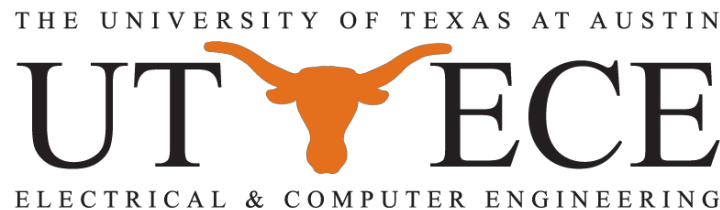# All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

Aseem Sayal, Shirin Fathima, S.S. Teja Nibhanupudi, Jaydeep P. Kulkarni

Department of Electrical and Computer Engineering

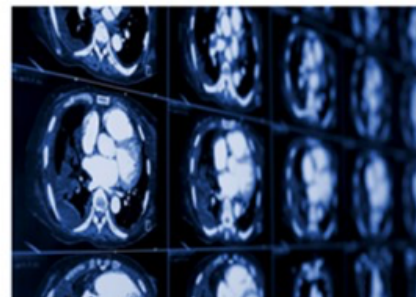The University of Texas at Austin, TX

THE UNIVERSITY OF TEXAS AT AUSTIN
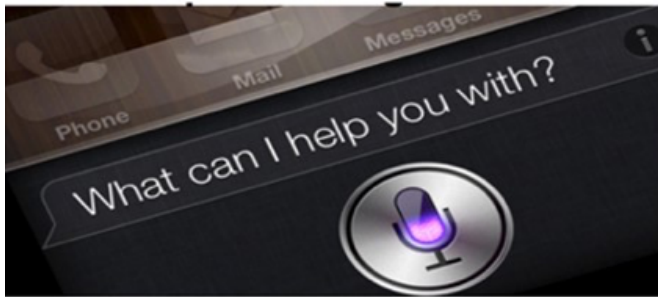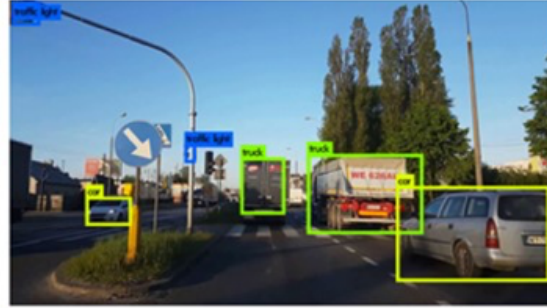
UT ECE

ELECTRICAL & COMPUTER ENGINEERING

# Outline

- Motivation

- Concept of Time-domain MAC

- Proposed CNN Engine Architecture

- Chip Implementation and Measurements

- Comparison with Prior Work

- Summary

# Need: Energy Efficient Edge Computing



## Applications

➢ Image Classification

➢ Speech Recognition

➢ Object Detection

**Energy efficient edge processing required for privacy and minimal data transmission**
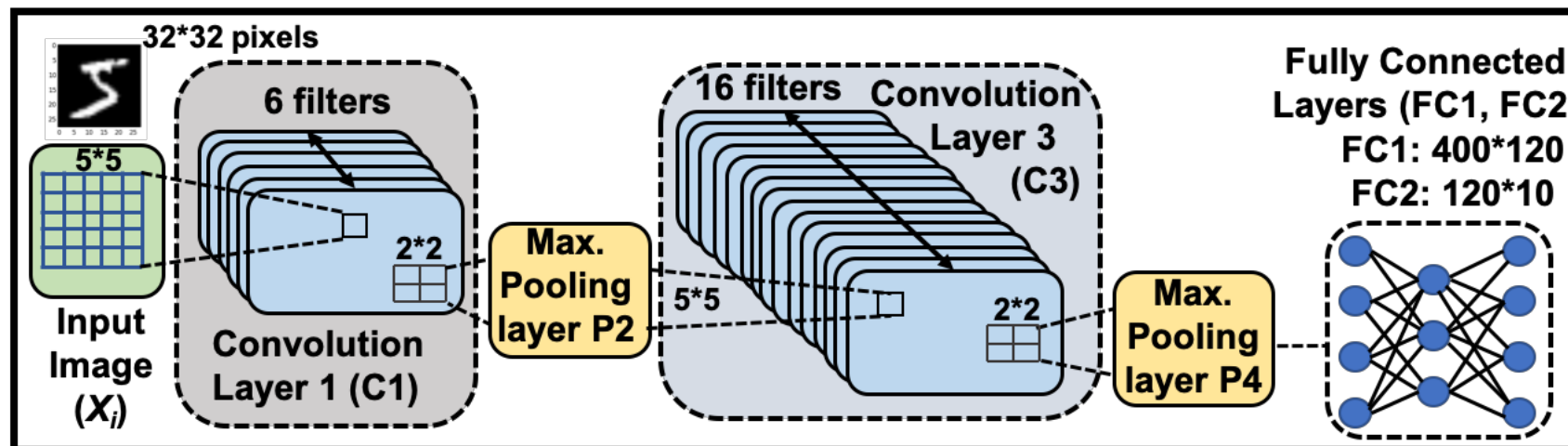
# Convolutional Neural Network (CNN)

- ➢ State-of-the-art classification accuracy for image/speech recognition

- ➢ Multiple filters to extract specific features

- ➢ Memory dominant and compute intensive

- ➢ Power hungry MAC operations in convolution

**Multiply-Accumulate-Average (MAV)**

$$Y = \frac{1}{N} \sum_{i=1}^{N} \left( X_i * w_i \right)$$

**N = # dot products per MAC**

**Output value**   **Input pixel**   **Weight**

**LeNet-5 CNN**

32*32 pixels

6 filters

5*5

2*2

**Input Image ($X_i$)**

**Convolution Layer 1 (C1)**

**Max. Pooling layer P2**

16 filters   **Convolution Layer 3 (C3)**

5*5   2*2

**Max. Pooling layer P4**

**Fully Connected Layers (FC1, FC2)**
**FC1: 400*120**
**FC2: 120*10**
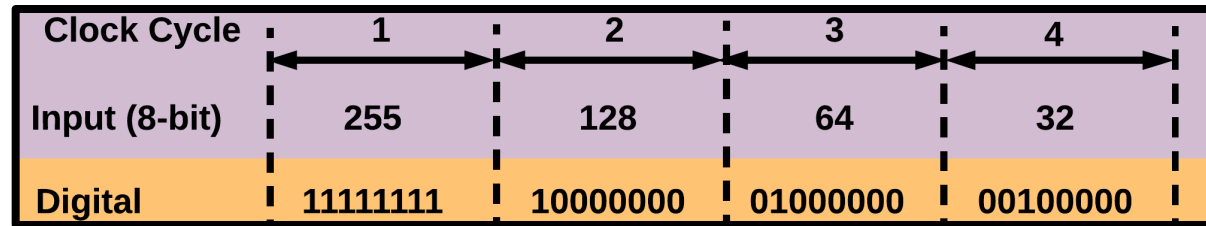
Ref:
LeCun et al., Gradient-based learning applied to document recognition. Proc. Of IEEE, 1998

# Why Time Domain Signal Processing?

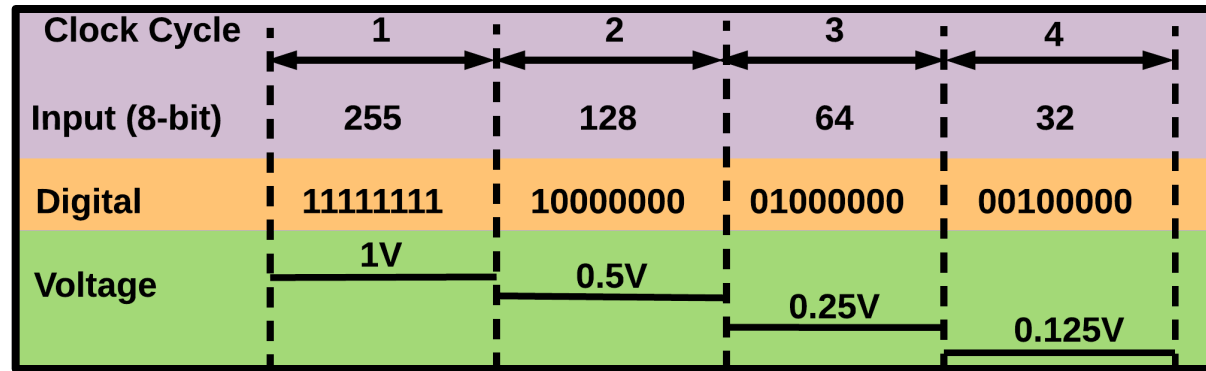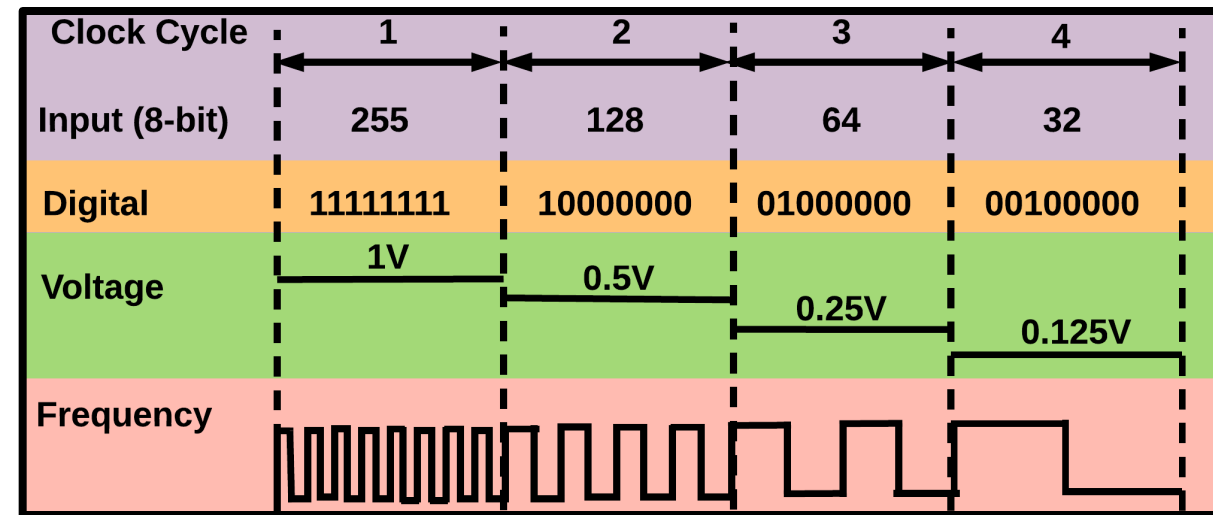| Clock Cycle | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Input (8-bit) | 255 | 128 | 64 | 32 |
| Digital | 11111111 | 10000000 | 01000000 | 00100000 |

