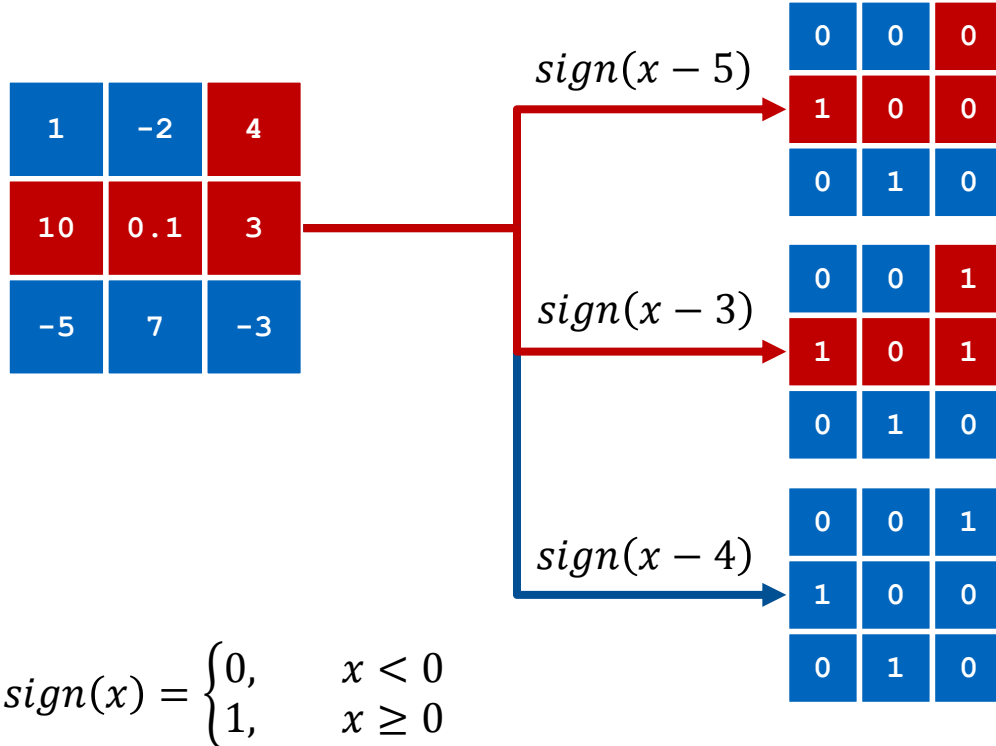
$$sign(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

Captured Information:

0.1 is less positive than **10**

Overview of Binary Neural Networks

Approximation through binary bases



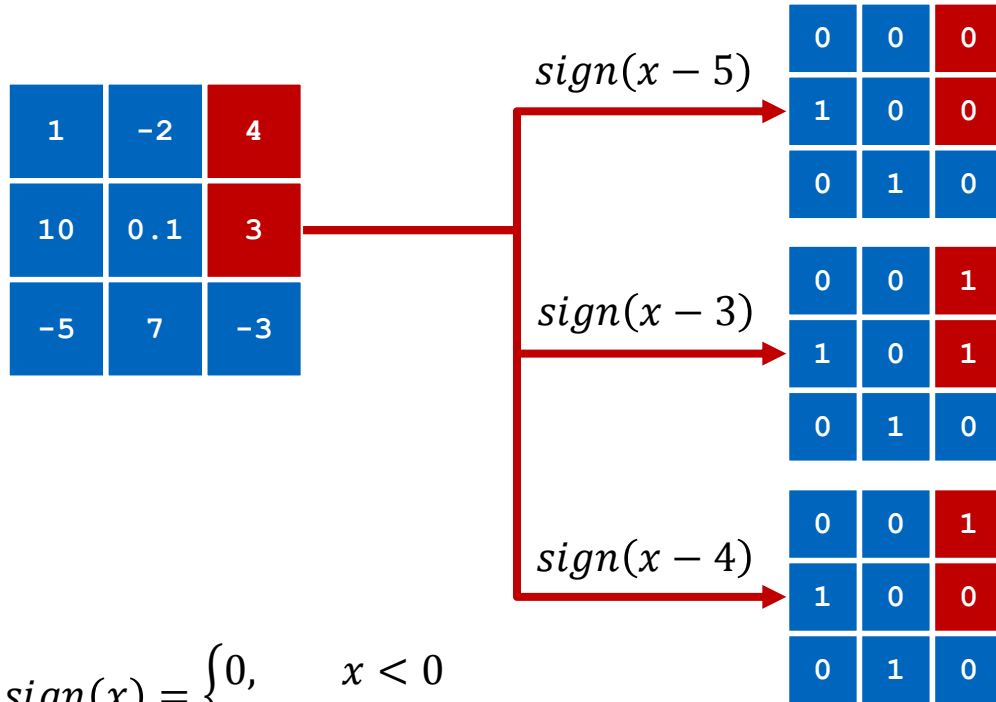
Captured Information:

0.1 is less positive than **10**

4 and **3** lie between **10** and **0.1**

Overview of Binary Neural Networks

Approximation through binary bases



$$sign(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

Captured Information:

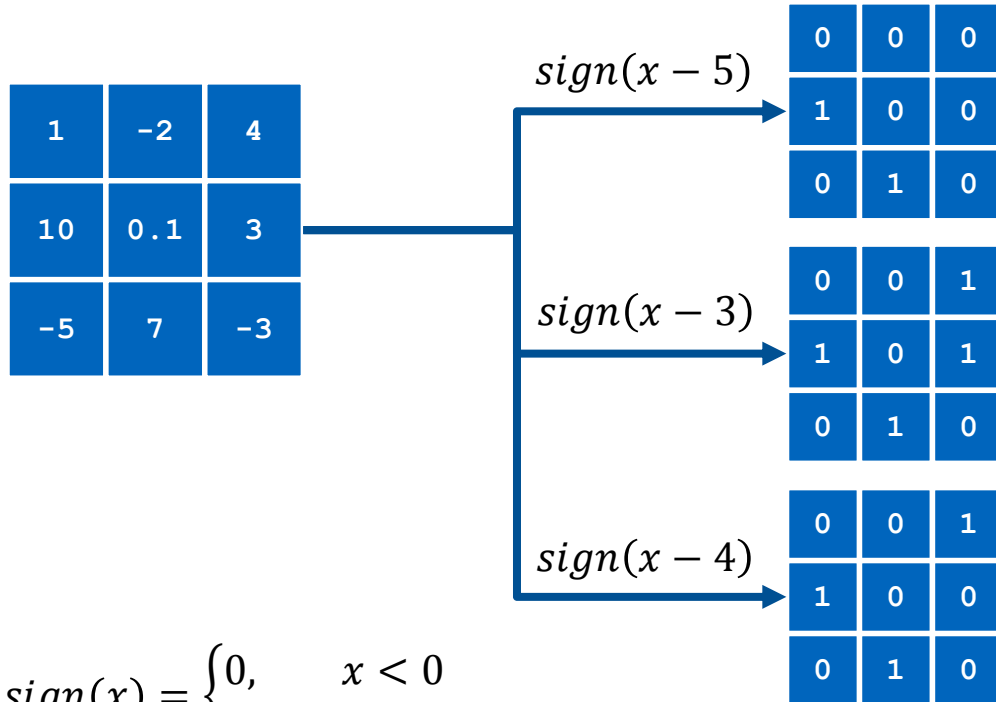
0.1 is less positive than **10**

4 and **3** lie between **10** and **0.1**

4 and **3** are single unit away from each other

Overview of Binary Neural Networks

Approximation through binary bases



$$sign(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

Captured Information:

