Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

Results
0000000

# A multiplier-free RNS-based CNN accelerator exploiting bit-level sparsity

V. Sakellariou[1], V. Paliouras[2], I. Kouretas[2], H. Saleh[1], T. Stouraitis[1]

[1]Department of Electrical Engineering and Computer Science
Khalifa University
[2]Department of Electrical and Computer Engineering
University of Patras

30th IEEE International Symposium on Computer Arithmetic, Portland, Oregon, USA. September 4-6, 2023.

# Table of Contents

# Presentation Outline

# RNS Basics

- In RNS, integers are represented by their residues with respect to a modulus set: $\mathcal{B} = \{m_1, m_2, \dots, m_K\}$

## RNS Basics

- In RNS, integers are represented by their residues with respect to a modulus set: $\mathcal{B} = \{m_1, m_2, \ldots, m_K\}$

$$X \mapsto (x_1, x_2, \ldots, x_K), \ x_i = \langle X \rangle_{m_i} \tag{1}$$

- Multiplication and addition are done independently and in parallel in each channel

$$a \oplus b = \left( \langle a_1 \oplus b_1 \rangle_{m_1}, \langle a_2 \oplus b_2 \rangle_{m_2}, \ldots, \langle a_N \oplus b_N \rangle_{m_K} \right)$$

## RNS Basics

- In RNS, integers are represented by their residues with respect to a modulus set: $\mathcal{B} = \{m_1, m_2, \ldots, m_K\}$
$$X \mapsto (x_1, x_2, \ldots, x_K), \; x_i = \langle X \rangle_{m_i} \qquad (1)$$

- Multiplication and addition are done independently and in parallel in each channel
$$a \oplus b = \left( \langle a_1 \oplus b_1 \rangle_{m_1}, \langle a_2 \oplus b_2 \rangle_{m_2}, \ldots, \langle a_N \oplus b_N \rangle_{m_K} \right)$$

- Large number computations are decomposed into smaller

- Smaller critical path $\rightarrow$ higher frequencies or reduced power dissipation

- Efficient implementation of MAC $\rightarrow$ good candidate for use in CNNs

**Introduction**
○●○○○

Proposed PE (exploiting bit-level sparsity)
○○○○○○○○○

Overall CNN Architecture
○○○

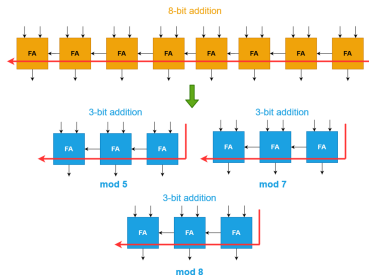Results
○○○○○○○

## RNS Basics

- In RNS, integers are represented by their residues with respect to a modulus set: $\mathcal{B} = \{m_1, m_2, \ldots, m_K\}$
$$X \mapsto (x_1, x_2, \ldots, x_K), \ x_i = \langle X \rangle_{m_i} \tag{1}$$

- Multiplication and addition are done independently and in parallel in each channel
$$a \oplus b = \left( \langle a_1 \oplus b_1 \rangle_{m_1}, \langle a_2 \oplus b_2 \rangle_{m_2}, \ldots, \langle a_N \oplus b_N \rangle_{m_K} \right)$$

- Large number computations are decomposed into smaller

- Smaller critical path → higher frequencies or reduced power dissipation

- Efficient implementation of MAC → good candidate for use in CNNs
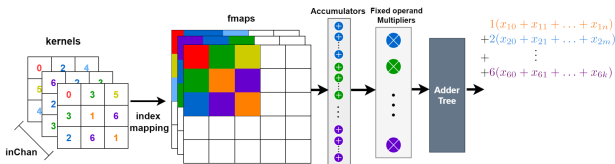
# Motivation: Multiplier-free PE



Figure: Modulo-7 Convolution

**Introduction**
○○○●○

Proposed PE (exploiting bit-level sparsity)
○○○○○○○○○

Overall CNN Architecture
○○○

Results
○○○○○○○

# Motivation: Multiplier-free PE
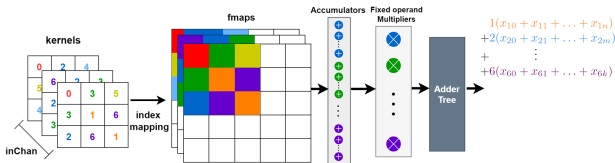


Figure: Modulo-7 Convolution



Figure: Direct implementation is the Multiplier-free
(MF) PE → 1 $\log_2 C$-bit adder and $C$ registers

**Introduction**
○○○●○

Proposed PE (exploiting bit-level sparsity)
○○○○○○○○○

Overall CNN Architecture
○○○

Results
○○○○○○○

# Motivation: Multiplier-free PE



Figure: Modulo-7 Convolution



Use many small word-length RNS channels

- small set of uniformly distributed weights
- increased number of common factors
- perform the additions first and then multiplications

Figure: Direct implementation is the Multiplier-free (MF) PE → 1 $\log_2 C$-bit adder and $C$ registers

Introduction
○○○●○

Proposed PE (exploiting bit-level sparsity)
○○○○○○○○○

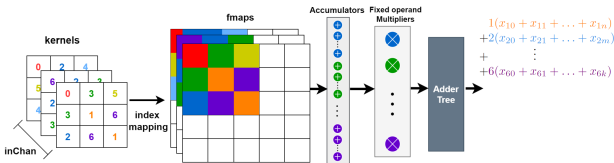Overall CNN Architecture
○○○

Results
○○○○○○○
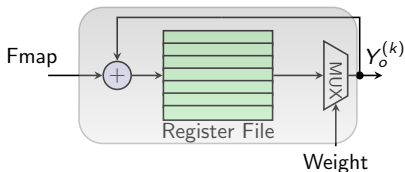
# Motivation: Multiplier-free PE



Figure: Modulo-7 Convolution



Figure: Direct implementation is the Multiplier-free (MF) PE → 1 $\log_2$ C-bit adder and C registers

Use many small word-length RNS channels

- small set of uniformly distributed weights
- increased number of common factors
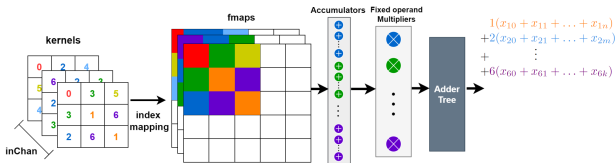- perform the additions first and then multiplications
- trivial/fixed operand mult. → simplified implementation

# Motivation: Multiplier-free Distributed PE



Figure: Distributed Multiplier-free (MF-D)
PE $\rightarrow \log_2 C$ adders and $\log_2 C$ registers.