➢ **Digital domain:** multi-bit digital vector – high $C_{DYN}$ and toggle activity ➔ consequently switching power
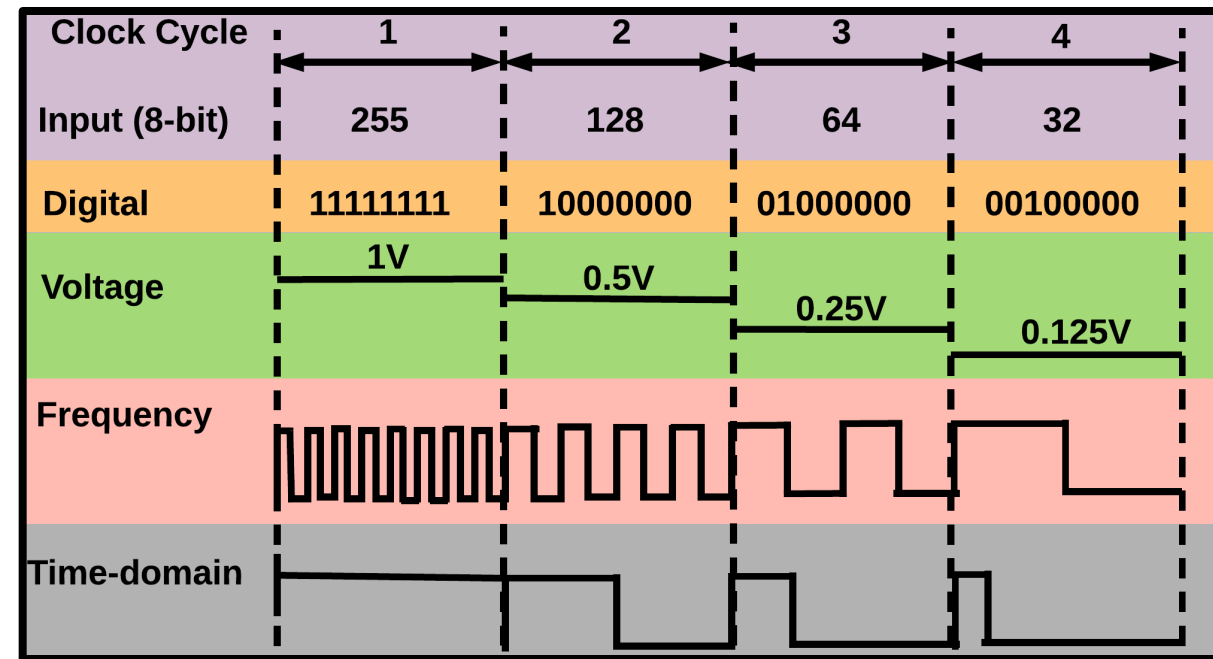
# Why Time Domain Signal Processing?

| Clock Cycle | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Input (8-bit) | 255 | 128 | 64 | 32 |
| Digital | 11111111 | 10000000 | 01000000 | 00100000 |
| Voltage | 1V | 0.5V | 0.25V | 0.125V |

➢ **Digital domain:** multi-bit digital vector – high $C_{DYN}$ and toggle activity ➔ consequently switching power

➢ **Analog Voltage domain:** continuously varying voltage signal – limited voltage scalability and need of power and area intensive ADCs/DACs

# Why Time Domain Signal Processing?

| Clock Cycle | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Input (8-bit) | 255 | 128 | 64 | 32 |
| Digital | 11111111 | 10000000 | 01000000 | 00100000 |
| Voltage | 1V | 0.5V | 0.25V | 0.125V |
| Frequency | | | | |

- **Digital domain:** multi-bit digital vector – high $C_{DYN}$ and toggle activity ➔ consequently switching power

- **Analog Voltage domain:** continuously varying voltage signal – limited voltage scalability and need of power and area intensive ADCs/DACs

- **Frequency domain:** frequency varying signal – need of accurate frequency generators/modulators ➔ limits performance scalability

# Why Time Domain Signal Processing?

| Clock Cycle | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Input (8-bit) | 255 | 128 | 64 | 32 |
| Digital | 11111111 | 10000000 | 01000000 | 00100000 |
| Voltage | 1V | 0.5V | 0.25V | 0.125V |
| Frequency | | | | |
| Time-domain | | | | |

- **Digital domain:** multi-bit digital vector – high $C_{DYN}$ and toggle activity → consequently switching power

- **Analog Voltage domain:** continuously varying voltage signal – limited voltage scalability and need of power and area intensive ADCs/DACs
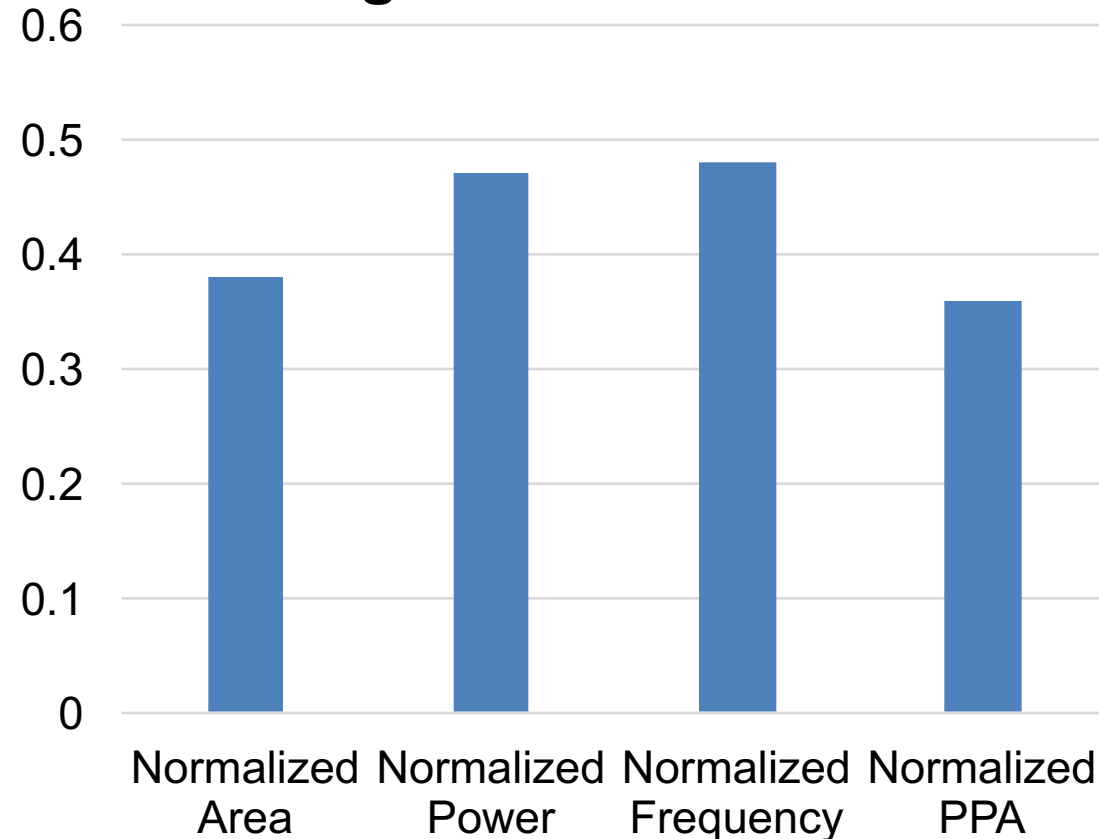
**Time domain:** multi-bit digital bit stream encoded as PWM signal – smaller $C_{DYN}$ and toggle activity, and ultra-low voltage operation resulting in low power, no need of multiple clocks/frequency modulators

- **Frequency domain:** frequency varying signal – need of accurate frequency generators/modulators → limits performance scalability

# Time Domain Approach *vs.* Digital Domain

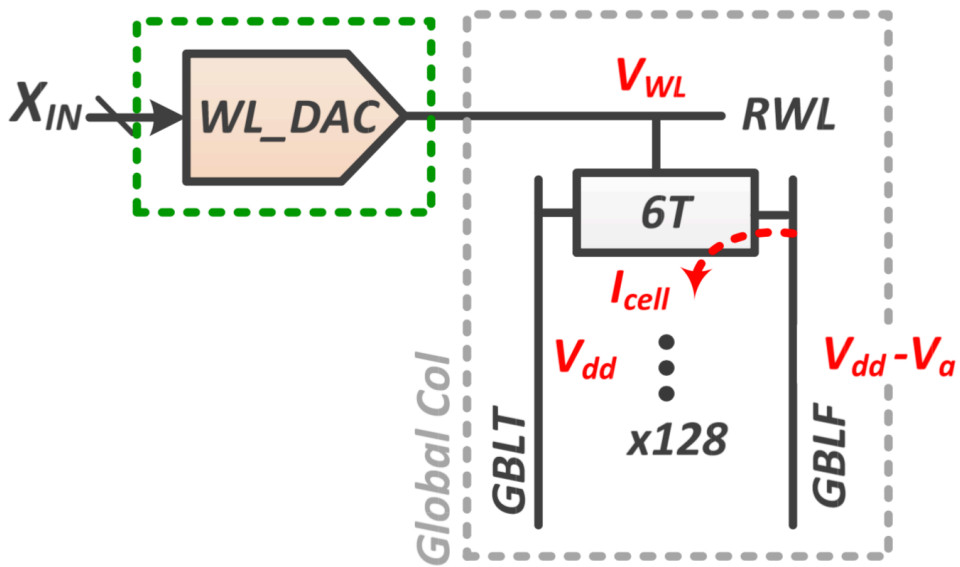| Design Metrics | Digital Domain | Proposed Time Domain |
|---|---|---|
| Configuration | 8b adder + 12b register for partial sum | 16-unit MDL + 12b counter for partial sum |
| Technology Node | 40nm | 40nm |
| Area (post-PNR) | $381\mu m^2$ | $144\mu m^2$ |
| $F_{MAX}$ supported | 1.67GHz | 0.8GHz |
| Power (post-PNR) | $54.06\mu W$ | $25.23\mu W$ |
| PPA/MAC ($\mu W*\mu m^2*ns$) | 12358 | 4541 |