0.1 is less positive than **10**

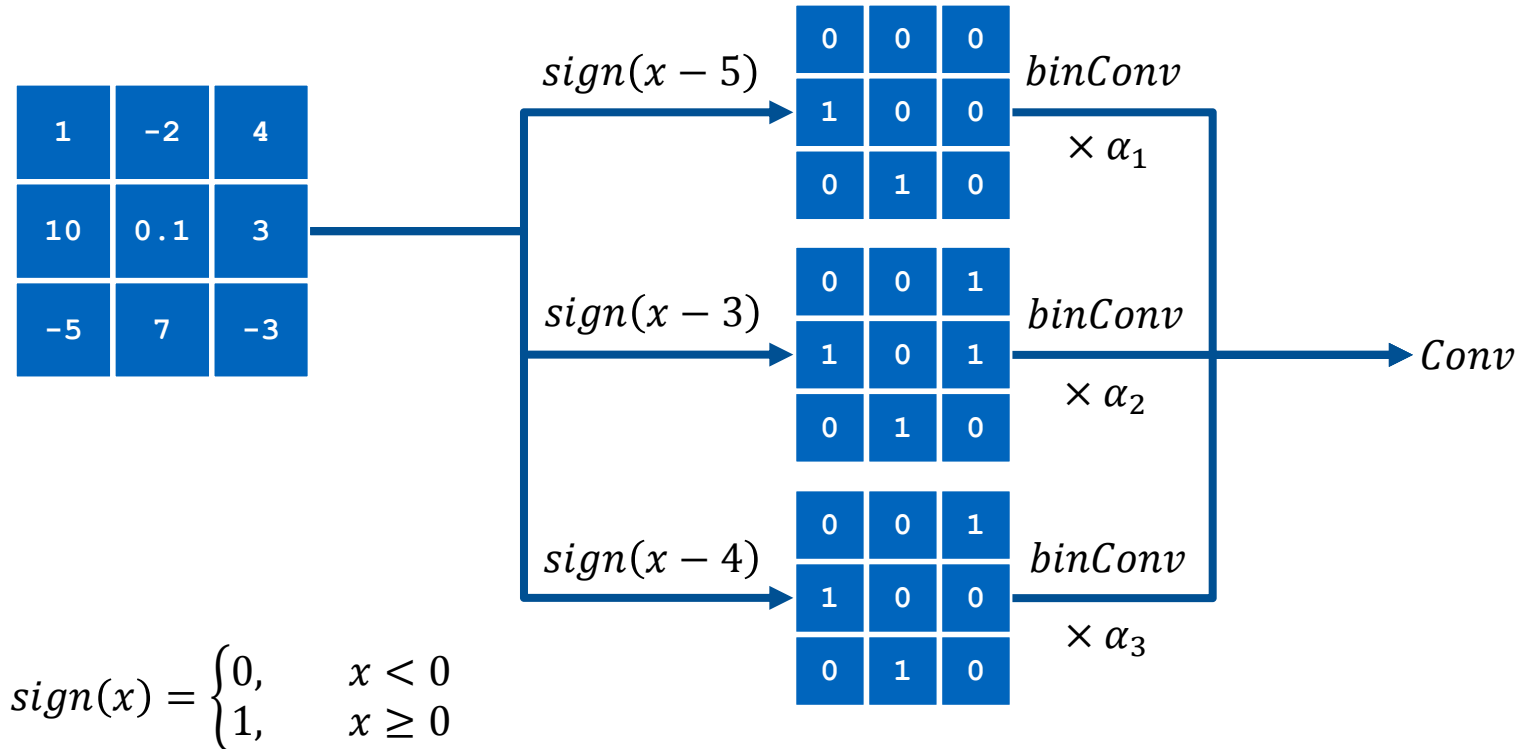
4 and 3 lie between **10** and **0.1**

4 and 3 are single unit away from each other

Information can be extracted for negative numbers, e.g. $sign(x + 3)$

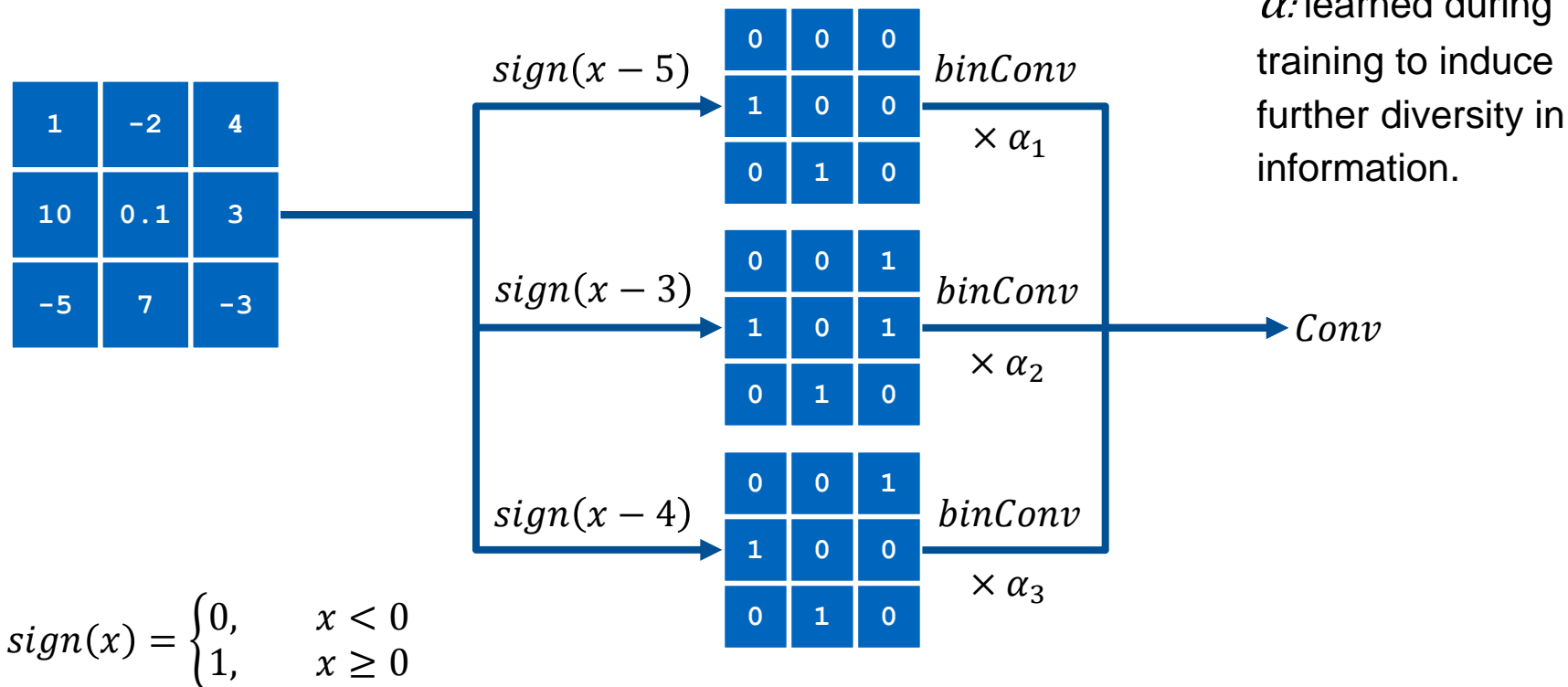
Overview of Binary Neural Networks

Approximation through binary bases



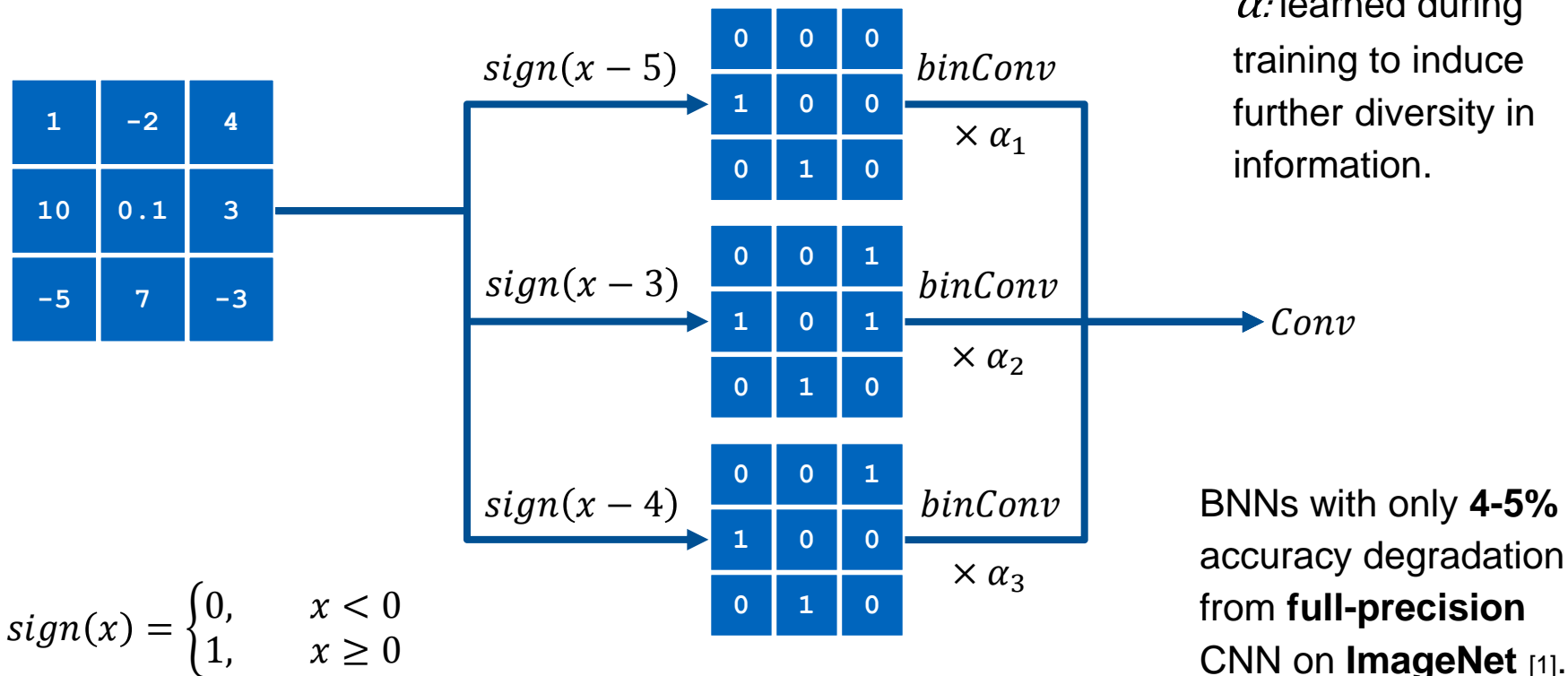
Overview of Binary Neural Networks

Approximation through binary bases



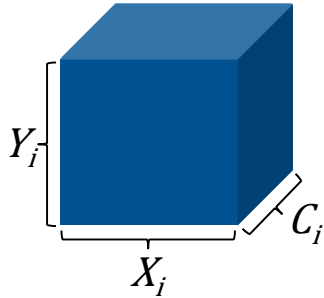
Overview of Binary Neural Networks

Approximation through binary bases

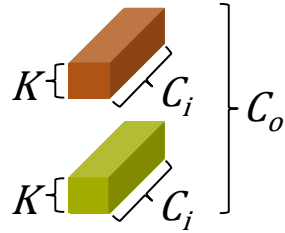


How Binary are Binary Neural Networks?