**Introduction**
oooo

Proposed PE (exploiting bit-level sparsity)
ooooooooo

Overall CNN Architecture
ooo

Results
ooooooo

# Motivation: Multiplier-free Distributed PE



Figure: Distributed Multiplier-free (MF-D)
PE $\rightarrow \log_2 C$ adders and $\log_2 C$ registers.

### Examples

$W = 5 = 101 \rightarrow$
Fmap is added to Y0,Y2

# Motivation: Multiplier-free Distributed PE



Weight $\{w_2 w_1 w_0\}^{(k)}$ Fmap

### Examples

$W = 5 = 101 \rightarrow$
Fmap is added to Y0,Y2

Figure: Distributed Multiplier-free (MF-D)
PE $\rightarrow \log_2 C$ adders and $\log_2 C$ registers.

### Problem

**Under-utilization** of adders in case of 0's
Throughput: **1 Input / cycle**

# Presentation Outline

Introduction
oooo

Proposed PE (exploiting bit-level sparsity)
o●ooooooo

Overall CNN Architecture
ooo

Results
ooooooo

## Increasing effective throughput



- To avoid under-utilization,
  **load two inputs per cycle**

### Examples

W1 = 100, W2=010 →
Fmap1 is added to Y2, Fmap2 is
added to Y1

## Increasing effective throughput



Weight $\{w_2 w_1 w_0\}^{(k)}$ Fmap

$w_0$   $w_1$   $w_2$   MUX

$Y_0^{(k)}$   $Y_1^{(k)}$   $Y_2^{(k)}$

- To avoid under-utilization, **load two inputs per cycle**
- Conflict if two 1's in the same digital position

### Examples

W1 = 100, W2=010 → Fmap1 is added to Y2, Fmap2 is added to Y1



Figure: $P_c$ of an MF-D PE channel for various $sp$ vs $l$.

$sp = \frac{\text{num. of zero weights}}{\text{total num. of weights}}$

- Conflict probability $P_c$ quickly increase as channel word-length l increases

## Stack-based distributed PE (MF-D-S PE)

- Use a stack-like storage element to store conflicted inputs and **minimize stalls**



Figure: MF-D-S PE . The PE receives two FMAPs $F_0^{(k)}$, $F_1^{(k)}$, one select signal ($sel$) for each of the adders and a write-enable signal for each of the stacks.

# Stack-based distributed PE (MF-D-S PE)

- Use a stack-like storage element to store conflicted inputs and **minimize stalls**
- In case of 2 '0's, adders **process previously conflicted inputs**



Figure: MF-D-S PE . The PE receives two FMAPs $F_0^{(k)}$, $F_1^{(k)}$, one select signal ($sel$) for each of the adders and a write-enable signal for each of the stacks.

# Stack-based distributed PE (MF-D-S PE)

- Use a stack-like storage element to store conflicted inputs and **minimize stalls**
- In case of 2 '0's, adders **process previously conflicted inputs**
- Higher sparsity translates into higher throughput



Figure: MF-D-S PE . The PE receives two FMAPs $F_0^{(k)}$, $F_1^{(k)}$, one select signal (*sel*) for each of the adders and a write-enable signal for each of the stacks.

## Stack-based distributed PE (MF-D-S PE)

- Use a stack-like storage element to store conflicted inputs and **minimize stalls**
- In case of 2 '0's, adders **process previously conflicted inputs**
- Higher sparsity translates into higher throughput
- Channel size also affects $P_c$ $\rightarrow$ RNS is naturally more suitable representation: lower stall probability and less hardware resources



Figure: MF-D-S PE . The PE receives two FMAPs $F_0^{(k)}$, $F_1^{(k)}$, one select signal (*sel*) for each of the adders and a write-enable signal for each of the stacks.

# Stack-based distributed PE (MF-D-S PE)

## Stack-based distributed PE (MF-D-S PE)

- Model as Markov chain to calculate probability of stalls $P_c$
- State represents num. of elements in stack

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000●00000

Overall CNN Architecture
000

Results
0000000

## Stack-based distributed PE (MF-D-S PE)

- Model as Markov chain to calculate probability of stalls $P_c$
- State represents num. of elements in stack

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000●00000

Overall CNN Architecture
000

Results
0000000

## Stack-based distributed PE (MF-D-S PE)

- Model as Markov chain to calculate probability of stalls $P_c$
- State represents num. of elements in stack



- $\boldsymbol{N} = \sum_{k=0}^{\infty} \boldsymbol{Q}^k = (\boldsymbol{I} - \boldsymbol{Q})^{-1}$, where Q is the transition matrix
- stalls $N_s = \frac{(S+1)N}{T_a}$, $T_a = \boldsymbol{N}\boldsymbol{1}$
- throughput $T = \frac{2N}{N+N_s}$

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

Results
0000000

# Stack-based distributed PE (MF-D-S PE)

- Model as Markov chain to calculate probability of stalls $P_c$
- State represents num. of elements in stack



- $\boldsymbol{N} = \sum_{k=0}^{\infty} \boldsymbol{Q}^k = (\boldsymbol{I} - \boldsymbol{Q})^{-1}$, where Q is the transition matrix
- stalls $N_s = \frac{(S+1)N}{T_a}$, $T_a = \boldsymbol{N1}$
- throughput $T = \frac{2N}{N+N_s}$

- $P_c$ and thus throughput depend on channel size $l$, sparsity $sp$ and stack size $S$.