PPA – Power Performance Area Product
PNR – Place and Route

## Time Domain Metrics Normalized *w.r.t.* Digital Domain Metrics



**2.72x PPA improvement**
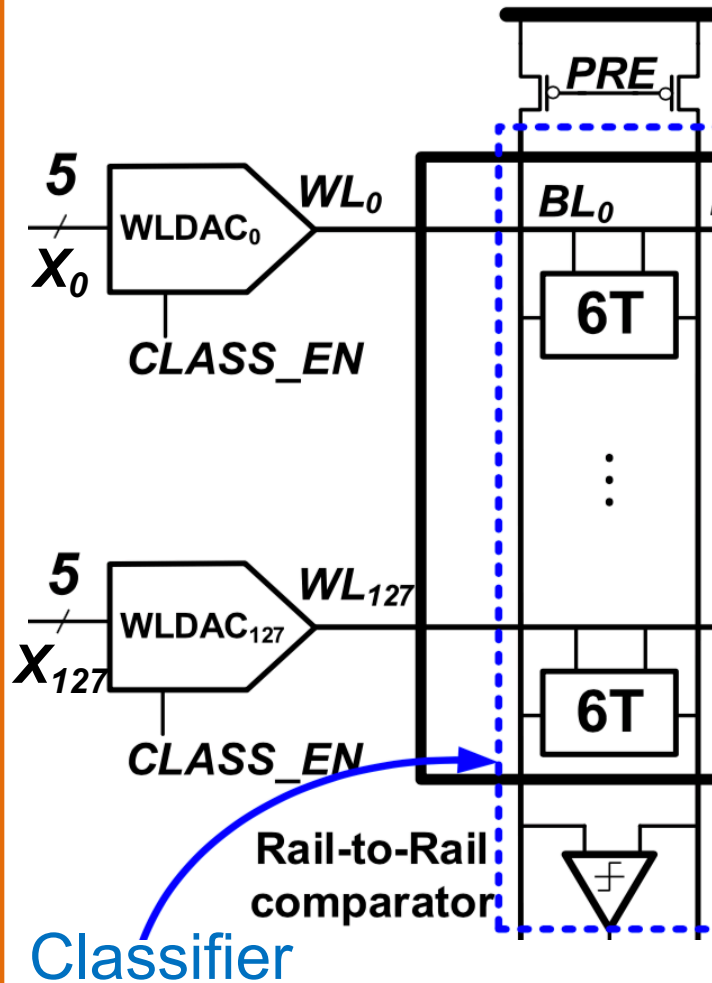
# Prior Work: Analog Voltage Approaches
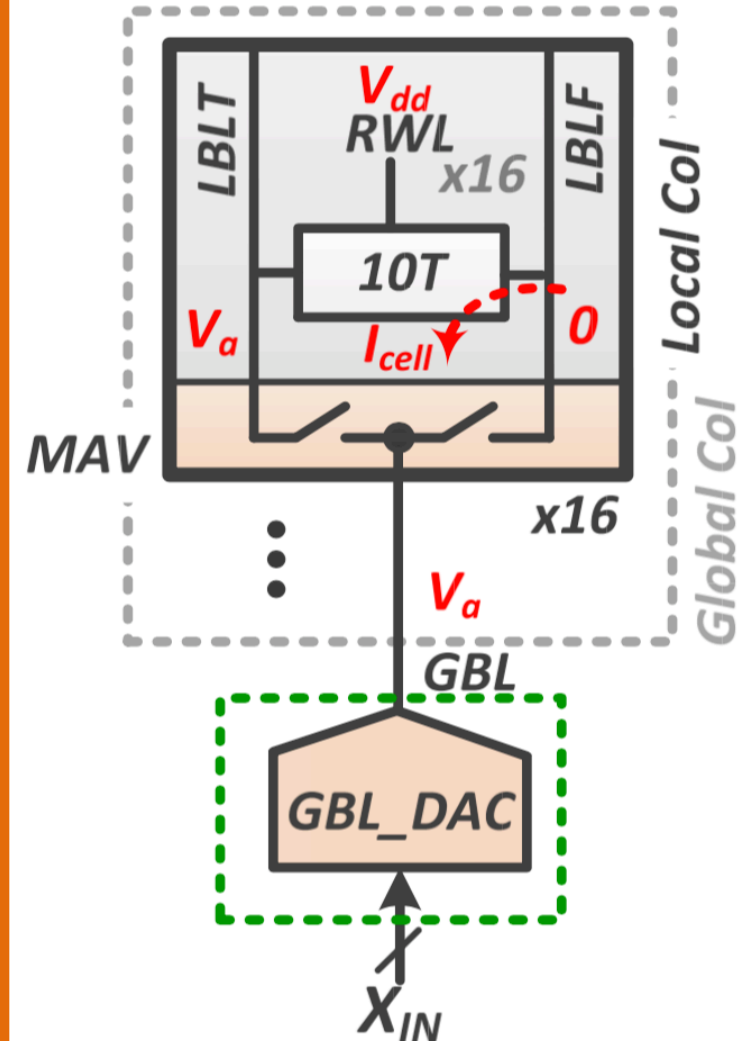


Bit-cell discharge current:

$$I_{cell} \propto (V_{WL} - V_t)$$

$$V_a \propto I_{cell}$$
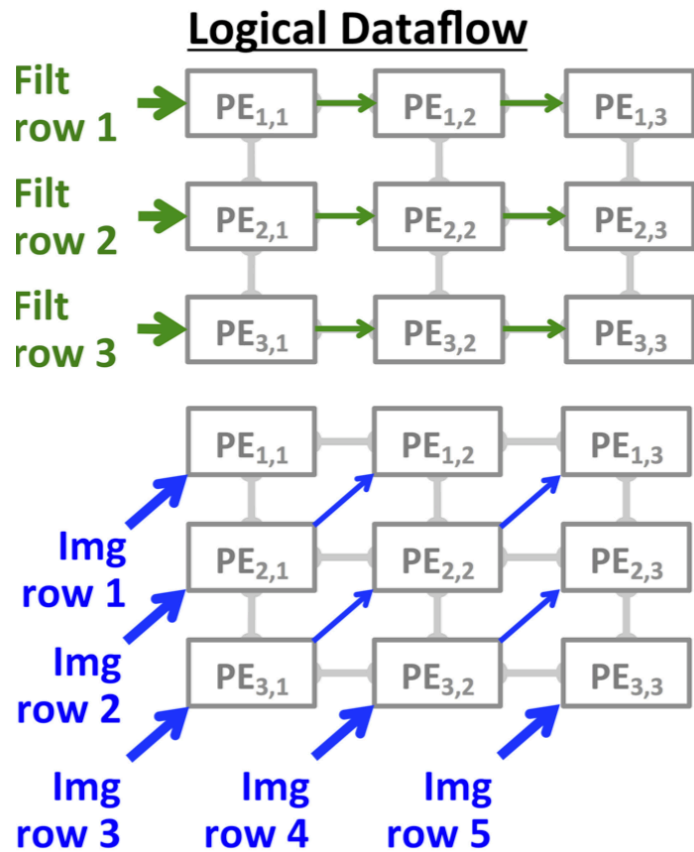
**M. Kang et al., VLSIC'16
(UIUC)**

**J. Zhang et al., JSSC'17
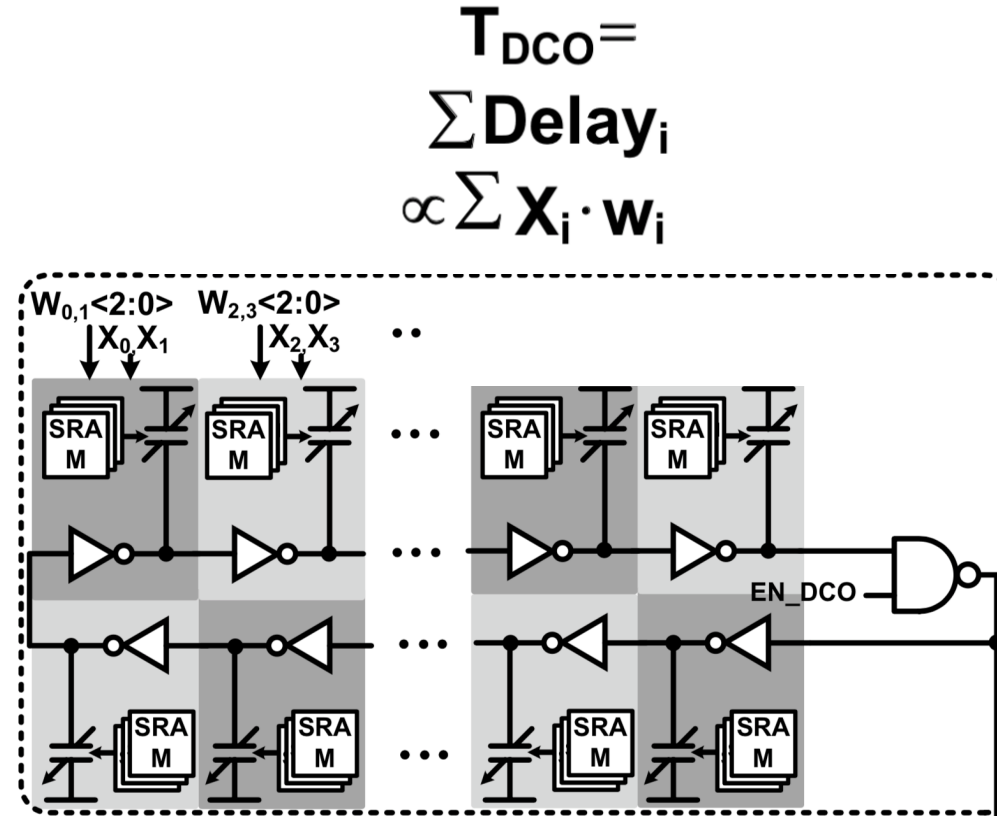(Princeton)**
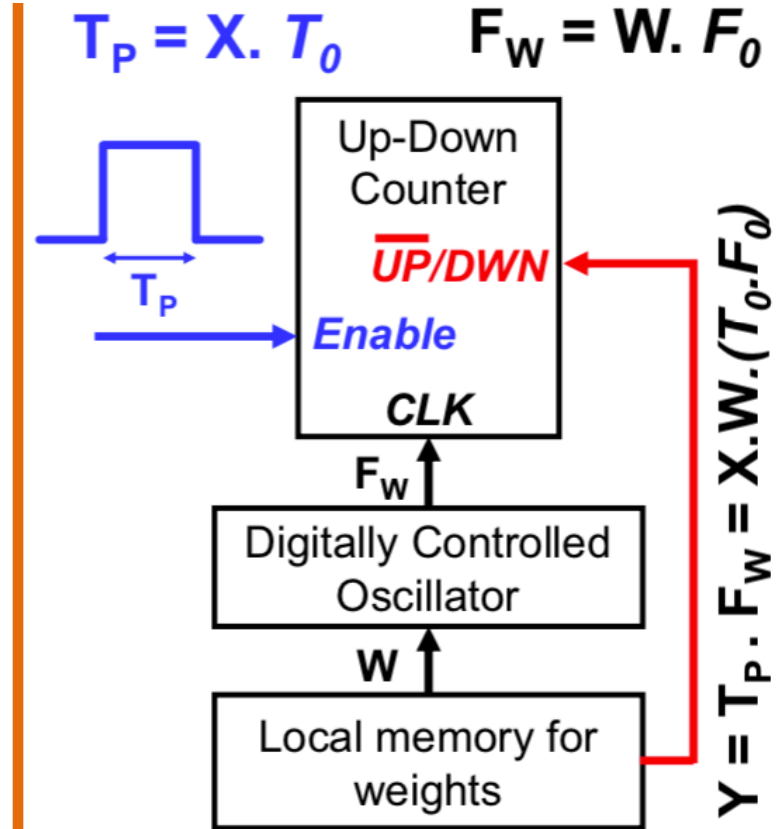
**A. Biswas et al., ISSCC'18
(MIT)**