$$A^{l-1} \in \mathbb{R}^{X_i \times Y_i \times C_i}$$



$$W^l \in \mathbb{R}^{K \times K \times C_i \times C_o}$$

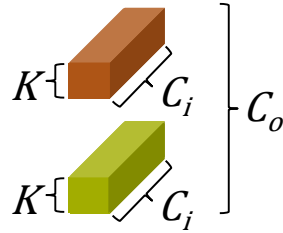
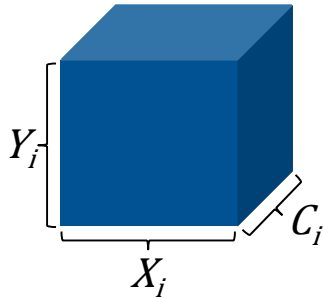


$$A^l = \text{Conv}(W^l, A^{l-1})$$

How Binary are Binary Neural Networks?

$$H^{l-1} \in \mathbb{B}^{X_i \times Y_i \times C_i \times N}$$

$$B^l \in \mathbb{B}^{K \times K \times C_i \times M \times C_o}$$



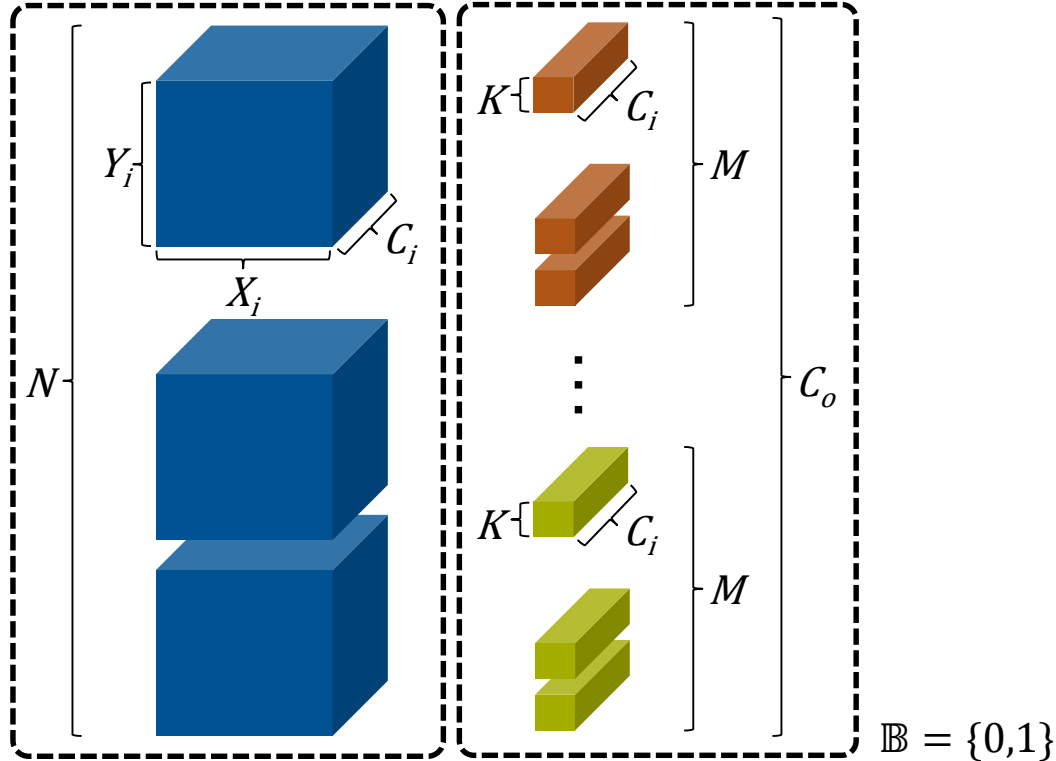
$$A^l = \text{Conv}(W^l, A^{l-1})$$

$$\mathbb{B} = \{0,1\}$$

How Binary are Binary Neural Networks?

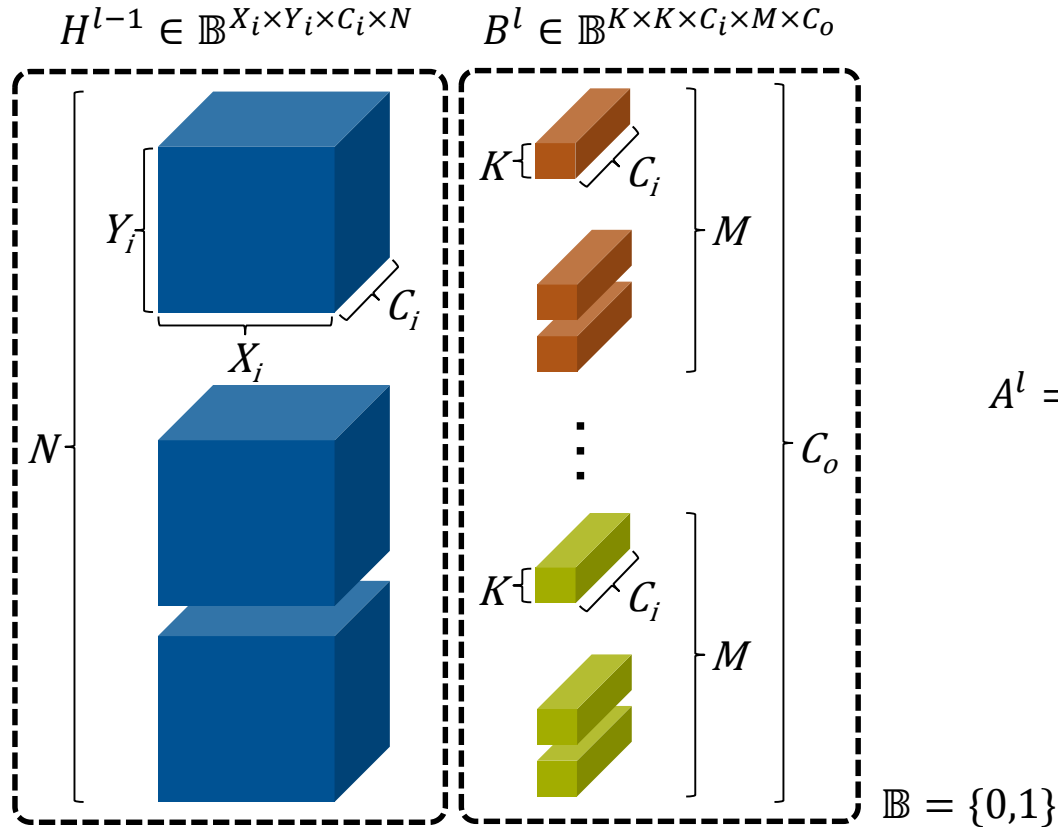
$$H^{l-1} \in \mathbb{B}^{X_i \times Y_i \times C_i \times N}$$

$$B^l \in \mathbb{B}^{K \times K \times C_i \times M \times C_o}$$



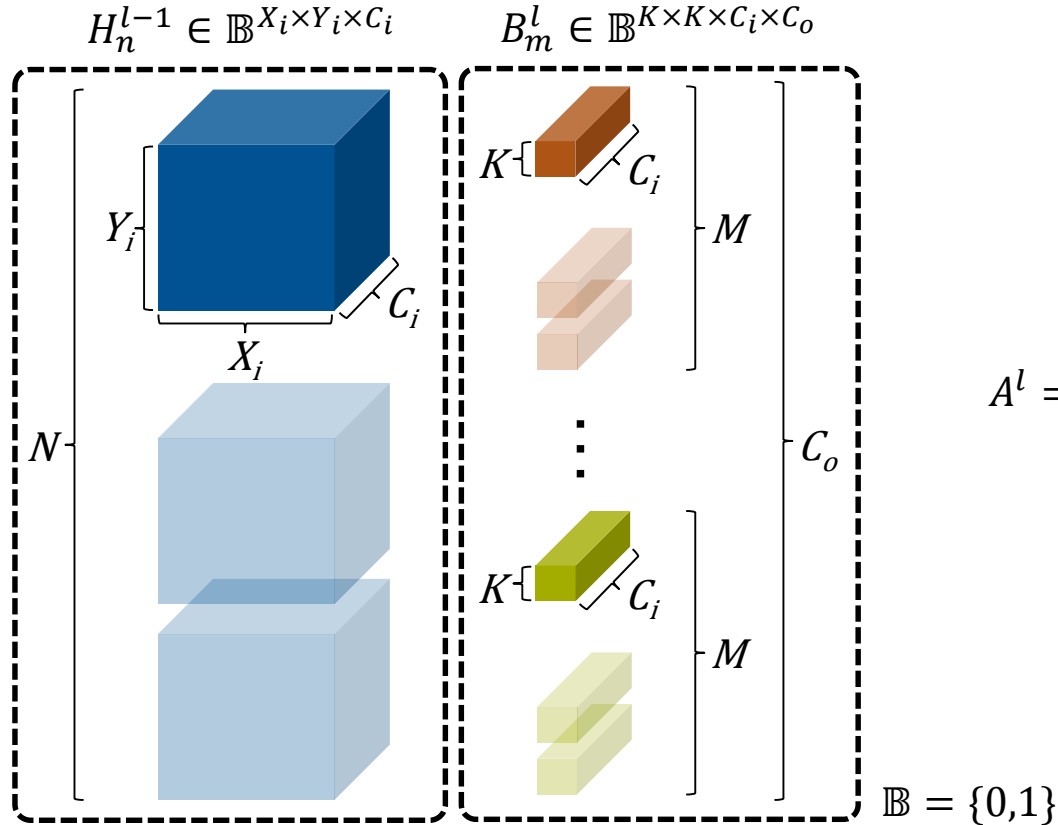
$$A^l = \text{Conv}(W^l, A^{l-1})$$

How Binary are Binary Neural Networks?