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
0000●0000

Overall CNN Architecture
000

Results
0000000

# Stack-based distributed PE (MF-D-S PE)

- Model as Markov chain to calculate probability of stalls $P_c$
- State represents num. of elements in stack



- $\boldsymbol{N} = \sum_{k=0}^{\infty} \boldsymbol{Q}^k = (\boldsymbol{I} - \boldsymbol{Q})^{-1}$, where Q is the transition matrix
- stalls $N_s = \frac{(S+1)N}{T_a}$, $T_a = \boldsymbol{N}\mathbf{1}$
- throughput $T = \frac{2N}{N+N_s}$

- $P_c$ and thus throughput depend on channel size $l$, sparsity $sp$ and stack size $S$.



### Examples

For a 5-bit channel with $sp = 0$:
T=1.12 (MF-D) $\rightarrow$
T=1.32 (MF-D-S, S=1)

11 / 25

## CSD encoding for reduced chance of stalls

- $P_c$ depends on the distribution of '1's and '0's in the input weight vector.

## CSD encoding for reduced chance of stalls

- $P_c$ depends on the distribution of '1's and '0's in the input weight vector.
- An encoding that reduces the probability of conflicts (two '1's at the same position) would increase its throughput.

## CSD encoding for reduced chance of stalls

- $P_c$ depends on the distribution of '1's and '0's in the input weight vector.
- An encoding that reduces the probability of conflicts (two '1's at the same position) would increase its throughput.
- The number of the non-zero elements of the weight representation can be minimized through the use of Canonical Signed Digit (CSD).

## CSD encoding for reduced chance of stalls

- $P_c$ depends on the distribution of '1's and '0's in the input weight vector.
- An encoding that reduces the probability of conflicts (two '1's at the same position) would increase its throughput.
- The number of the non-zero elements of the weight representation can be minimized through the use of Canonical Signed Digit (CSD).
- In CSD, the value of each digit can be either 0, 1, or -1

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000●000

Overall CNN Architecture
000

Results
0000000

# CSD encoding for reduced chance of stalls

- $P_c$ depends on the distribution of '1's and '0's in the input weight vector.
- An encoding that reduces the probability of conflicts (two '1's at the same position) would increase its throughput.
- The number of the non-zero elements of the weight representation can be minimized through the use of Canonical Signed Digit (CSD).
- In CSD, the value of each digit can be either 0, 1, or -1
- Overhead: XOR gates needed to support subtractions

### Examples

$01110 \rightarrow 1\ 0\ 0\ (-1)\ 0$
$01101 \rightarrow 0\ 1\ 1\quad 0\ 1$    stall is avoided

Introduction
oooo

Proposed PE (exploiting bit-level sparsity)
ooooooo●oo

Overall CNN Architecture
ooo

Results
ooooooo

# CSD encoding for reduced chance of stalls



Figure: Conflict probability $P_c$ of an MF-D PE channel for various sparsity levels $sp$ and word lengths $l$. Blue color denotes binary encoding while red denotes CSD encoding.

# Optimal encoding of pairs of weights

- CSD doesn't consider relative digit positions in weight pairs

## Optimal encoding of pairs of weights

- CSD doesn't consider relative digit positions in weight pairs
- Jointly encode pairs of weights to minimize conflicts/stalls

$T(W)$ denotes the position of the trailing non-zero digit of $W$

Introduction
oooo

Proposed PE (exploiting bit-level sparsity)
oooooooo●o

Overall CNN Architecture
ooo

Results
ooooooo

## Optimal encoding of pairs of weights

- CSD doesn't consider relative digit positions in weight pairs
- Jointly encode pairs of weights to minimize conflicts/stalls

$T(W)$ denotes the position of the trailing non-zero digit of $W$

### Lemma

*A signed-digit encoding (not necessarily canonical)*
*$E(W_a, W_b) : (W_a, W_b) \mapsto (\hat{W}_a, \hat{W}_b)$, such that $C(\hat{W}_a, \hat{W}_b) = 0$,*
*for two non-zero weights, exists, if and only if:*

$$T(W_a) \neq T(W_b). \tag{2}$$

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

Results
0000000

## Optimal encoding of pairs of weights

- CSD doesn't consider relative digit positions in weight pairs
- Jointly encode pairs of weights to minimize conflicts/stalls

$T(W)$ denotes the position of the trailing non-zero digit of $W$

### Lemma

*A signed-digit encoding (not necessarily canonical)*
$E(W_a, W_b) : (W_a, W_b) \mapsto (\hat{W}_a, \hat{W}_b)$, such that $C(\hat{W}_a, \hat{W}_b) = 0$,
*for two non-zero weights, exists, if and only if:*
$$T(W_a) \neq T(W_b). \tag{2}$$

### Lemma

*The probability $P_c$ that no zero-conflict encoding for two weights*
$W_a, W_b$ exists, i.e., $T(W_a) = T(W_b)$, is given by
$$P_c = \frac{1}{3} - \frac{1}{3 \cdot 4^n}. \tag{3}$$

## Encoder scheme comparison

Table: Probability of a conflict $P_c$ for $\mod 2^n$, $\mod(2^n - 1)$ and $\mod(2^n + 1)$ and channel sizes $n = 4$ and $n = 5$

| $n$ | $\mod 2^n$ | | | | | $\mod(2^n - 1)$ | | | | | $\mod(2^n + 1)^{\ddagger}$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | no enc | CSD | CSD/bin | combined | Opt. | no enc | CSD$^{\dagger}$ | CSD$^{\dagger}$/bin | combined$^{\dagger}$ | Opt.$^{\dagger}$ | no enc | CSD$^{\dagger}$ | CSD$^{\dagger}$/bin | combined$^{\dagger}$ | Opt.$^{\ddagger}$ |
| 4 | 0.683 | 0.457 | 0.425 | 0.425 | 0.332 | 0.648 | 0.515 | 0.462 | 0.328 | 0.328 | 0.605 | 0.439 | 0.401 | 0.335 | 0.308 |
| 5 | 0.762 | 0.516 | 0.476 | 0.476 | 0.333 | 0.749 | 0.578 | 0.521 | 0.413 | 0.339 | 0.717 | 0.529 | 0.482 | 0.413 | 0.327 |

no enc: both weights in binary  CSD: independent CSD encoding of weights  $^{\dagger}$: no re-encoding after EAC  Opt.: Optimal
CSD/bin: $a$, CSD; $b$, either binary or CSD  combined: $a$, CSD, $b$, binary, $b_{CSD}$, or $b'_{CSD}$  $^{\ddagger}$: diminished-1 representation