# Prior Work: Digital, Time and Frequency domain



**Y. Chen et al., ISSCC'16 (MIT)**

**M. Liu et al., CICC'17 (Univ. of Minnesota)**

**A. Amravati et al., ISSCC'18 (Georgia Tech)**

# Outline

- Motivation

- Concept of Time-domain MAC

- Proposed CNN Engine Architecture

- Chip Implementation and Measurements

- Comparison with Prior Work

- Summary

# Concept: Time-domain MAC

**Image Pixels (X)**

| 123 | 60 | 43 |
|-----|-----|-----|
| 128 | 255 | 16 |
| 89 | 209 | 76 |

●

**Weights (w)**

| 1 | 0 | 1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | 0 | 1 |

MAC = 123(1) + 60(0) + 43(1) + 128(0) +255(1) + 16(1) + 89(0) + 209(0) + 76(1)
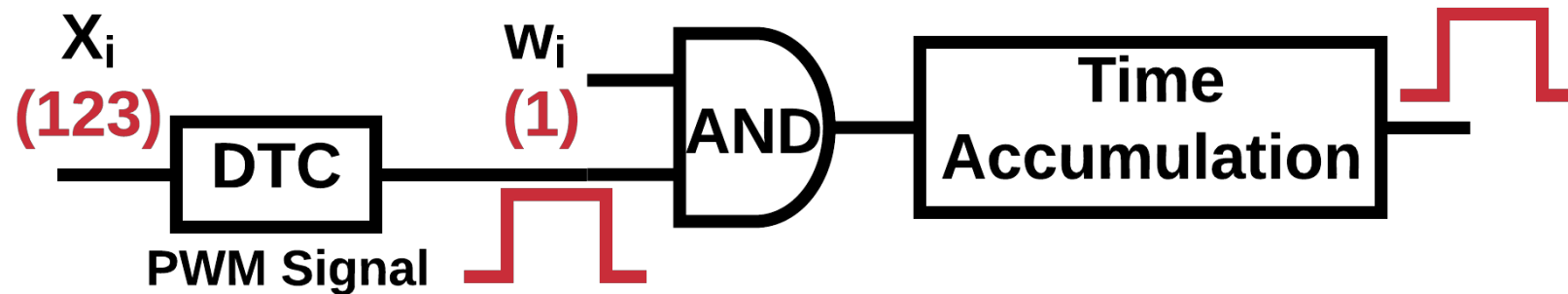
# Concept: Time-domain MAC

**Image Pixels (X)**

| 123 | 60 | 43 |
|-----|-----|-----|
| 128 | 255 | 16 |
| 89 | 209 | 76 |

●

**Weights (w)**

| 1 | 0 | 1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | 0 | 1 |

MAC = **123(1)** + 60(0) + 43(1) + 128(0) +255(1) + 16(1) + 89(0) + 209(0) + 76(1)

$X_i$
**(123)** → DTC → $w_i$ **(1)** → AND → Time Accumulation

PWM Signal

DTC – Digital to Time Converter

© 2019 IEEE
International Solid-State Circuits Conference

**14.4:** All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

*14 of 67*

# Concept: Time-domain MAC

**Image Pixels (X)**

| 123 | 60 | 43 |
|-----|-----|-----|
| 128 | 255 | 16 |
| 89 | 209 | 76 |

●

**Weights (w)**

| 1 | 0 | 1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | 0 | 1 |

MAC = 123(1) + 60(0) + 43(1) + 128(0) +255(1) + 16(1) + 89(0) + 209(0) + 76(1)

$X_i$
**(60)**

DTC

**PWM Signal**

$w_i$
**(0)**

AND

Time Accumulation

+

# Concept: Time-domain MAC



**Image Pixels (X)**

| 123 | 60 | 43 |
|-----|-----|-----|
| 128 | 255 | 16 |
| 89 | 209 | 76 |

●

**Weights (w)**

| 1 | 0 | 1 |
|-----|-----|-----|
| 0 | 1 | 1 |
| 0 | 0 | 1 |

MAC = 123(1) + 60(0) + 43(1) + 128(0) +255(1) + 16(1) + 89(0) + 209(0) + 76(1)

$X_i$ (76) → DTC → PWM Signal

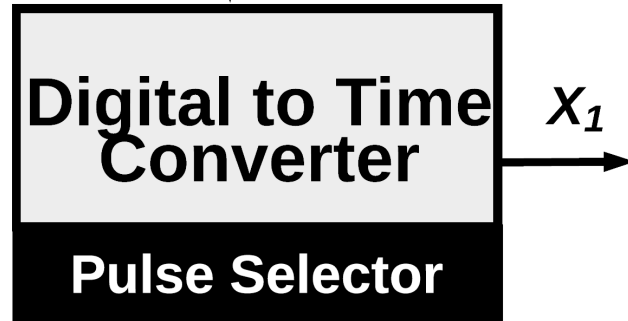$w_i$ (1) → AND → Time Accumulation → + ‒ + + + + + + + + = MAC

# Outline

- Motivation

- Concept of Time-domain MAC

- Proposed CNN Engine Architecture

- Chip Implementation and Measurements

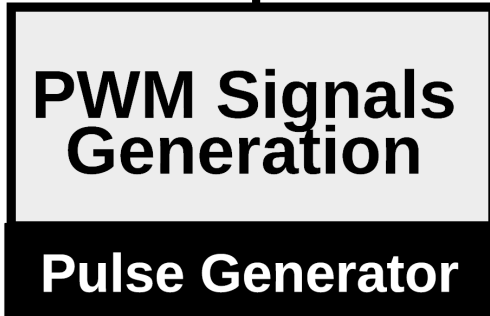- Comparison with Prior Work

- Summary

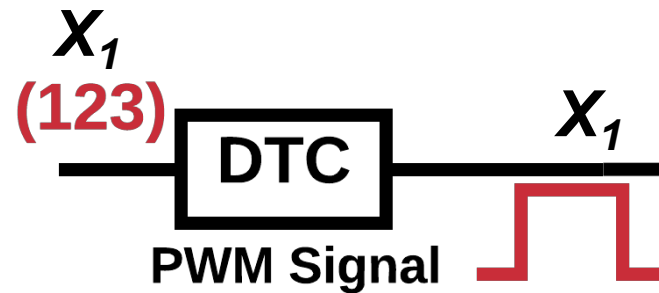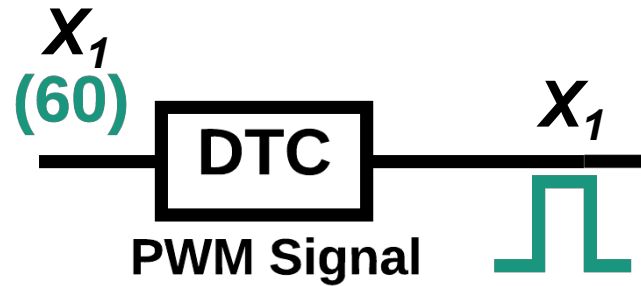# Architecture of Proposed CNN Engine

## Input image pixel

$X_1$ 8-bits

Digital to Time Converter

Pulse Selector

$X_1$

16

PWM Signals Generation

Pulse Generator

## Examples

$X_1$
(60)

DTC

$X_1$

PWM Signal

$X_1$
(123)

DTC

$X_1$

PWM Signal

# Architecture of Proposed CNN Engine

**Input image pixel**

**Weight Register**

**Parallel In, Serial Out Shift Register**

$X_1$ **8-bits**

$w_1$

**Digital to Time Converter**

$X_1$

**Dot product $(X_1*w_1)$**

**Pulse Selector**

**AND gate**

**16**

**PWM Signals Generation**

**Pulse Generator**

**14.4:** All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

*19 of 67*

# Architecture of Proposed CNN Engine

**Input image pixel**    **Weight Register**

$X_1$ **8-bits**

$w_1$

**Digital to Time Converter**

**Pulse Selector**

$X_1$

**Dot product** $(X_1 \ast w_1)$

**AND gate**

**Memory Delay Line (MDL)**

**Accumulation (MAC)**

**16**

**PWM Signals Generation**

**Pulse Generator**

# Architecture of Proposed CNN Engine

Input image pixel

**Weight Register**

$\nearrow X_1$ 8-bits

$w_1$

| Digital to Time Converter | $X_1$ | Dot product $(X_1*w_1)$ | Memory Delay Line (MDL) | Time to Digital Converter |
|---|---|---|---|---|
| **Pulse Selector** | | **AND gate** | **Accumulation (MAC)** | **Pos-edge Counter** |