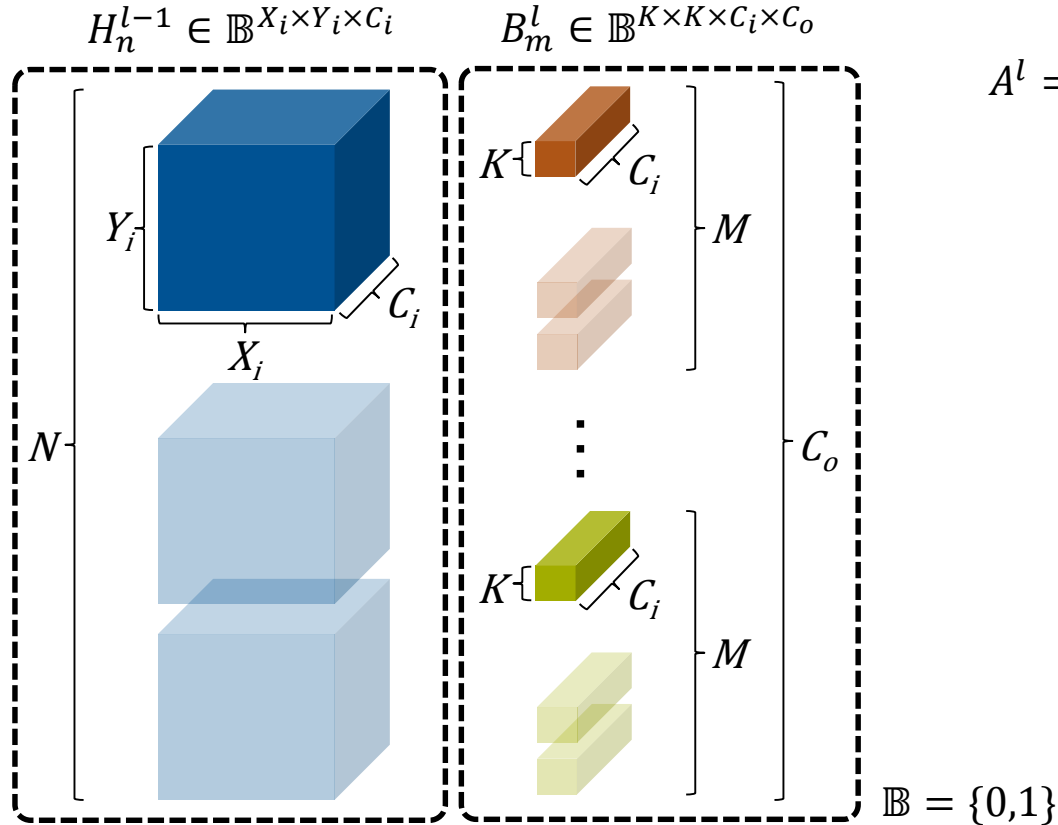
$$A^l = \sum_{m=1}^M \sum_{n=1}^N \alpha_m \beta_n \text{BinConv}(B_m^l, H_n^{l-1})$$

How Binary are Binary Neural Networks?



$$A^l = \sum_{m=1}^M \sum_{n=1}^N \alpha_m \beta_n \text{BinConv}(B_m^l, H_n^{l-1})$$

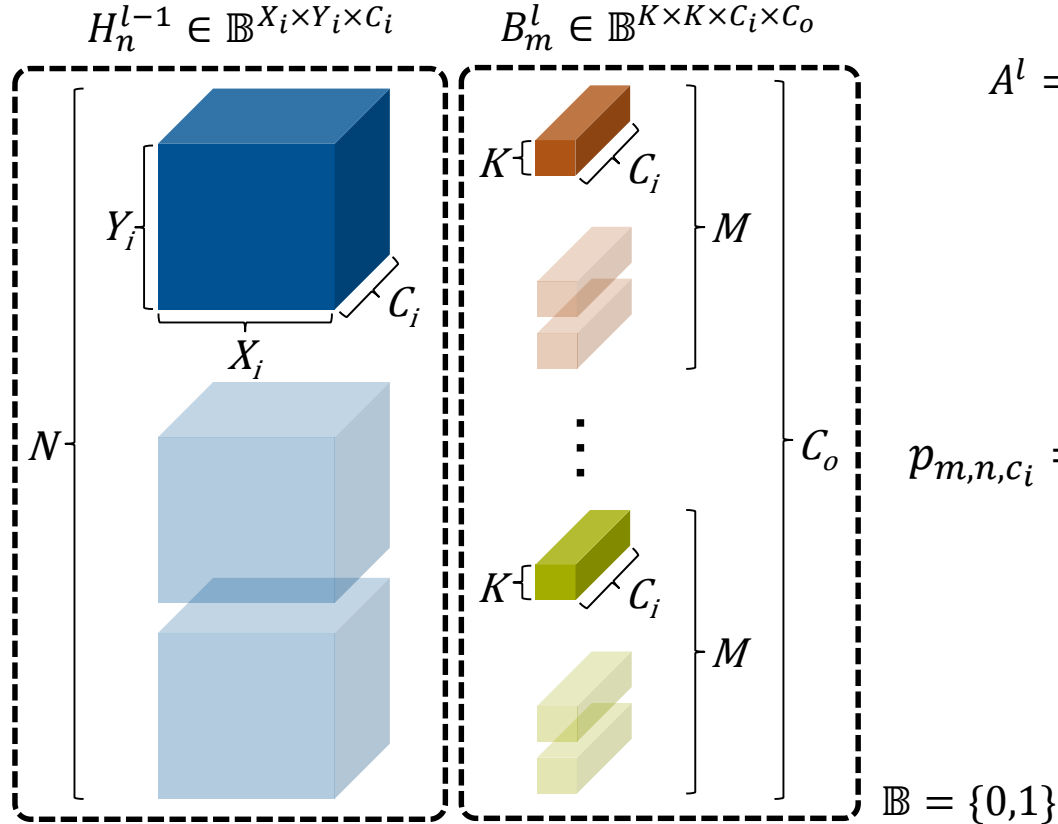
How Binary are Binary Neural Networks?



$$A^l = \sum_{m=1}^M \sum_{n=1}^N \alpha_m \beta_n \text{BinConv}(B_m^l, H_n^{l-1})$$

$$a_{m,n} = \sum_{c_i=1}^{C_i} (p_{m,n,c_i})$$

How Binary are Binary Neural Networks?