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
00000000●

Overall CNN Architecture
000

Results
0000000

## Encoder scheme comparison

Table: Probability of a conflict $P_c$ for mod $2^n$, mod$(2^n - 1)$ and mod$(2^n + 1)$ and channel sizes $n = 4$ and $n = 5$

| $n$ | mod $2^n$ | | | | | mod $(2^n - 1)$ | | | | | mod $(2^n + 1)^{\ddagger}$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | no enc | CSD | CSD/bin | combined | Opt. | no enc | CSD$^{\dagger}$ | CSD$^{\dagger}$/bin | combined$^{\dagger}$ | Opt.$^{\dagger}$ | no enc | CSD$^{\dagger}$ | CSD$^{\dagger}$/bin | combined$^{\dagger}$ | Opt.$^{\ddagger}$ |
| 4 | 0.683 | 0.457 | 0.425 | 0.425 | 0.332 | 0.648 | 0.515 | 0.462 | 0.328 | 0.328 | 0.605 | 0.439 | 0.401 | 0.335 | 0.308 |
| 5 | 0.762 | 0.516 | 0.476 | 0.476 | 0.333 | 0.749 | 0.578 | 0.521 | 0.413 | 0.339 | 0.717 | 0.529 | 0.482 | 0.413 | 0.327 |

no enc: both weights in binary  CSD: independent CSD encoding of weights  $^{\dagger}$: no re-encoding after EAC  Opt.: Optimal
CSD/bin: $a$, CSD; $b$, either binary or CSD  combined: $a$, CSD, $b$, binary, $b_{CSD}$, or $b'_{CSD}$  $^{\ddagger}$: diminished-1 representation

Table: Speedup ($\times$) for a mod-32 channel vs stack size $S$ and $sp$

| $sp$ | $S = 0$ | | | $S = 1$ | | | $S = 2$ | ZS$^{\dagger}$ |
|---|---|---|---|---|---|---|---|---|
| | Bin. | CSD | Opt. | Bin. | CSD | Opt. | Bin. | |
| 0 | 1.13 | 1.32 | 1.50 | 1.32 | 1.65 | **1.74** | 1.37 | 1 |
| 0.1 | 1.23 | 1.43 | 1.57 | 1.43 | 1.75 | **1.84** | 1.53 | 1.11 |
| 0.2 | 1.34 | 1.52 | 1.64 | 1.55 | 1.83 | **1.91** | 1.72 | 1.25 |
| 0.3 | 1.45 | 1.61 | 1.71 | 1.68 | 1.89 | **1.95** | 1.85 | 1.42 |
| 0.4 | 1.57 | 1.70 | 1.78 | 1.79 | 1.94 | **1.98** | 1.94 | 1.66 |
| 0.5 | 1.67 | 1.78 | 1.84 | 1.89 | 1.97 | **1.99** | 1.98 | 2 |

$^{\dagger}$: Zero skipping

Introduction
oooo
Proposed PE (exploiting bit-level sparsity)
ooooooooo●
Overall CNN Architecture
ooo
Results
ooooooo

## Encoder scheme comparison

Table: Probability of a conflict $P_c$ for mod $2^n$, mod($2^n - 1$) and mod($2^n + 1$) and channel sizes $n = 4$ and $n = 5$

| $n$ | mod $2^n$ | | | | | mod $(2^n - 1)$ | | | | | mod $(2^n + 1)^\ddagger$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | no enc | CSD | CSD/bin | combined | Opt. | no enc | CSD$^\dagger$ | CSD$^\dagger$/bin | combined$^\dagger$ | Opt.$^\dagger$ | no enc | CSD$^\dagger$ | CSD$^\dagger$/bin | combined$^\dagger$ | Opt.$^\ddagger$ |
| 4 | 0.683 | 0.457 | 0.425 | 0.425 | 0.332 | 0.648 | 0.515 | 0.462 | 0.328 | 0.328 | 0.605 | 0.439 | 0.401 | 0.335 | 0.308 |
| 5 | 0.762 | 0.516 | 0.476 | 0.476 | 0.333 | 0.749 | 0.578 | 0.521 | 0.413 | 0.339 | 0.717 | 0.529 | 0.482 | 0.413 | 0.327 |

no enc: both weights in binary CSD: independent CSD encoding of weights $^\dagger$: no re-encoding after EAC Opt.: Optimal
CSD/bin: $a$, CSD; $b$, either binary or CSD combined: $a$, CSD, $b$, binary, $b_{CSD}$, or $b'_{CSD}$ $^\ddagger$: diminished-1 representation

Table: Speedup ($\times$) for a mod-32 channel vs stack size $S$ and $sp$

| $sp$ | $S = 0$ | | | $S = 1$ | | | $S = 2$ | ZS$^\dagger$ |
|---|---|---|---|---|---|---|---|---|
| | Bin. | CSD | Opt. | Bin. | CSD | Opt. | Bin. | |
| 0 | 1.13 | 1.32 | 1.50 | 1.32 | 1.65 | **1.74** | 1.37 | 1 |
| 0.1 | 1.23 | 1.43 | 1.57 | 1.43 | 1.75 | **1.84** | 1.53 | 1.11 |
| 0.2 | 1.34 | 1.52 | 1.64 | 1.55 | 1.83 | **1.91** | 1.72 | 1.25 |
| 0.3 | 1.45 | 1.61 | 1.71 | 1.68 | 1.89 | **1.95** | 1.85 | 1.42 |
| 0.4 | 1.57 | 1.70 | 1.78 | 1.79 | 1.94 | **1.98** | 1.94 | 1.66 |
| 0.5 | 1.67 | 1.78 | 1.84 | 1.89 | 1.97 | **1.99** | 1.98 | 2 |