$\nearrow$ 16

**PWM Signals Generation**

**Pulse Generator**

20

# Architecture of Proposed CNN Engine

Input image pixel

**Weight Register**

$X_1$ 8-bits

$w_1$

| Digital to Time Converter | $X_1$ | Dot product $(X_1 * w_1)$ | | Memory Delay Line (MDL) | | Time to Digital Converter |
|---|---|---|---|---|---|---|
| **Pulse Selector** | | **AND gate** | | **Accumulation (MAC)** | | **Pos-edge Counter** |

16

| PWM Signals Generation |
|---|
| **Pulse Generator** |

$$Y = \frac{1}{N}\sum_{i=1}^{N}\left(X_i * w_i\right)$$

| Bi-directional barrel shifter |
|---|
| **Averaging/Scaling** |

8

20

**MAV - Multiply Accumulate Average**

**Implementing MAV operation**

# Architecture of Proposed CNN Engine

© 2019 IEEE
International Solid-State Circuits Conference

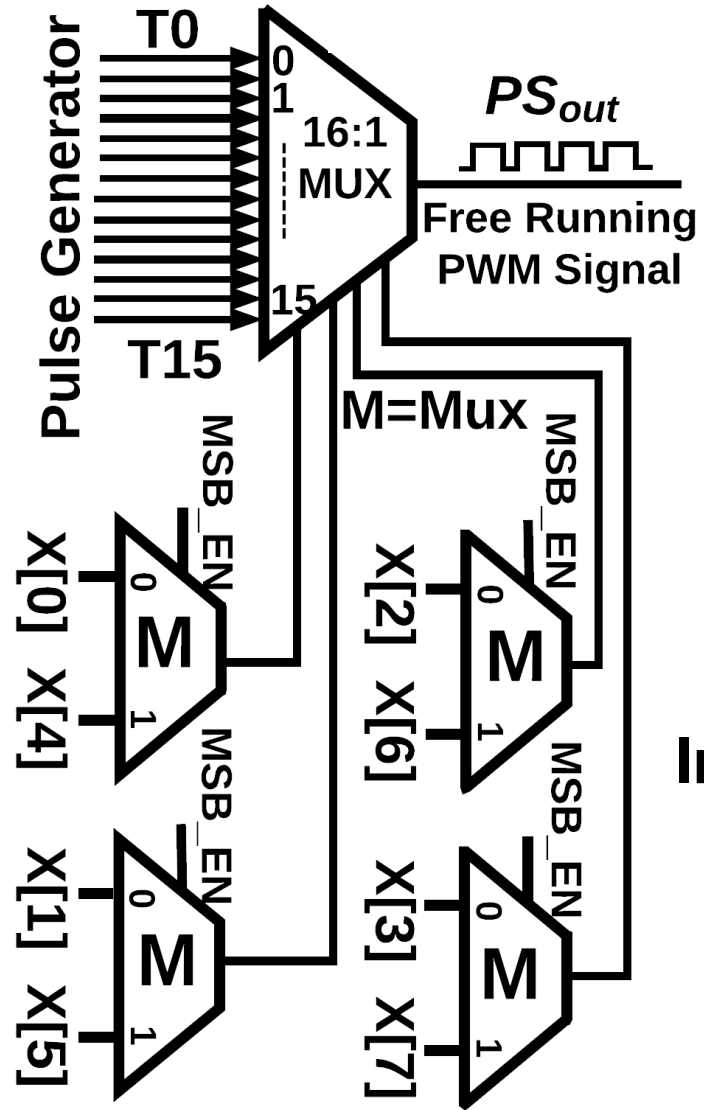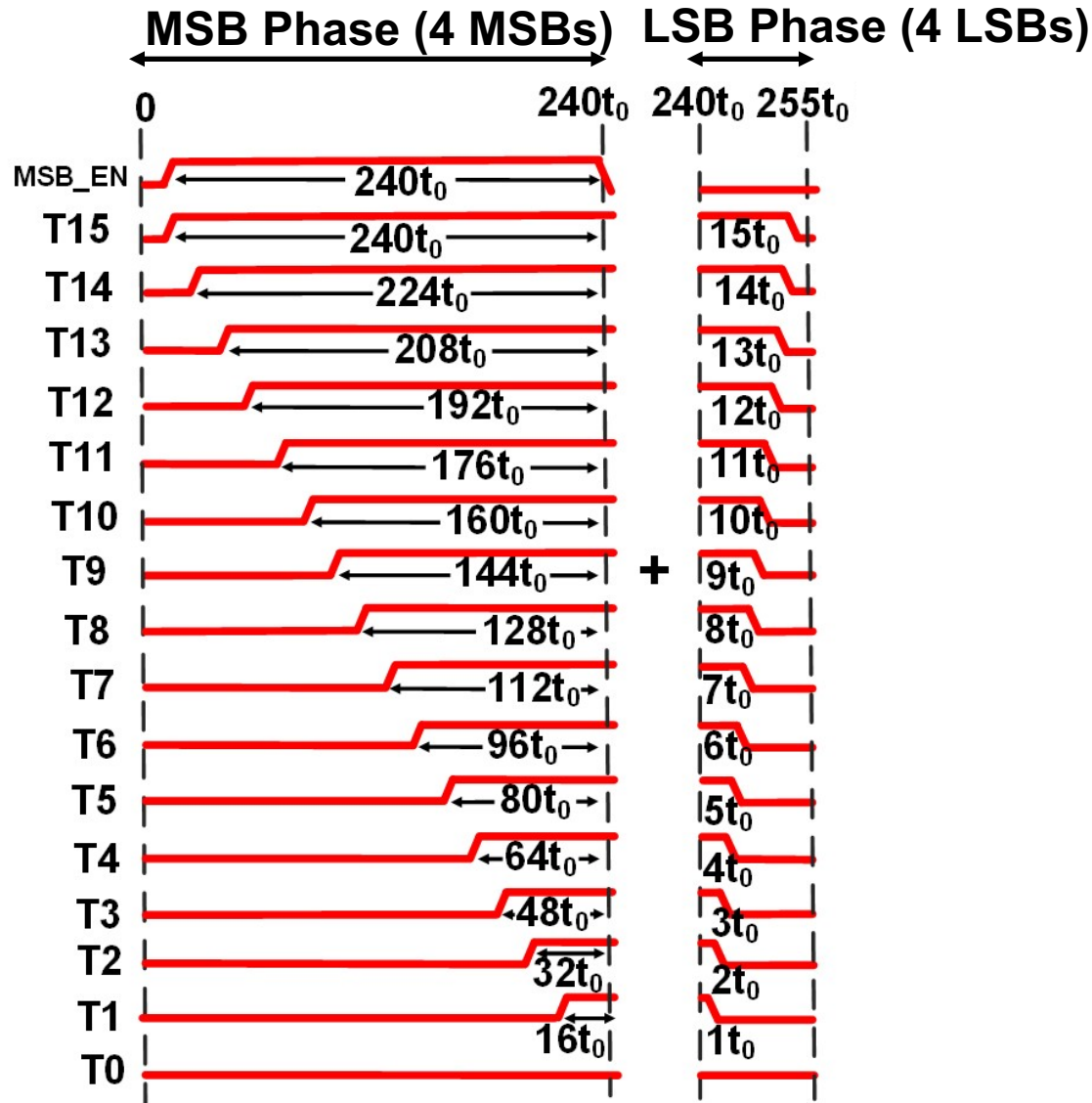*14.4:* **All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing**