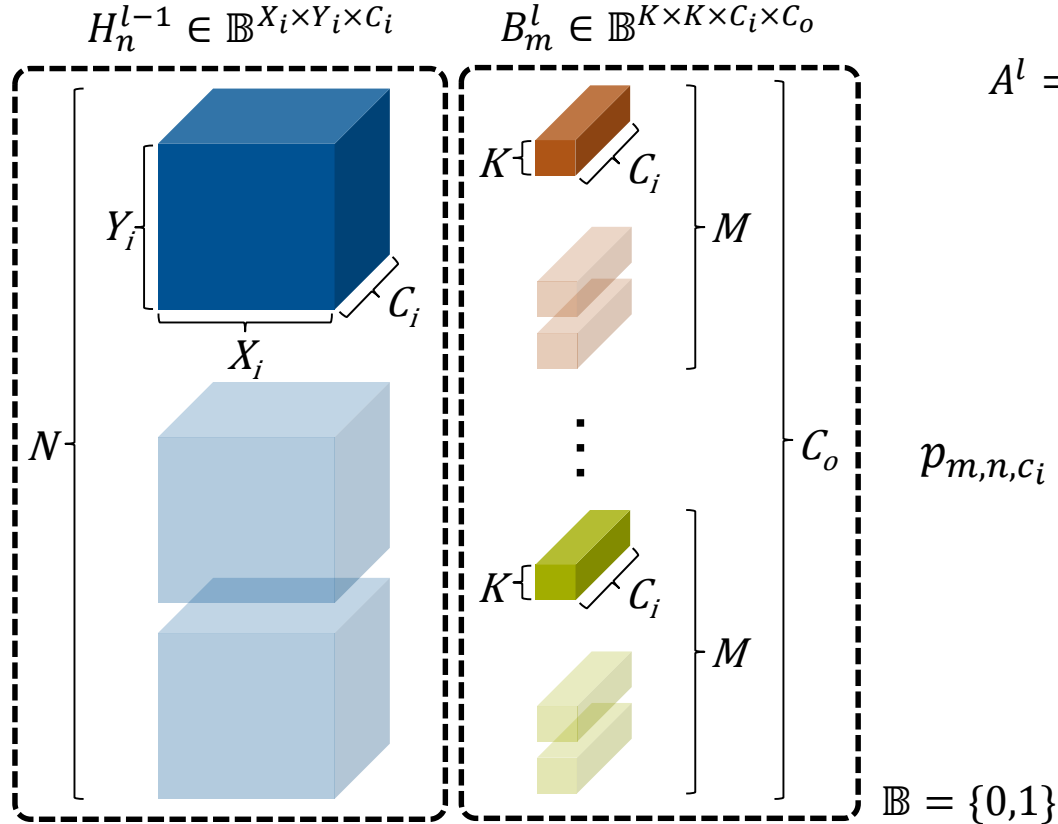


$$A^l = \sum_{m=1}^M \sum_{n=1}^N \alpha_m \beta_n \text{BinConv}(B_m^l, H_n^{l-1})$$

$$a_{m,n} = \sum_{c_i=1}^{C_i} (p_{m,n,c_i})$$

$$p_{m,n,c_i} = \sum_{k_x=1}^K \sum_{k_y=1}^K \text{xnor}(b_{k_x, k_y}, h_{x_i+k_x, y_i+k_y})$$

How Binary are Binary Neural Networks?

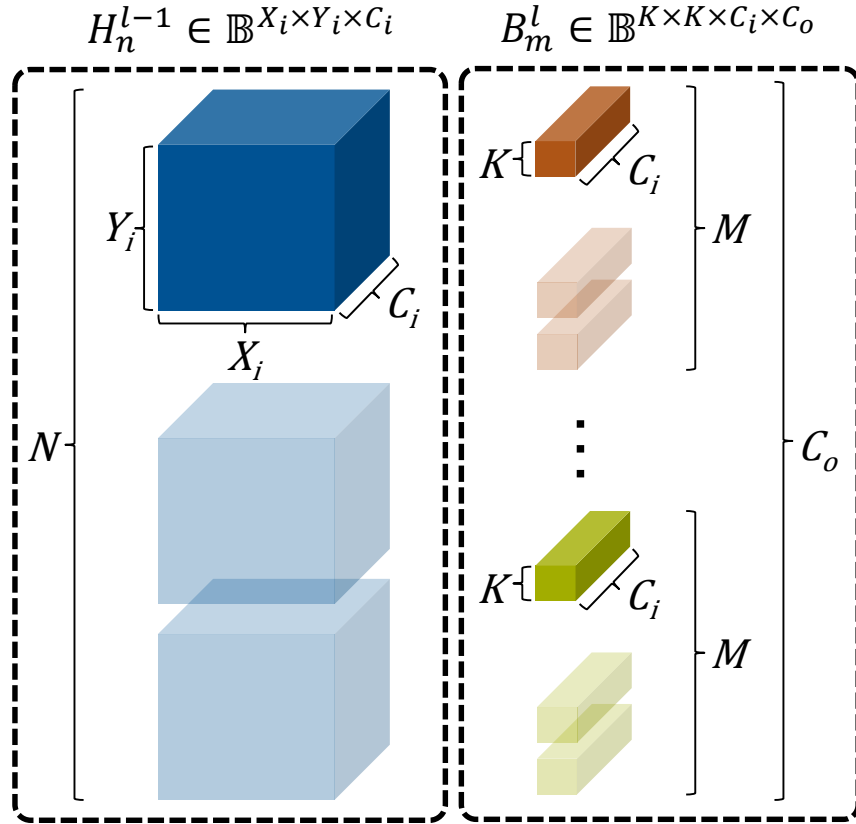


$$A^l = \sum_{m=1}^M \sum_{n=1}^N \alpha_m \beta_n \text{BinConv}(B_m^l, H_n^{l-1})$$

$$a_{m,n} = \sum_{c_i=1}^{C_i} (p_{m,n,c_i})$$

$$p_{m,n,c_i} = \text{popcnt}(\text{xnor}(b_{k_x, k_y}, h_{x_i+k_x, y_i+k_y}))$$

How Binary are Binary Neural Networks?



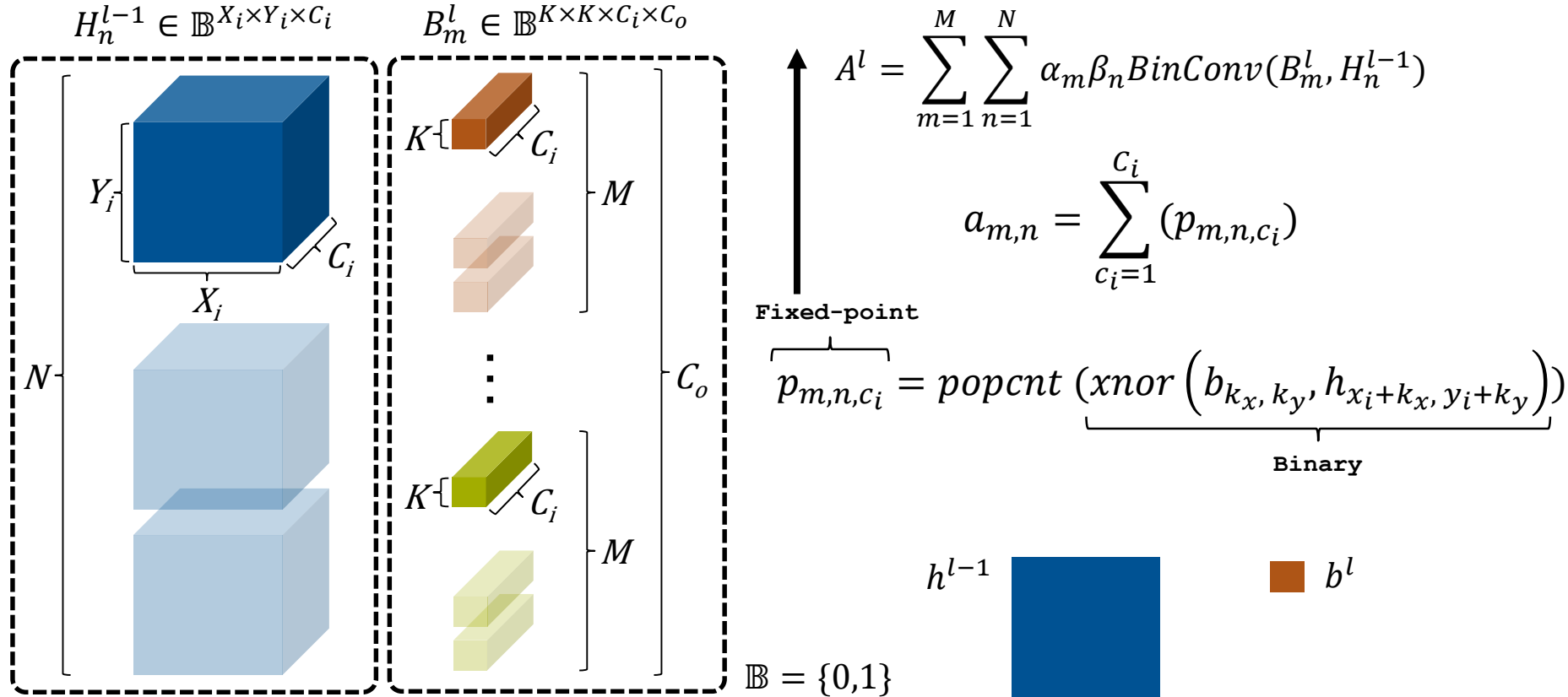
$$A^l = \sum_{m=1}^M \sum_{n=1}^N \alpha_m \beta_n \text{BinConv}(B_m^l, H_n^{l-1})$$

$$a_{m,n} = \sum_{c_i=1}^{C_i} (p_{m,n,c_i})$$

$$p_{m,n,c_i} = \text{popcnt}(\text{xnor}(b_{K_x, K_y}, h_{x_i+K_x, y_i+K_y}))$$



How Binary are Binary Neural Networks?