$^\dagger$: Zero skipping

Table: Hardware complexity (number of gates) of the different encoders

| Encoding | mod $2^n$ | mod $(2^n - 1)$ |
|---|---|---|
| CSD | 21 | 29 |
| CSD/bin | 50 | 71 |
| Combined | 80 | 112 |
| Optimal | 495 | 515 |

# Presentation Outline

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
0●0

Results
0000000

## Processing Core

- PEs that share the same weight are grouped together

## Processing Core

- PEs that share the same weight are grouped together
- Since the weight determines their operation, all PEs work in a synchronized manner

1. No additional storage elements/scheduling to support the different processing rates of PEs

2. Control logic is amortized over 16 PEs

## Processing Core

- PEs that share the same weight are grouped together
- Since the weight determines their operation, all PEs work in a synchronized manner



1. No additional storage elements/scheduling to support the different processing rates of PEs

2. Control logic is amortized over 16 PEs

3. The CSD/optimal encoder is also amortized; its overhead becomes negligible

Core overhead:

- shift-add units
- one additional base extension unit
- larger shift register array
- weight FIFO buffers

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
00●

Results
0000000

## System-level Architecture and Dataflow



- Data is represented in a reduced RNS base of 8 or 12 bits

- 16 RNS Cores

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
00●

Results
0000000

## System-level Architecture and Dataflow



- Data is represented in a reduced RNS base of 8 or 12 bits
- FMEM: stores intermediate layer results in the reduced RNS base, WMEM: stores weights

- 16 RNS Cores
- $\mathcal{B} = (5, 7, 31, 32, 33)$

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
00●

Results
0000000

## System-level Architecture and Dataflow



- Data is represented in a reduced RNS base of 8 or 12 bits
- FMEM: stores intermediate layer results in the reduced RNS base, WMEM: stores weights
- Base Extension: extend data to the full RNS base (22 bits)

- 16 RNS Cores
- $\mathcal{B} = (5, 7, 31, 32, 33)$

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
00●

Results
0000000

## System-level Architecture and Dataflow



- Data is represented in a reduced RNS base of 8 or 12 bits
- FMEM: stores intermediate layer results in the reduced RNS base, WMEM: stores weights
- Base Extension: extend data to the full RNS base (22 bits)
- A 2-D input block is fetched and broadcast to all cores

- 16 RNS Cores
- $\mathcal{B} = (5, 7, 31, 32, 33)$

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
00●

Results
0000000

## System-level Architecture and Dataflow



- Data is represented in a reduced RNS base of 8 or 12 bits
- FMEM: stores intermediate layer results in the reduced RNS base, WMEM: stores weights
- Base Extension: extend data to the full RNS base (22 bits)
- A 2-D input block is fetched and broadcast to all cores
- One Activation/Scaling Unit per two cores

- 16 RNS Cores
- $\mathcal{B} = (5, 7, 31, 32, 33)$

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

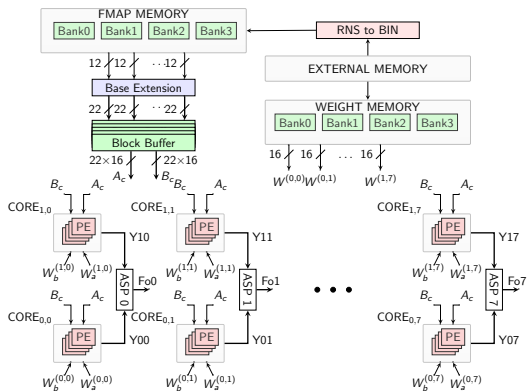Overall CNN Architecture
00●

Results
0000000

## System-level Architecture and Dataflow



- 16 RNS Cores
- $\mathcal{B} = (5, 7, 31, 32, 33)$

- Data is represented in a reduced RNS base of 8 or 12 bits
- FMEM: stores intermediate layer results in the reduced RNS base, WMEM: stores weights
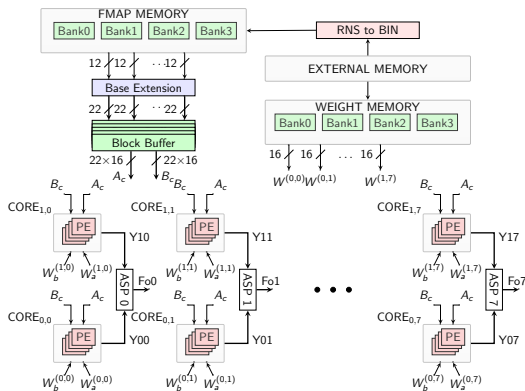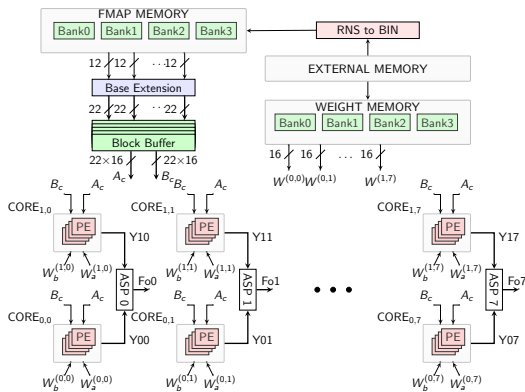- Base Extension: extend data to the full RNS base (22 bits)
- A 2-D input block is fetched and broadcast to all cores
- One Activation/Scaling Unit per two cores
- Block Buffer to decouple access to FMEM