*23 of 67*

# Architecture of Proposed CNN Engine

# PWM Pulse Generation & Selection



**Digital to Time Converter (DTC)**

**Input: X[7:0]**
**Output: PS_out**

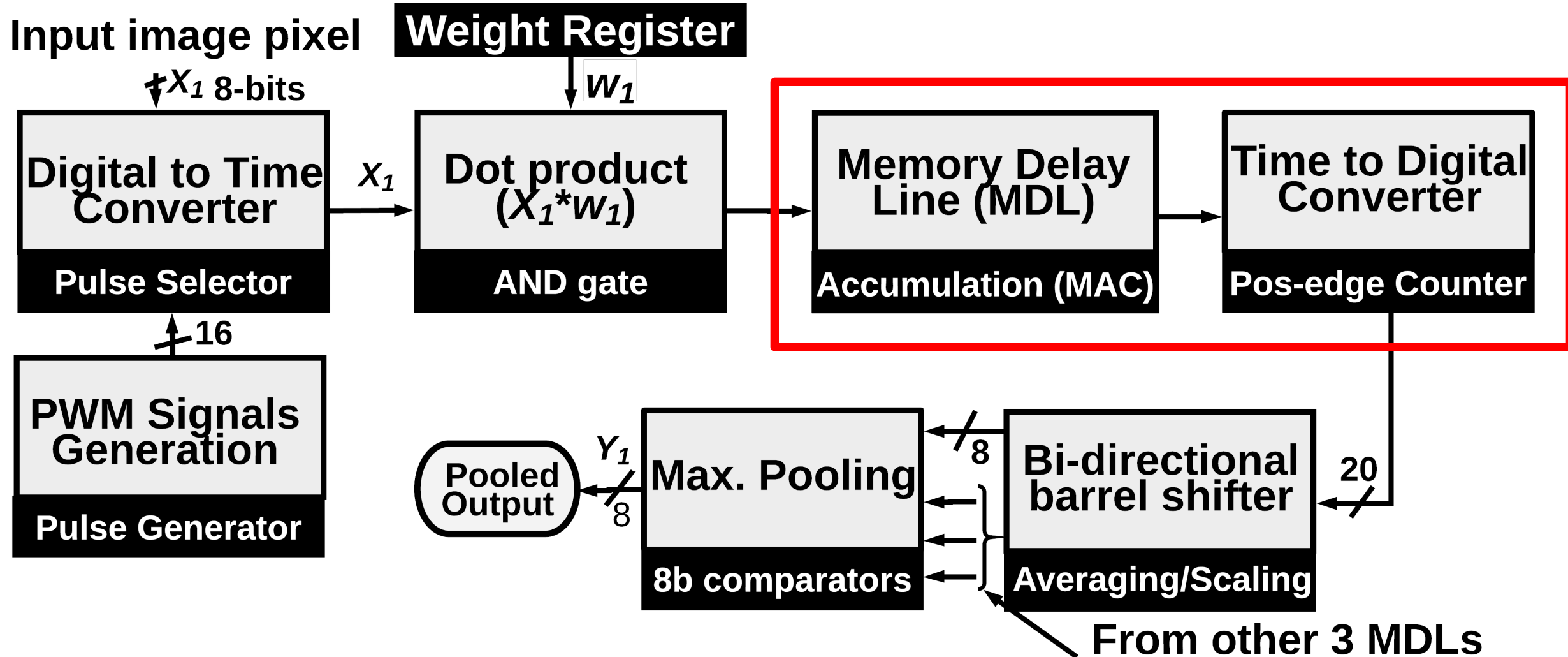**PWM Signal selection based on input Image pixel value**
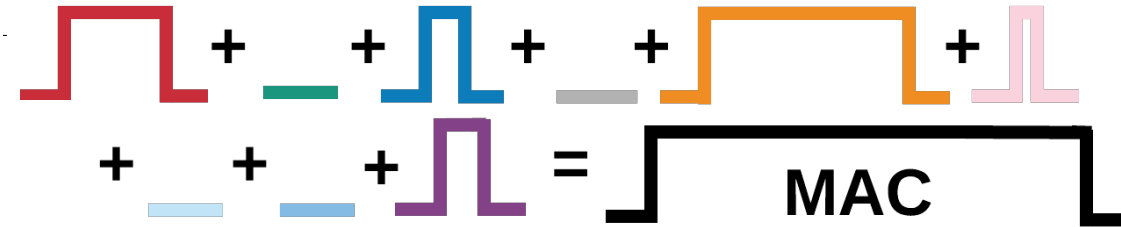
**Ref:**
**A. Biswas et al., ISSCC'18 (MIT)**
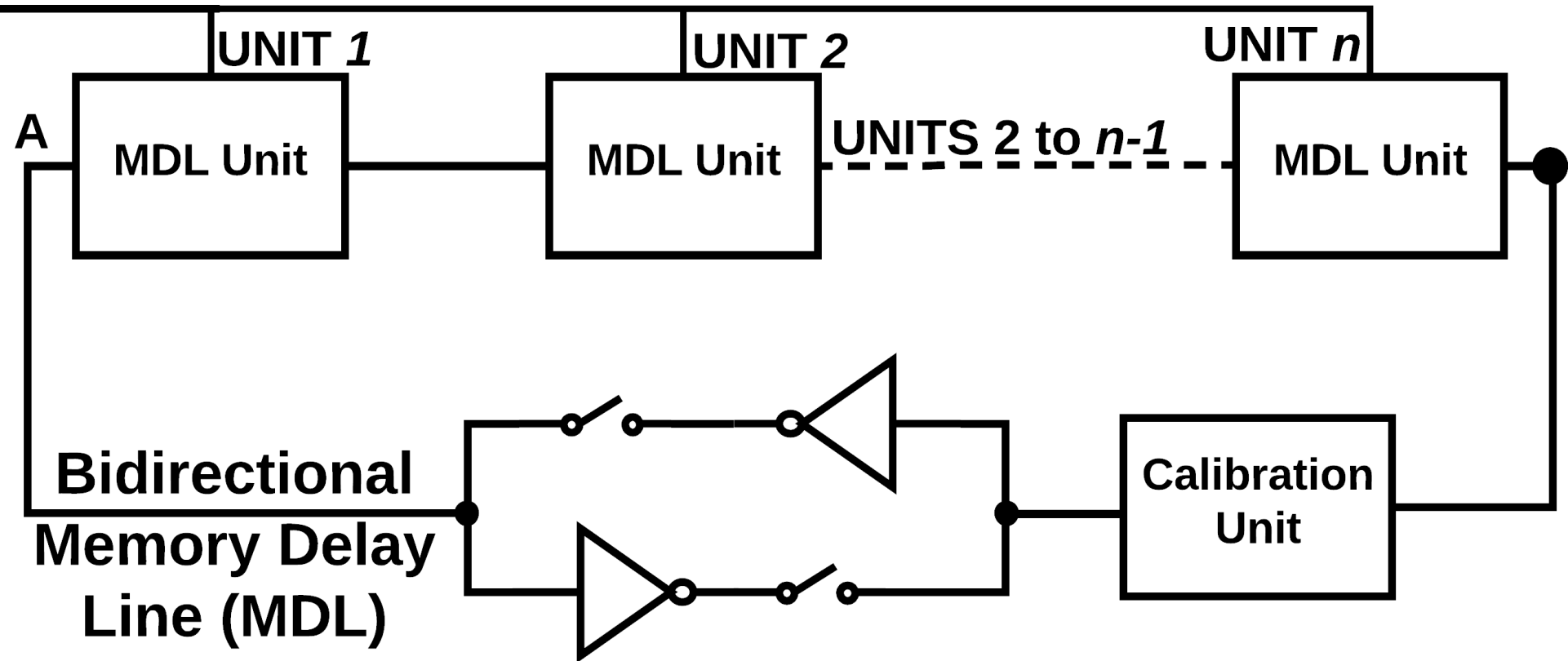
14.4: All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

# Architecture of Proposed CNN Engine

© 2019 IEEE
International Solid-State Circuits Conference

**14.4:** All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

*26 of 67*

# Bi-directional Memory Delay Line

**Time Accumulation**



MAC

$EN = X_i * w_i$

UNIT *1*  UNIT *2*  UNIT *n*

A

MDL Unit    MDL Unit    UNITS 2 to *n-1*    MDL Unit

**Bidirectional Memory Delay Line (MDL)**

Calibration Unit

*14.4:* All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

# Bi-directional Memory Delay Line



Time Accumulation

MAC

EN = $X_i * w_i$

UNIT 1     UNIT 2     UNIT n

A

MDL Unit   MDL Unit   UNITS 2 to n-1   MDL Unit

A

Down

Up

Up/Down +ve Edge Triggered Counter

**Bidirectional Memory Delay Line (MDL)**

Calibration Unit

**Time-to-Digital Converter (TDC)**

# MDL Approach inspired by Time Register
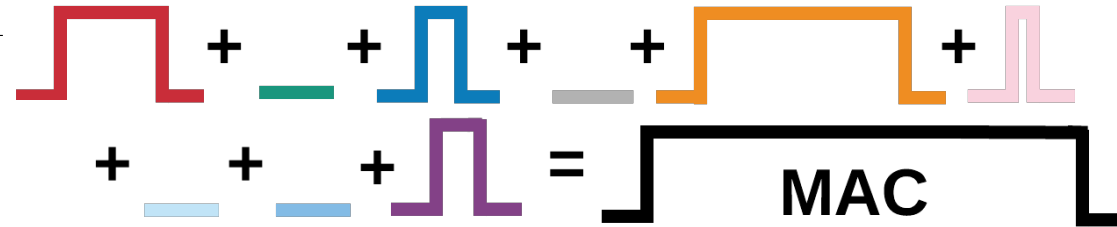


**Gated Delay Cell**

Ref: K. Kim et al., VLSIC'13

➢ Phase increases when EN is high, and is held when the EN is low

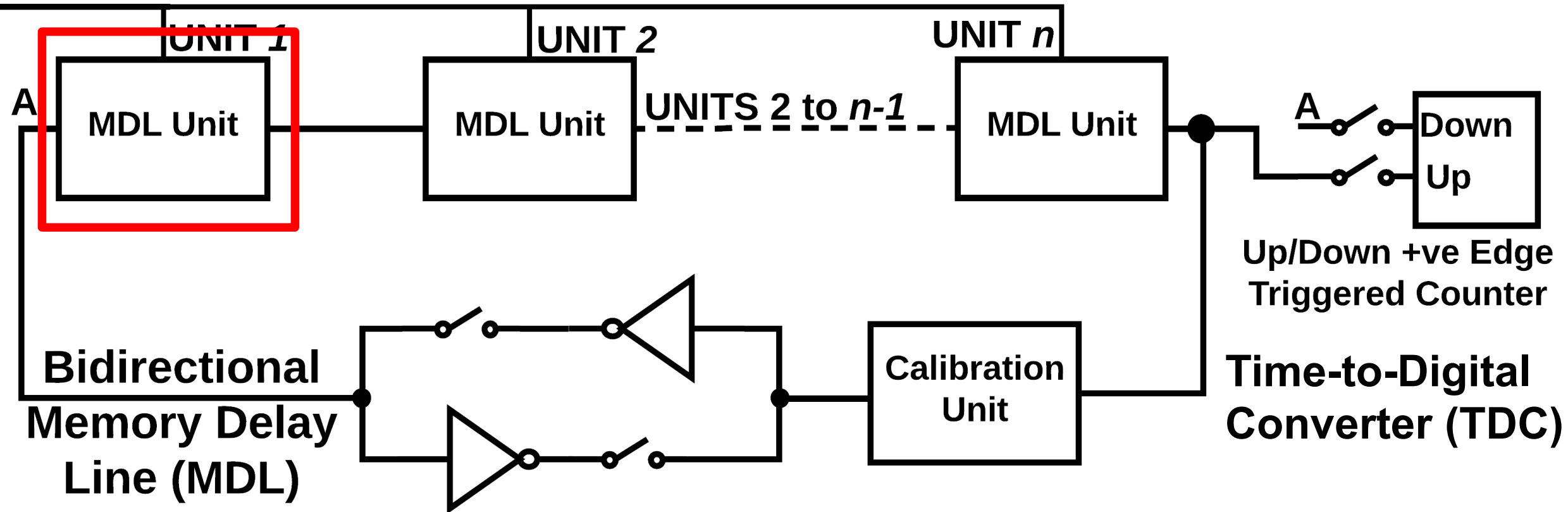➢ Metastability risk ➔ tri-stated output when long period of EN disable