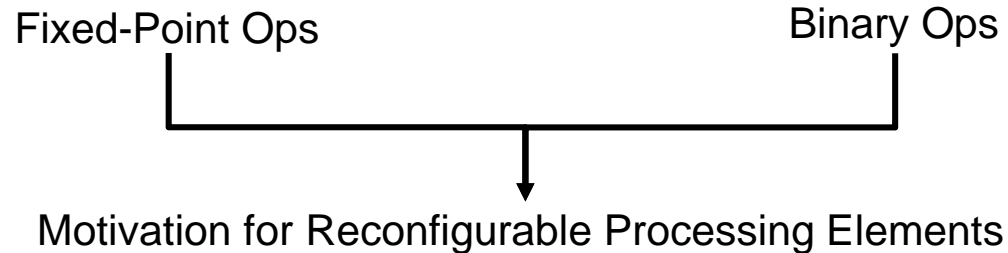


More FP in Binary Neural Networks

- Binary Weight and Activation Bases (Scale/Shift)
- First Layer Remains Non-Binarized
- Batch Normalization

More FP in Binary Neural Networks

- Binary Weight and Activation Bases (Scale/Shift)
- First Layer Remains Non-Binarized
- Batch Normalization



OrthrusPE: Runtime Reconfigurable PEs for BNNs



Dual Modes

- **Fixed-precision mode:** First layer, Batch-norm, Scale and Shift Operations
- **Binary mode:** SIMD Binary Hadamard Products and Popcounts
- Achieved with high resource reuse

OrthrusPE: Binary Mode

Efficient SIMD Binary Hadamard Product Execution

Input: $h^{l-1} \in H_n^{l-1}$

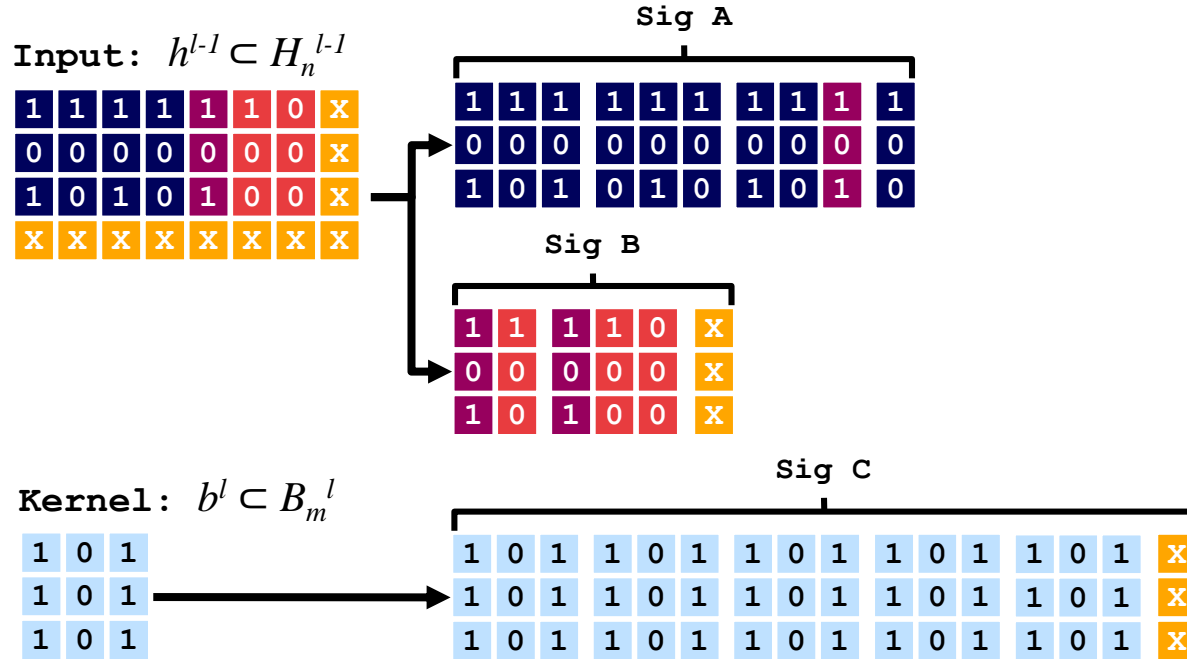
1	1	1	1	1	1	0	x
0	0	0	0	0	0	0	x
1	0	1	0	1	0	0	x
x	x	x	x	x	x	x	x