# Presentation Outline

# Sparsity and throughput exploration on CNN benchmarks

Table: Sparsity of CNN benchmarks and expected speedup

| Network (8-bit quant.) | Weight Sparsity (%) | MF-D-S Theoretical Speedup* | MF-D-S Simulation Speedup* | ZS[†] |
|---|---|---|---|---|
| VGG19 | 42 | $1.82\times/1.95\times$ | $1.79\times/1.88\times$ | $1.72\times$ |
| ResNet50 | 12 | $1.42\times/1.77\times$ | $1.47\times/1.76\times$ | $1.13\times$ |
| Yolo3 | 40 | $1.80\times/1.94\times$ | $1.80\times/1.88\times$ | $1.67\times$ |
| InceptionV3 | 8 | $1.40\times/1.73\times$ | $1.43\times/1.73\times$ | $1.08\times$ |
| MobileNet | 7 | $1.39\times/1.73\times$ | $1.41\times/1.73\times$ | $1.03\times$ |

*: In entries of the form $x/y$, $x$ refers to binary encoding and
$y$ refers to CSD encoding ($S = 1$)

## Sparsity and throughput exploration on CNN benchmarks

Table: Sparsity of CNN benchmarks and expected speedup

| Network (8-bit quant.) | Weight Sparsity (%) | MF-D-S Theoretical Speedup* | MF-D-S Simulation Speedup* | ZS[†] |
|---|---|---|---|---|
| VGG19 | 42 | $1.82\times/1.95\times$ | $1.79\times/1.88\times$ | $1.72\times$ |
| ResNet50 | 12 | $1.42\times/1.77\times$ | $1.47\times/1.76\times$ | $1.13\times$ |
| Yolo3 | 40 | $1.80\times/1.94\times$ | $1.80\times/1.88\times$ | $1.67\times$ |
| InceptionV3 | 8 | $1.40\times/1.73\times$ | $1.43\times/1.73\times$ | $1.08\times$ |
| MobileNet | 7 | $1.39\times/1.73\times$ | $1.41\times/1.73\times$ | $1.03\times$ |

*: In entries of the form $x/y$, $x$ refers to binary encoding and
$y$ refers to CSD encoding $(S=1)$

- Proposed method results in $1.73\times$ to $1.88\times$ speedup

## Sparsity and throughput exploration on CNN benchmarks

Table: Sparsity of CNN benchmarks and expected speedup

| Network (8-bit quant.) | Weight Sparsity (%) | MF-D-S Theoretical Speedup* | MF-D-S Simulation Speedup* | ZS[†] |
|---|---|---|---|---|
| VGG19 | 42 | $1.82\times/1.95\times$ | $1.79\times/1.88\times$ | $1.72\times$ |
| ResNet50 | 12 | $1.42\times/1.77\times$ | $1.47\times/1.76\times$ | $1.13\times$ |
| Yolo3 | 40 | $1.80\times/1.94\times$ | $1.80\times/1.88\times$ | $1.67\times$ |
| InceptionV3 | 8 | $1.40\times/1.73\times$ | $1.43\times/1.73\times$ | $1.08\times$ |
| MobileNet | 7 | $1.39\times/1.73\times$ | $1.41\times/1.73\times$ | $1.03\times$ |

*: In entries of the form $x/y$, $x$ refers to binary encoding and $y$ refers to CSD encoding ($S = 1$)

- Proposed method results in $1.73\times$ to $1.88\times$ speedup

### Takeaway

Unlike other sparse processing CNN architectures that rely on zero-skipping and require high sparsity levels to become efficient, the proposed method achieves gains with zero word-level sparsity (exploits bit-level sparsity)

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

Results
0000000

## Single PE comparisons

Table: PE comparison for various target clock periods (0.5 V)

| $T_{clk}$ | Area ($\mu m^2$) | | | | Power ($\mu W$) | | | |
|---|---|---|---|---|---|---|---|---|
| (ns) | BNS | RNS | MF-D-S$^\dagger$ | | BNS | RNS | MF-D-S$^\dagger$ | |
| | | | $S = 0$ | $S = 1$ | | | $S = 0$ | $S = 1$ |
| 0.8 | - | - | 420 | 592/625 | - | - | 137 | 153/158 |
| 0.9 | - | 340 | 400 | 580/610 | - | 107 | 125 | 131/140 |
| 1.0 | - | 334 | 391 | 571/585 | - | 93 | 110 | 116/130 |
| 1.1 | 350 | 329 | 384 | 556/581 | 128 | 84 | 102 | 106/112 |
| 1.2 | 336 | 308 | 379 | 542/581 | 110 | 74 | 83 | 99/103 |

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

Results
0000000

## Single PE comparisons

Table: PE comparison for various target clock periods (0.5 V)

| $T_{clk}$ | Area ($\mu m^2$) | | | | Power ($\mu W$) | | | |
|---|---|---|---|---|---|---|---|---|
| (ns) | BNS | RNS | MF-D-S[†] | | BNS | RNS | MF-D-S[†] | |
| | | | $S=0$ | $S=1$ | | | $S=0$ | $S=1$ |
| 0.8 | - | - | 420 | 592/625 | - | - | 137 | 153/158 |
| 0.9 | - | 340 | 400 | 580/610 | - | 107 | 125 | 131/140 |
| 1.0 | - | 334 | 391 | 571/585 | - | 93 | 110 | 116/130 |
| 1.1 | 350 | 329 | 384 | 556/581 | 128 | 84 | 102 | 106/112 |
| 1.2 | 336 | 308 | 379 | 542/581 | 110 | 74 | 83 | 99/103 |

- $1.85\times$ and $1.54\times$ more energy efficient processing compared to binary and conventional RNS

- The MF-D-S (S=1) PE achieves higher energy efficiency gains
  $\frac{P_{RNS}}{P_{MF-D-S}} \times speedup$ as the clock period becomes smaller and sparsity increases.