# Bi-directional Memory Delay Line

**Time Accumulation**



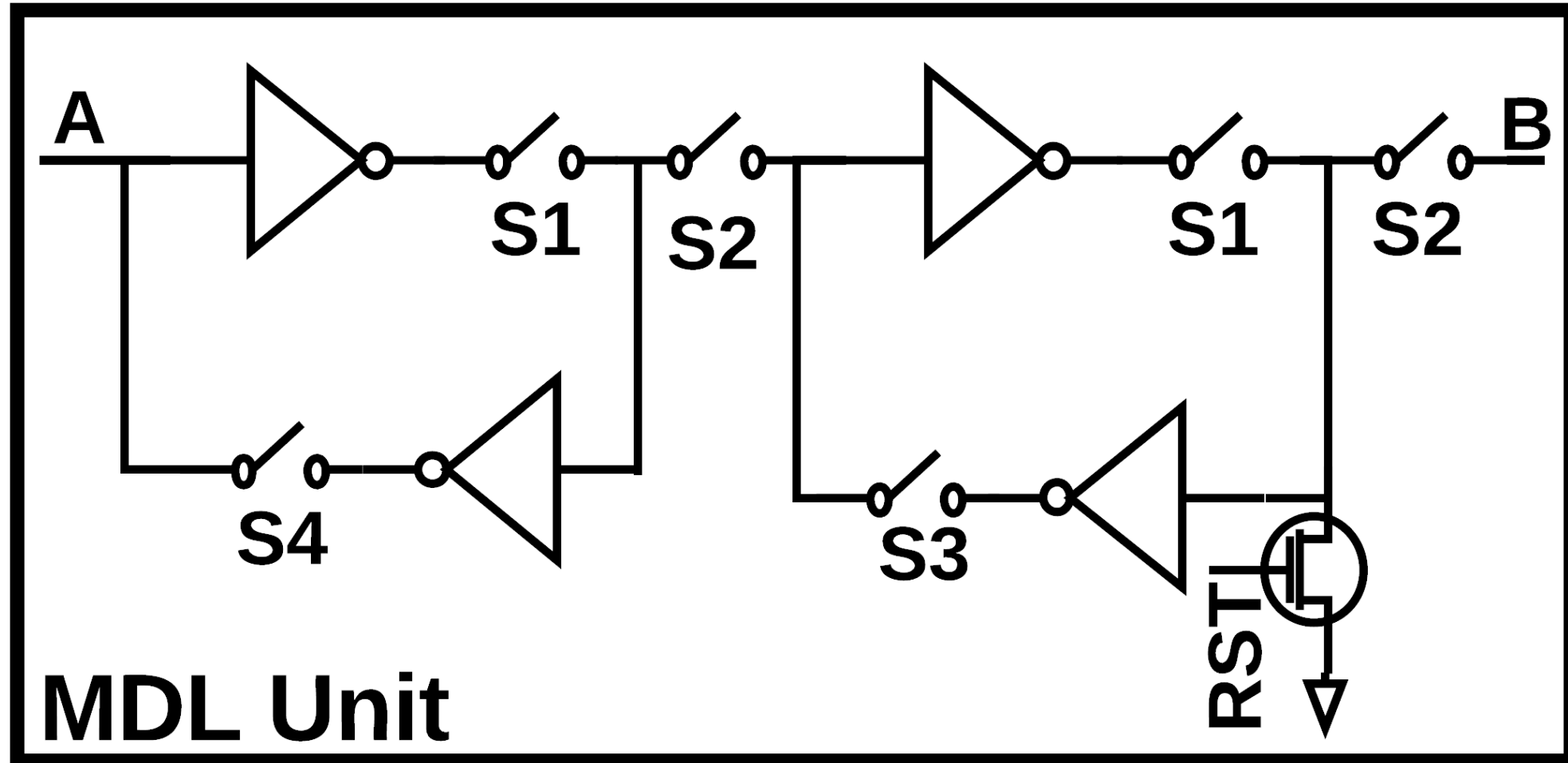MAC

EN = $X_i * w_i$



UNIT *1*

UNIT *2*

UNIT *n*

A

MDL Unit

MDL Unit

UNITS 2 to *n-1*

MDL Unit

A

Down

Up

**Up/Down +ve Edge Triggered Counter**

**Bidirectional Memory Delay Line (MDL)**

Calibration Unit

**Time-to-Digital Converter (TDC)**

# Bi-directional Memory Delay Line Unit

# MDL Unit – Delay Line (Negative weight)



**A**    **B**

S1    S2    S1    S2

S4    S3

RST

**DELAY LINE (-ve weight)**

**MDL Unit**

**S2,S3,S4 - ON**

**S1 - OFF**
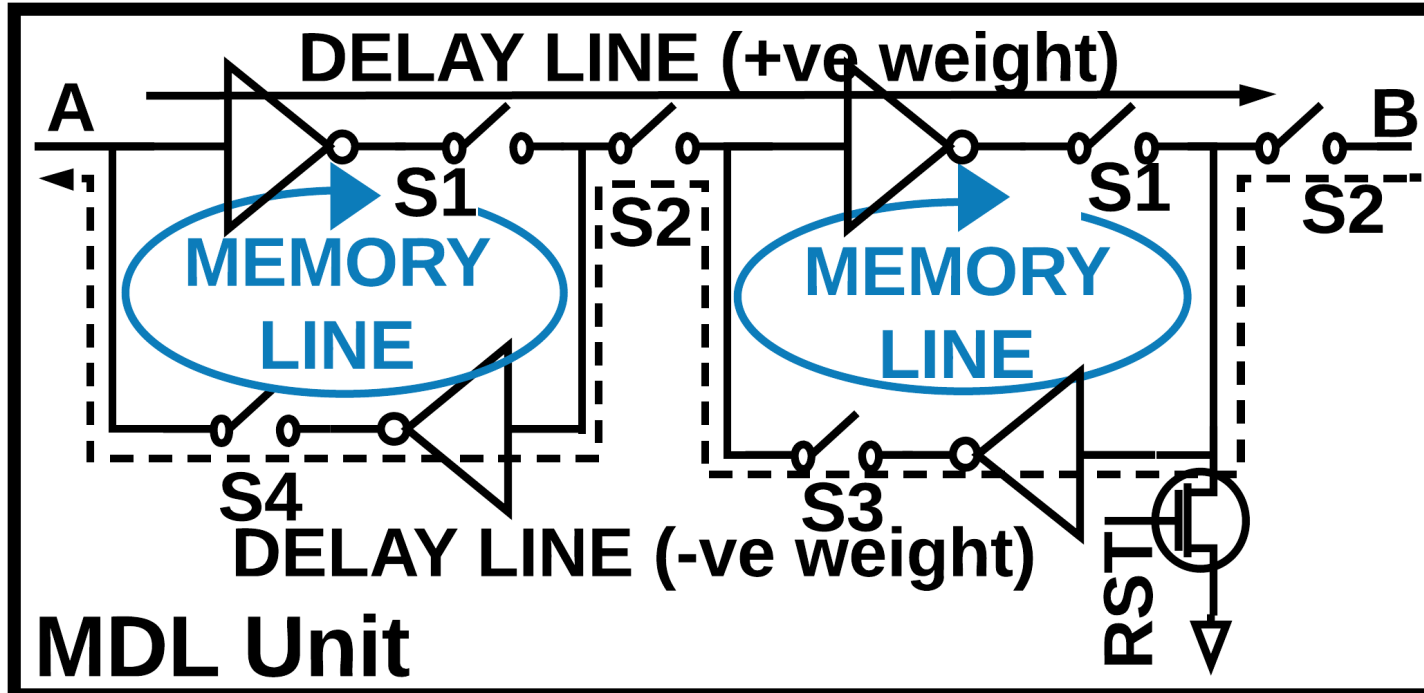
# MDL Unit – Memory Line



No tristate cells ➔ State of the delay line is retained in memory phase

# Bi-directional Memory Delay Line Unit



| Switch | Value |
|--------|-------|
| S1 | (EN).(SIGN) |
| S2 | EN |
| S3 | (EN') OR (EN).(SIGN)' |
| S4 | (EN).(SIGN)' |

if EN==1,
   MDL ➜ Delay Line
if EN==0,
   MDL ➜ Memory Line

if weight == +1,
   SIGN = 1
if weight == -1,
   SIGN = 0

# Bi-directional MDL Implementation

$X_i * w_i$



Up/Down +ve Edge Triggered Counter

# MDL units per Block = 16
# MDL blocks in MDL = 1 or 2 or 3 or 4

| Switch | Value | Weight |
|--------|-------|--------|
| S5 | SIGN | +ve |
| S6 | (SIGN)' | -ve |

➔ **Clockwise propagation in MDL**

➔ **Anti-clockwise propagation in MDL**

# Bi-directional MDL Implementation



$X_i*w_i$

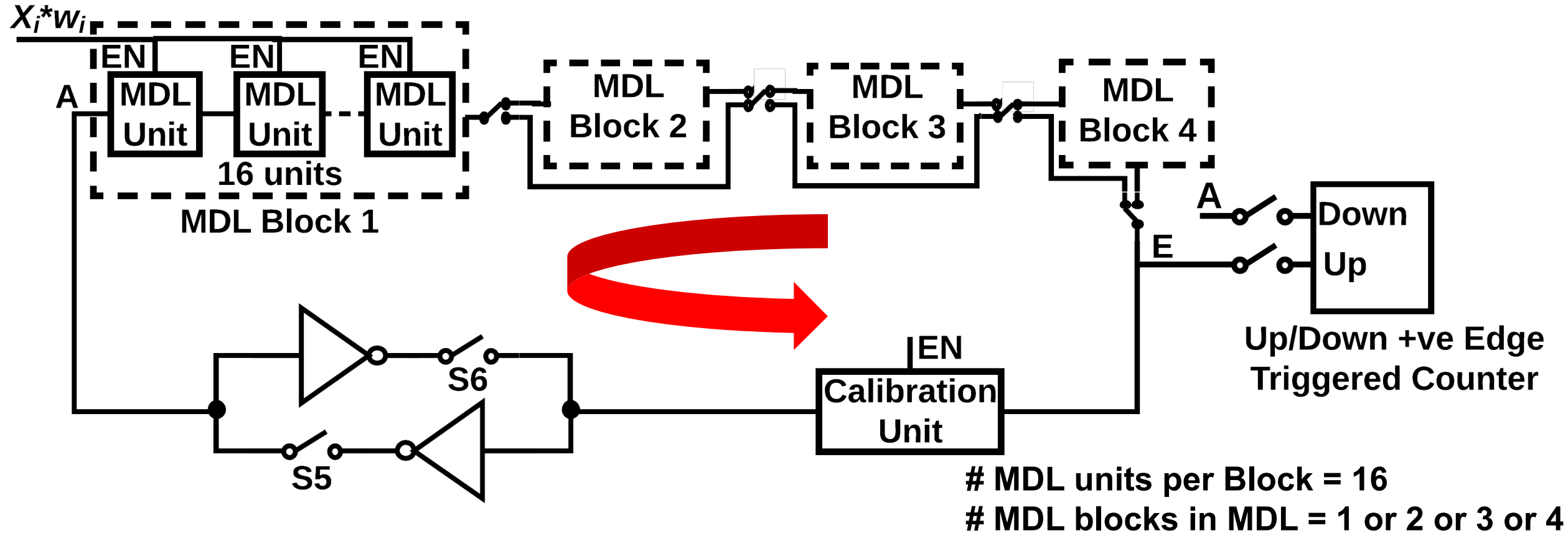MDL Block 1 — 16 units — EN MDL Unit, EN MDL Unit, EN MDL Unit