Kernel: $b^l \in B_m^l$

1	0	1
1	0	1
1	0	1

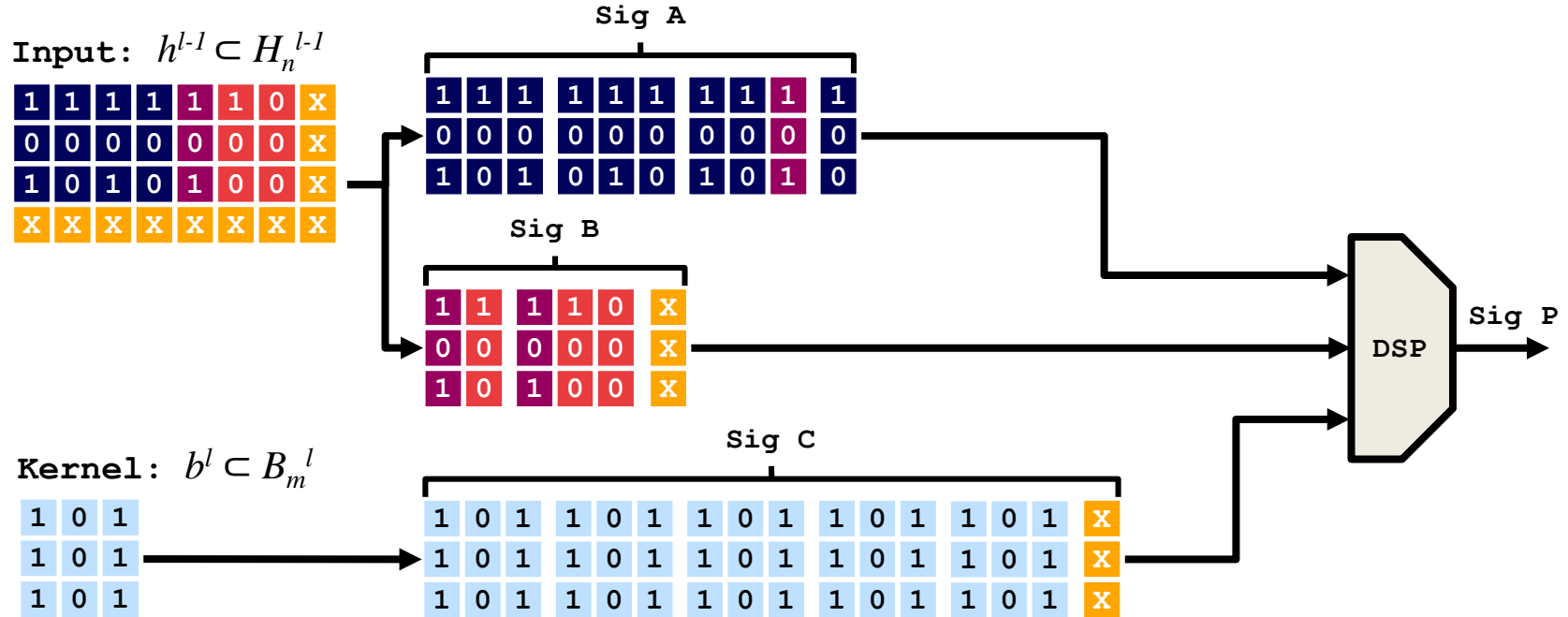
OrthrusPE: Binary Mode

Efficient SIMD Binary Hadamard Product Execution



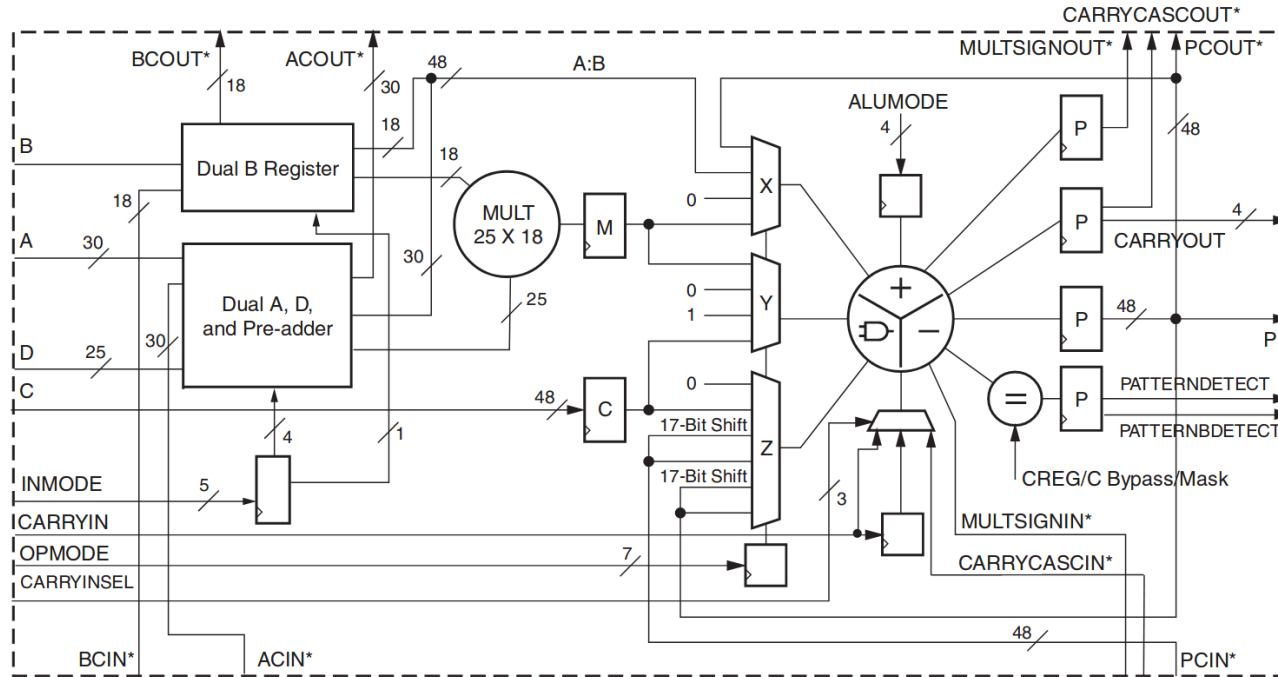
OrthrusPE: Binary Mode

Efficient SIMD Binary Hadamard Product Execution



OrthrusPE: Dual Modes

Runtime Reconfigurability



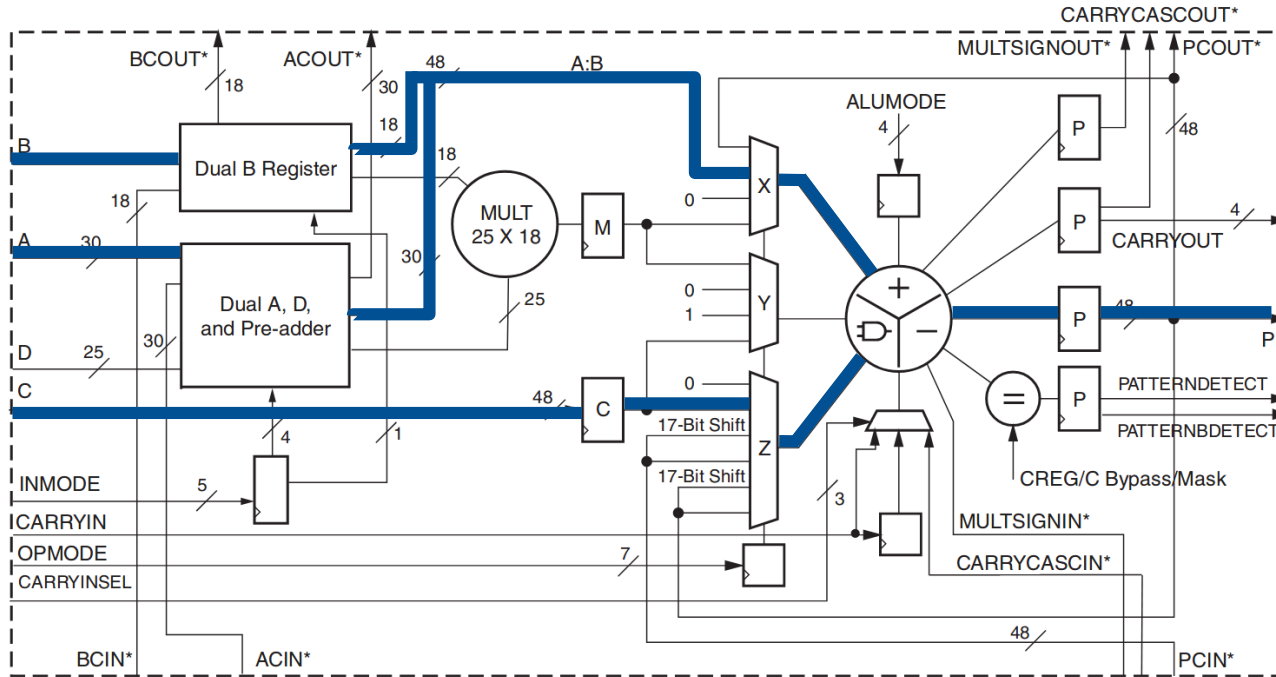
- Binary Mode
- Fixed-Precision Mode
- Reconfiguration Signals

Using the same hardware resource for two distinct, critical BNN operations

*These signals are dedicated routing paths internal to the DSP48E1 column. They are not accessible via fabric routing resources.

OrthrusPE: Dual Modes

Runtime Reconfigurability



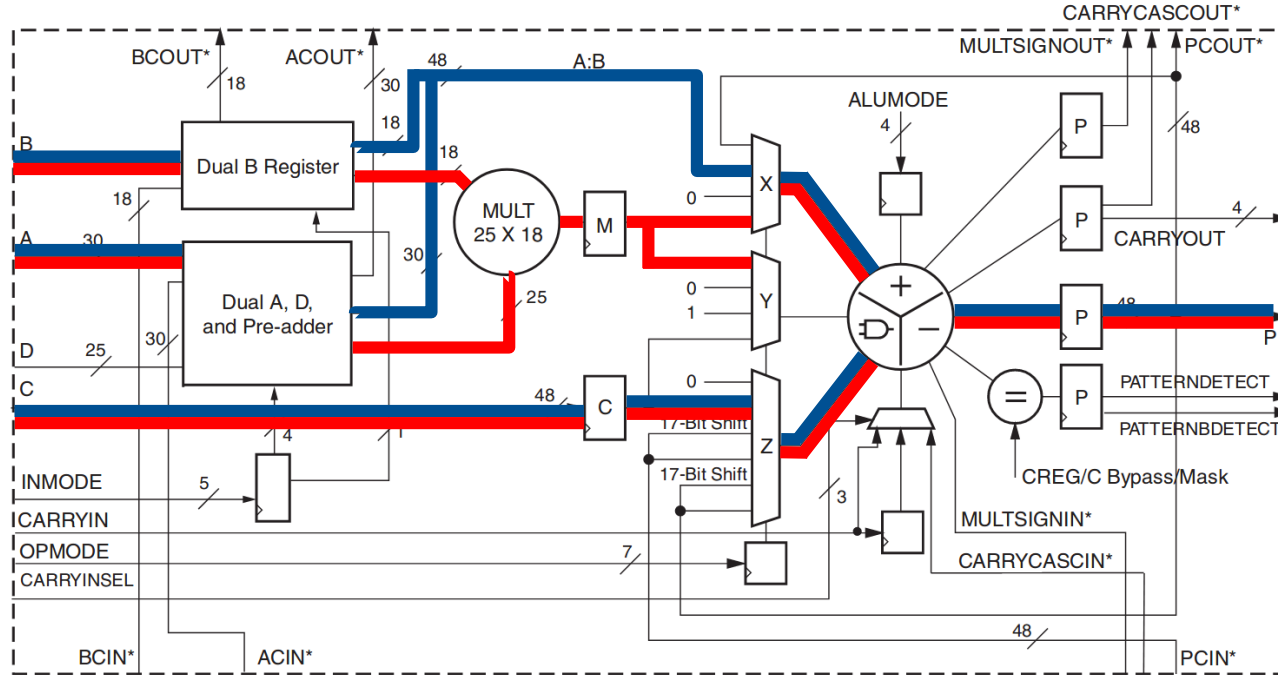
- Binary Mode
- Fixed-Precision Mode
- Reconfiguration Signals

Using the same hardware resource for two distinct, critical BNN operations

*These signals are dedicated routing paths internal to the DSP48E1 column. They are not accessible via fabric routing resources.

OrthrusPE: Dual Modes

Runtime Reconfigurability



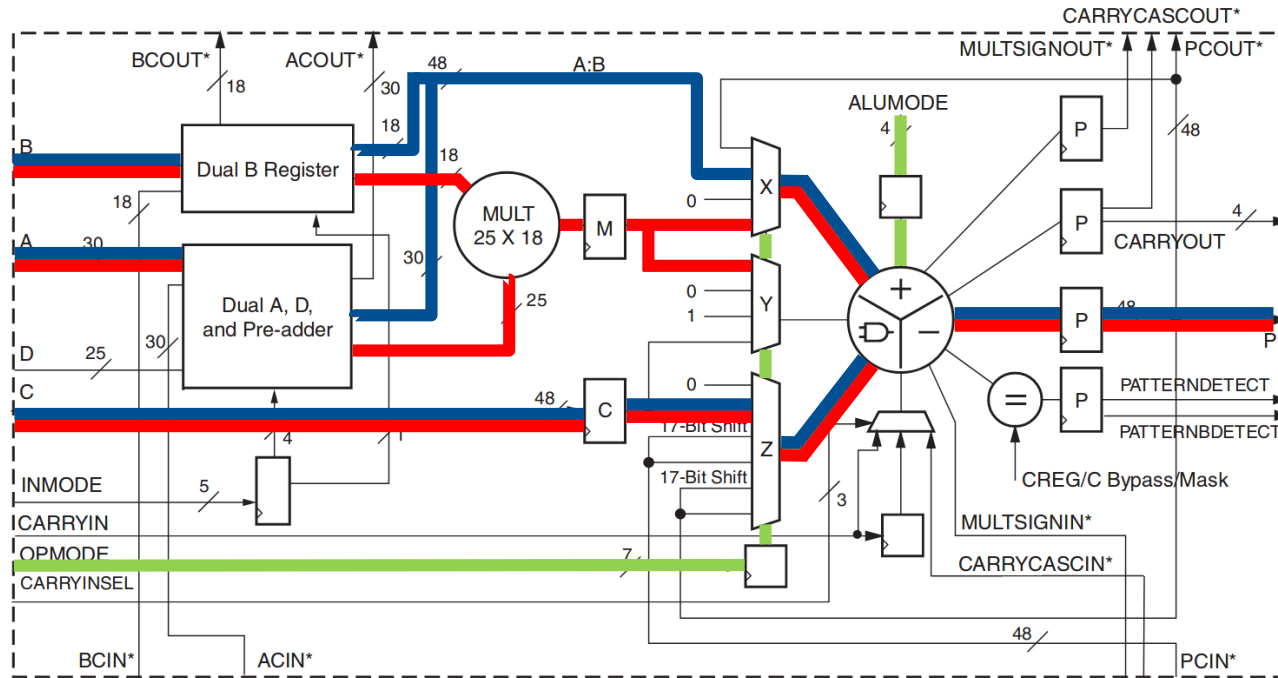
- Binary Mode
- Fixed-Precision Mode
- Reconfiguration Signals

Using the same hardware resource for two distinct, critical BNN operations

*These signals are dedicated routing paths internal to the DSP48E1 column. They are not accessible via fabric routing resources.