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

**Results**
0000000

## Single PE comparisons

Table: PE comparison for various target clock periods (0.5 V)

| $T_{clk}$ | Area ($\mu m^2$) | | | | Power ($\mu W$) | | | |
|---|---|---|---|---|---|---|---|---|
| (ns) | BNS | RNS | MF-D-S[†] | | BNS | RNS | MF-D-S[†] | |
| | | | $S=0$ | $S=1$ | | | $S=0$ | $S=1$ |
| 0.8 | - | - | 420 | 592/625 | - | - | 137 | 153/158 |
| 0.9 | - | 340 | 400 | 580/610 | - | 107 | 125 | 131/140 |
| 1.0 | - | 334 | 391 | 571/585 | - | 93 | 110 | 116/130 |
| 1.1 | 350 | 329 | 384 | 556/581 | 128 | 84 | 102 | 106/112 |
| 1.2 | 336 | 308 | 379 | 542/581 | 110 | 74 | 83 | 99/103 |

- $1.85\times$ and $1.54\times$ more energy efficient processing compared to binary and conventional RNS

- The MF-D-S ($S=1$) PE achieves higher energy efficiency gains
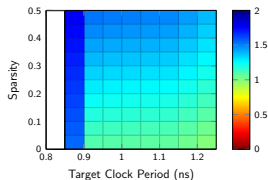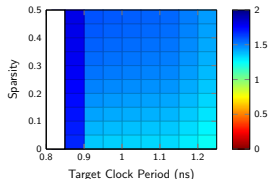  $\frac{P_{RNS}}{P_{MF-D-S}} \times speedup$ as the clock period becomes smaller and sparsity increases.

Energy Efficiency Gains (0.5 V - binary)



Energy Efficiency Gains (0.5 V - CSD)



- CSD encoding can further increase energy efficiency.

## Core-level comparisons

Table: Area and power breakdown of the various components

| Units | Components | Area[†] (µm²) | | | Power[†] (µW) | | |
|---|---|---|---|---|---|---|---|
| | | **BNS** | **RNS** | **MF-D-S** | **BNS** | **RNS** | **MF-D-S** |
| PE | | 350/277 | 329/279 | 556/529 | 128/220 | 93/140 | 120/229 |
| Shift-add Unit | | - | - | 160/110 | - | - | 20/25 |
| FMAP Base Ext. | | - | 250/190 | 250/190 | - | 61/86 | 61/86 |
| Weight Base Ext. | | - | 107/86 | 107/86 | - | 25/36 | 25/36 |
| ASP Unit | Scaling + ReLU + Pooling | 622/600 | 5761/5450 | 5761/5450 | 118/202 | 384/615 | 384/615 |
| **Core** | 4×4 PE array + BE (+shift-add unit) | **6.27/5.3**×10³ | **6.30/5.43**×10³ | **10.35/9.23**×10³ | **1.98/3.50**×10³ | **1.63/2.73**×10³ | **2.04/3.36**×10³ |

†: In entries of the form $x/y$, $x$ refers to a supply voltage of 0.5 V and $y$ refers to 0.65 V

Introduction
0000
Proposed PE (exploiting bit-level sparsity)
000000000
Overall CNN Architecture
000
Results
0000000

## Core-level comparisons

Table: Area and power breakdown of the various components

| Units | Components | Area$^\dagger$ ($\mu m^2$) | | | Power$^\dagger$ ($\mu W$) | | |
|---|---|---|---|---|---|---|---|
| | | BNS | RNS | MF-D-S | BNS | RNS | MF-D-S |
| PE | | 350/277 | 329/279 | 556/529 | 128/220 | 93/140 | 120/229 |
| Shift-add Unit | | - | - | 160/110 | - | - | 20/25 |
| FMAP Base Ext. | | - | 250/190 | 250/190 | - | 61/86 | 61/86 |
| Weight Base Ext. | | - | 107/86 | 107/86 | - | 25/36 | 25/36 |
| ASP Unit | Scaling + ReLU + Pooling | 622/600 | 5761/5450 | 5761/5450 | 118/202 | 384/615 | 384/615 |
| Core | 4×4 PE array + BE (+shift-add unit) | $6.27/5.3 \times 10^3$ | $6.30/5.43 \times 10^3$ | $10.35/9.23 \times 10^3$ | $1.98/3.50 \times 10^3$ | $1.63/2.73 \times 10^3$ | $2.04/3.36 \times 10^3$ |

$^\dagger$: In entries of the form $x/y$, $x$ refers to a supply voltage of 0.5 V and $y$ refers to 0.65 V

- Area overhead is completely compensated by the increased processing rate
- Power efficiency can be increased by 31%–54% (57%–85%) depending on the sparsity of the weights, compared to conventional RNS (BNS)

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

Results
0000●00

# System-level comparisons

### Table: Comparisons to state-of-the-art implementations

| | Eyerissv2[1] | ISSCC'20[2] | ISSCC'22[3] | RNSDNN[4] | This work RNS | **This work MF-D-S** |
|---|---|---|---|---|---|---|
| Process | 65 nm | 7 nm | 65 nm | 45 nm | 22 nm | **22 nm** |
| Supply voltage (V) | N/G | 0.575–0.825 | 1 | 1 | 0.65 | **0.65** |
| Frequency | 200 MHz | 290–880 MHz | 400 MHz | 1.2 GHz | 1 GHz | **1 GHz** |
| On-chip Memory (KB) | 246 | 2176 | 150 KB | N/G | 448 | **448** |
| Bit Precision (FMAP,wgt) | 8 | 8 | 8 | 16,8 | 12,8 | **12,8** |
| Network | AlexNet | MobileNet-v1 | VGG16 | VGG16 | VGG16 | **VGG16** |
| Performance (GOPS[†]) | 153.6 | 3604 | N/G | 134 | 220 | **364** |
| Area (10^6 Gates) | 2.69 | N/G | N/G | 2.18 | 4.74 | **5.25** |
| Area (mm²) | N/G | 3.04 | 4.47 | N/G | 0.94 | **1.04** |
| Area Eff.(GOPS[†]/10^6 Gates) | 57.1 | N/G | N/G | 62 | 46.4 | **69.3** |
| Power Eff. (TOPS[†]/W) | 0.253 - 0.962[‡] | 3.28 - 6.66[‡] | 1.82 | 0.223 | 1.74/2.36* | **1.98/3.12*** |

*: full system/on-chip power cons.       [†]1 OP = 1 MAC    [‡]: for dense - sparse network    N/G: not given

Introduction
0000
Proposed PE (exploiting bit-level sparsity)
000000000
Overall CNN Architecture
000
Results
0000●00

## System-level comparisons

Table: Comparisons to state-of-the-art implementations

|  | Eyerissv2[1] | ISSCC'20[2] | ISSCC'22[3] | RNSDNN[4] | This work RNS | **This work MF-D-S** |
|---|---|---|---|---|---|---|
| Process | 65 nm | 7 nm | 65 nm | 45 nm | 22 nm | **22 nm** |
| Supply voltage (V) | N/G | 0.575–0.825 | 1 | 1 | 0.65 | **0.65** |
| Frequency | 200 MHz | 290–880 MHz | 400 MHz | 1.2 GHz | 1 GHz | **1 GHz** |
| On-chip Memory (KB) | 246 | 2176 | 150 KB | N/G | 448 | **448** |
| Bit Precision (FMAP,wgt) | 8 | 8 | 8 | 16,8 | 12,8 | **12,8** |
| Network | AlexNet | MobileNet-v1 | VGG16 | VGG16 | VGG16 | **VGG16** |
| Performance (GOPS[†]) | 153.6 | 3604 | N/G | 134 | 220 | **364** |
| Area ($10^6$ Gates) | 2.69 | N/G | N/G | 2.18 | 4.74 | **5.25** |
| Area (mm$^2$) | N/G | 3.04 | 4.47 | N/G | 0.94 | **1.04** |
| Area Eff.(GOPS[†]/$10^6$ Gates) | 57.1 | N/G | N/G | 62 | 46.4 | **69.3** |
| Power Eff. (TOPS[†]/W) | 0.253 - 0.962[‡] | 3.28 - 6.66[‡] | 1.82 | 0.223 | 1.74/2.36* | **1.98/3.12*** |

*: full system/on-chip power cons.　　　[†]1 OP = 1 MAC　　　[‡]: for dense - sparse network　　N/G: not given

- 32% energy efficiency increase compared to RNS Counterpart
- 8.87× more energy efficient that state-of-the-art RNS CNN accelerator

## System-level comparisons

Table: Comparisons to state-of-the-art implementations

|  | Eyerissv2[1] | ISSCC'20[2] | ISSCC'22[3] | RNSDNN[4] | This work RNS | **This work MF-D-S** |
|---|---|---|---|---|---|---|
| Process | 65 nm | 7 nm | 65 nm | 45 nm | 22 nm | **22 nm** |
| Supply voltage (V) | N/G | 0.575–0.825 | 1 | 1 | 0.65 | **0.65** |
| Frequency | 200 MHz | 290–880 MHz | 400 MHz | 1.2 GHz | 1 GHz | **1 GHz** |
| On-chip Memory (KB) | 246 | 2176 | 150 KB | N/G | 448 | **448** |
| Bit Precision (FMAP,wgt) | 8 | 8 | 8 | 16,8 | 12,8 | **12,8** |
| Network | AlexNet | MobileNet-v1 | VGG16 | VGG16 | VGG16 | **VGG16** |
| Performance (GOPS[†]) | 153.6 | 3604 | N/G | 134 | 220 | **364** |
| Area ($10^6$ Gates) | 2.69 | N/G | N/G | 2.18 | 4.74 | **5.25** |
| Area (mm$^2$) | N/G | 3.04 | 4.47 | N/G | 0.94 | **1.04** |
| Area Eff.(GOPS[†]/$10^6$ Gates) | 57.1 | N/G | N/G | 62 | 46.4 | **69.3** |
| Power Eff. (TOPS[†]/W) | 0.253 - 0.962[‡] | 3.28 - 6.66[‡] | 1.82 | 0.223 | 1.74/2.36* | **1.98/3.12*** |

*: full system/on-chip power cons.    [†]1 OP = 1 MAC    [‡]: for dense - sparse network    N/G: not given

- 32% energy efficiency increase compared to RNS Counterpart
- 8.87× more energy efficient that state-of-the-art RNS CNN accelerator
- 2.05× more energy efficient that sparse version of Eyeriss

## References

📄 Y.-H. Chen, T.-J. Yang, J. Emer, and V. Sze, "Eyeriss v2: A flexible accelerator for emerging deep neural networks on mobile devices," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 9, no. 2, pp. 292–308, 2019.

📄 C.-H. Lin, C.-C. Cheng, Y.-M. Tsai, S.-J. Hung, Y.-T. Kuo, P. H. Wang, P.-K. Tsung, J.-Y. Hsu, W.-C. Lai, C.-H. Liu, S.-Y. Wang, C.-H. Kuo, C.-Y. Chang, M.-H. Lee, T.-Y. Lin, and C.-C. Chen, "7.1 a 3.4-to-13.3tops/w 3.6tops dual-core deep-learning accelerator for versatile ai applications in 7nm 5g smartphone soc," in *2020 IEEE International Solid- State Circuits Conference - (ISSCC)*, 2020, pp. 134–136.

📄 Y. Ju and J. Gu, "A 65nm systolic neural cpu processor for combined deep learning and general-purpose computing with 95locality and enhanced end-to-end performance," in *2022 IEEE International Solid- State Circuits Conference (ISSCC)*, vol. 65, 2022, pp. 1–3.

📄 N. Samimi, M. Kamal, A. Afzali-Kusha, and M. Pedram, "Res-DNN: A Residue Number System-Based DNN Accelerator Unit," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 67, no. 2, pp. 658–671, 2020.

Introduction
0000

Proposed PE (exploiting bit-level sparsity)
000000000

Overall CNN Architecture
000

Results
000000●

## Q & A

**Thank you!**