OrthrusPE: Dual Modes

Runtime Reconfigurability



*These signals are dedicated routing paths internal to the DSP48E1 column. They are not accessible via fabric routing resources.

- Binary Mode
- Fixed-Precision Mode
- Reconfiguration Signals

Using the same hardware resource for two distinct, critical BNN operations

Experimentation and Evaluation

Synthesized Four Throughput-Equivalent Configurations:

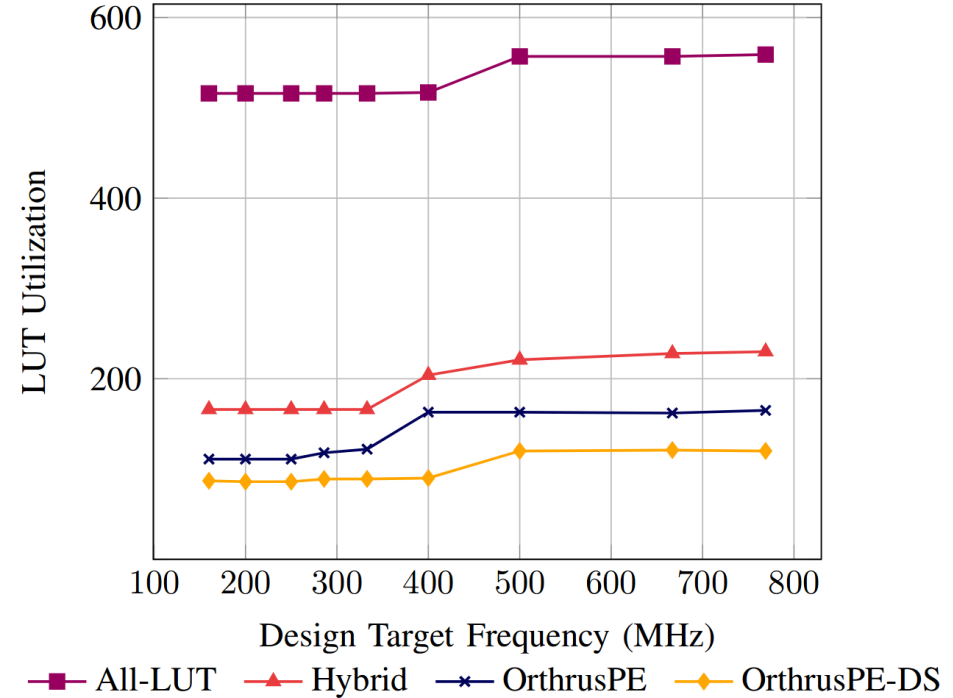
- OrthrusPE
- OrthrusPE-DS (Dual-Static): SIMD Binary Hadamard Products on Static DSP
- Hybrid (Common): Binary operations on LUTs, FP operations on DSP
- All-LUT: Execution restricted to LUTs

Experimentation and Evaluation

Resource Utilization

- OrthrusPE and OrthrusPE-DS are more resource efficient across all target accelerator frequencies.

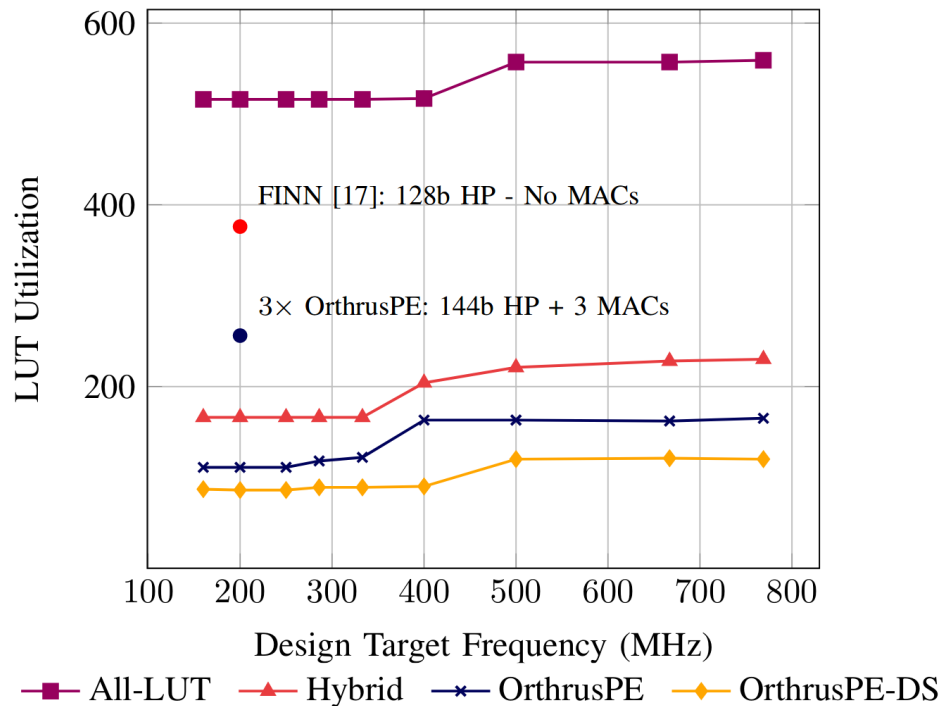
Implementation	F=770MHz			F=160MHz		
	LUTs	FF	DSP	LUTs	FF	DSP
All-LUT	559	160	0	516	160	0
Hybrid (Common)	230	253	1	166	253	1
OrthrusPE	165	210	1	111	210	1
OrthrusPE-DS	120	229	2	87	229	2



Experimentation and Evaluation

Resource Utilization

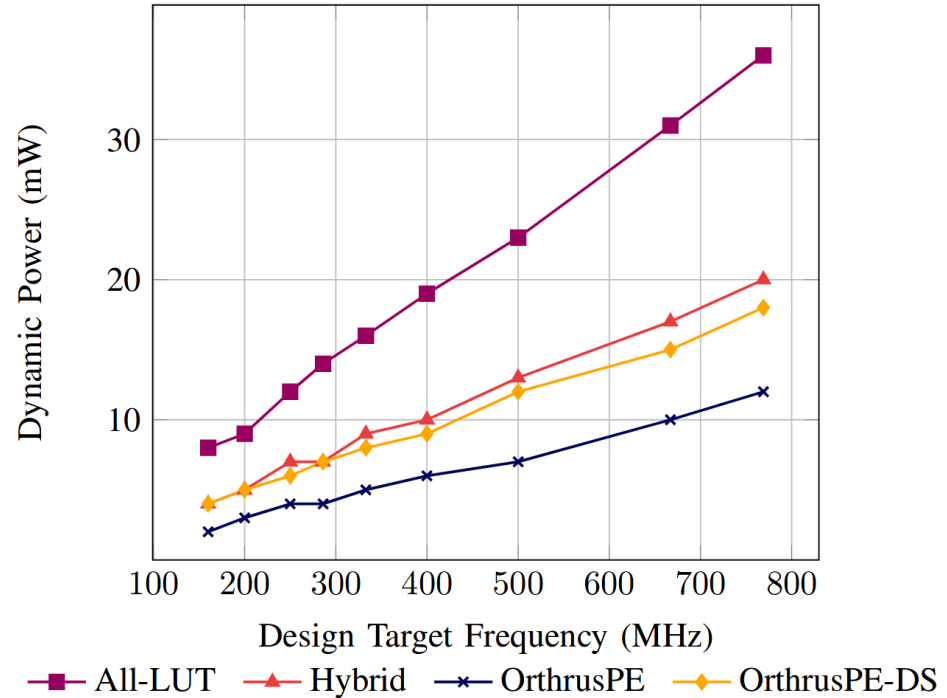
- OrthrusPE's closest FINN configuration @200MHz
 - 16 Extra Bit Accumulations
 - 3 MACs (through reconfigurability)
 - 32% fewer LUTs



Experimentation and Evaluation

Dynamic Power Estimation

- OrthrusPE more efficient across all frequencies
- Results scale as accelerators use 100-1000s of PEs



Conclusion and Future Work

- Accurate BNNs cannot be achieved without fixed-point operations and reliance on DSP blocks.
- OrthrusPE improves the efficiency of computation by executing both fixed-point and binary ops on FPGA hard blocks.
- Accurate BNNs solve many of the computation and memory challenges for deep neural network workloads on edge devices.

Thank you for your attention