MDL Block 2, MDL Block 3, MDL Block 4

A

Down / Up

Up/Down +ve Edge Triggered Counter

EN Calibration Unit

S5, S6

E

# MDL units per Block = 16
# MDL blocks in MDL = 1 or 2 or 3 or 4

| Switch | Value | Weight |
|--------|--------|--------|
| S5 | SIGN | +ve |
| S6 | (SIGN)' | -ve |

➔ **Clockwise propagation in MDL**

➔ **Anti-clockwise propagation in MDL**

# MDL Timing Diagram



**MAC_CLK PERIOD**

$X_0 * w_0$

**EN**

**DELAY LINE PHASE**

**MEMORY STORAGE PHASE**

**A**

**E**

**Possibility of Metastability**

**MDL State Vector**

**EN** UNIT *1* UNIT *2* UNIT *n*

**A** MEMORY DELAY LINE UNIT — MEMORY DELAY LINE UNIT — — MEMORY DELAY LINE UNIT **E** **A**

Down / Up

**Up/Down +ve Edge Triggered Counter**

Calibration Unit

# MDL Timing Diagram

# MDL Timing Diagram

# MDL Timing Diagram

# Mitigating Process Variations



Switch S9 value == SIGN, Switch S10 value == (SIGN)'

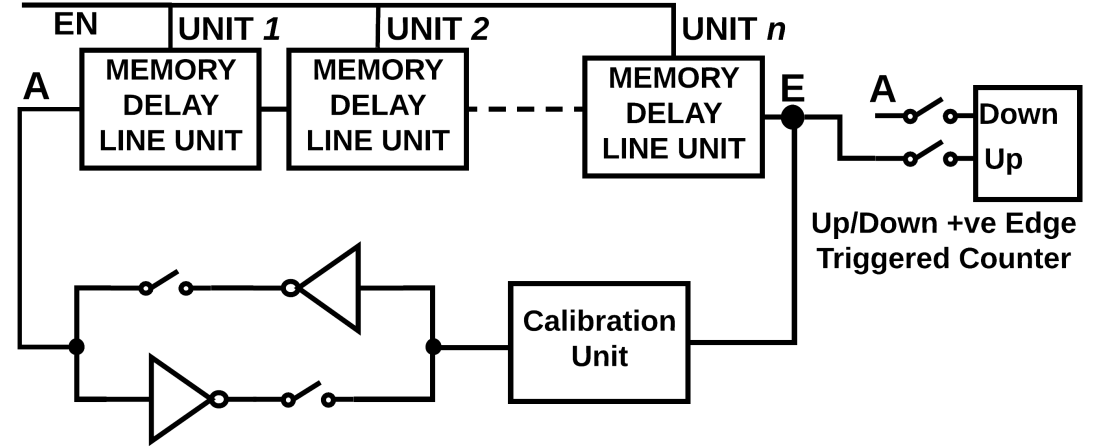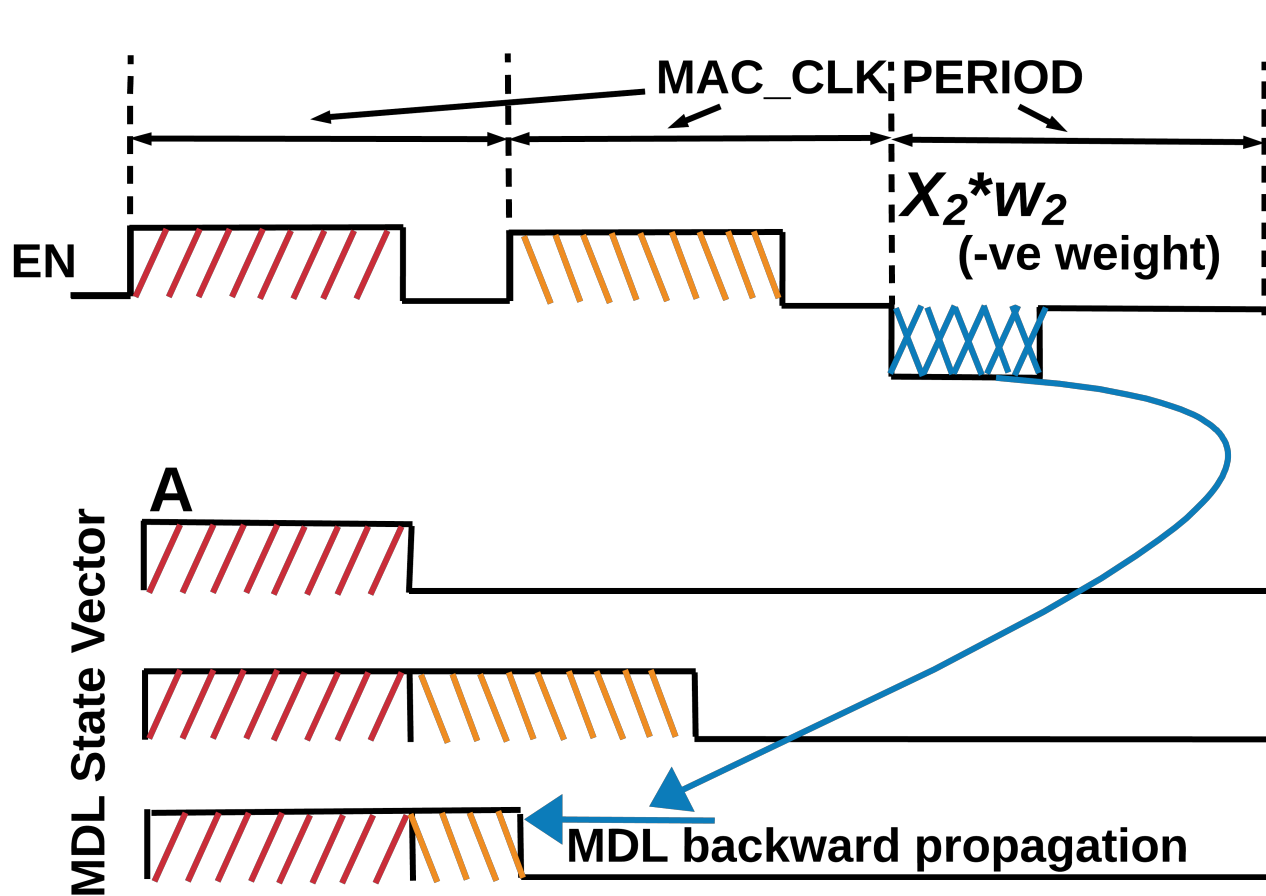➢ Calibration unit is added to all the 4 MDLs for each filter. Due to process variations, delay mismatch among 4 MDLs for each filter are compensated by by tuning *cal_bit[0:2]* and *cal_enable* signals.

# Increasing Throughput: Speedup modes



| Speedup Mode | 1x | 4x | 8x | 16x |
|---|---|---|---|---|
| Input clock period | $2t_o$ | $8t_o$ | $16t_o$ | $32t_o$ |
| MAC clock period | $256t_o$ | $256t_o$ | $256t_o$ | $256t_o$ |
| # input clock cycles/MAC clock cycle | 128 | 32 | 16 | 8 |
| Quantization error (input: 0-255) | 0 | $\pm 2$ | $\pm 4$ | $\pm 8$ |

# Proposed CNN Engine Implementation



**4 input image pixels each 8b**
$X_1 X_2 X_3 X_4$

**Filter 16**
**Filter 1**
Parallel In **25** → **Weight Register**
(4) (3) (2) (1) **4 MAC Blocks**

**Pulse Selector**
**Digital to Time Converter (DTC)**
$PS_{out}$ $x_1$ $x_2$ $x_3$ $x_4$
**16 ↗ T0-T15**

**Pulse Generator**
**4 precision (speed-up) modes**

**Pooled output of 16 Filters**
$Y_1$ **8**

Serial Out $w^i$

**Pulse Gating Logic**
$x^i_1$
**Pass once**

**Weight multiplication**
$(X_1^i * w^i)$
**AND gate**

**Memory Delay Line (MDL)**
**Time Accumulation**

**Max. Pooling**
C3 C1
C2
**3 8b Comparators**

**8**

**Bi-Directional Barrel Shifter**
0-7b shifts
**Averaging/Scaling**

**20**

**Counter (TDC) pos-edge**
**Digital MAC value**

# # MAC Blocks per Filter = 4
# # Filters = 16

# Outline

- Motivation

- Concept of Time-domain MAC

- Proposed CNN Engine Architecture

- <span style="color:orange">Chip Implementation and Measurements</span>

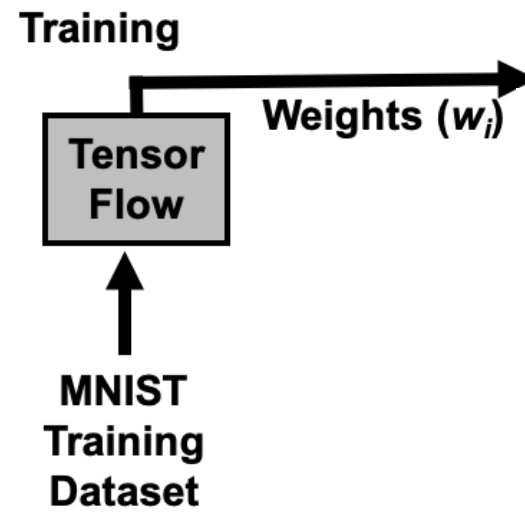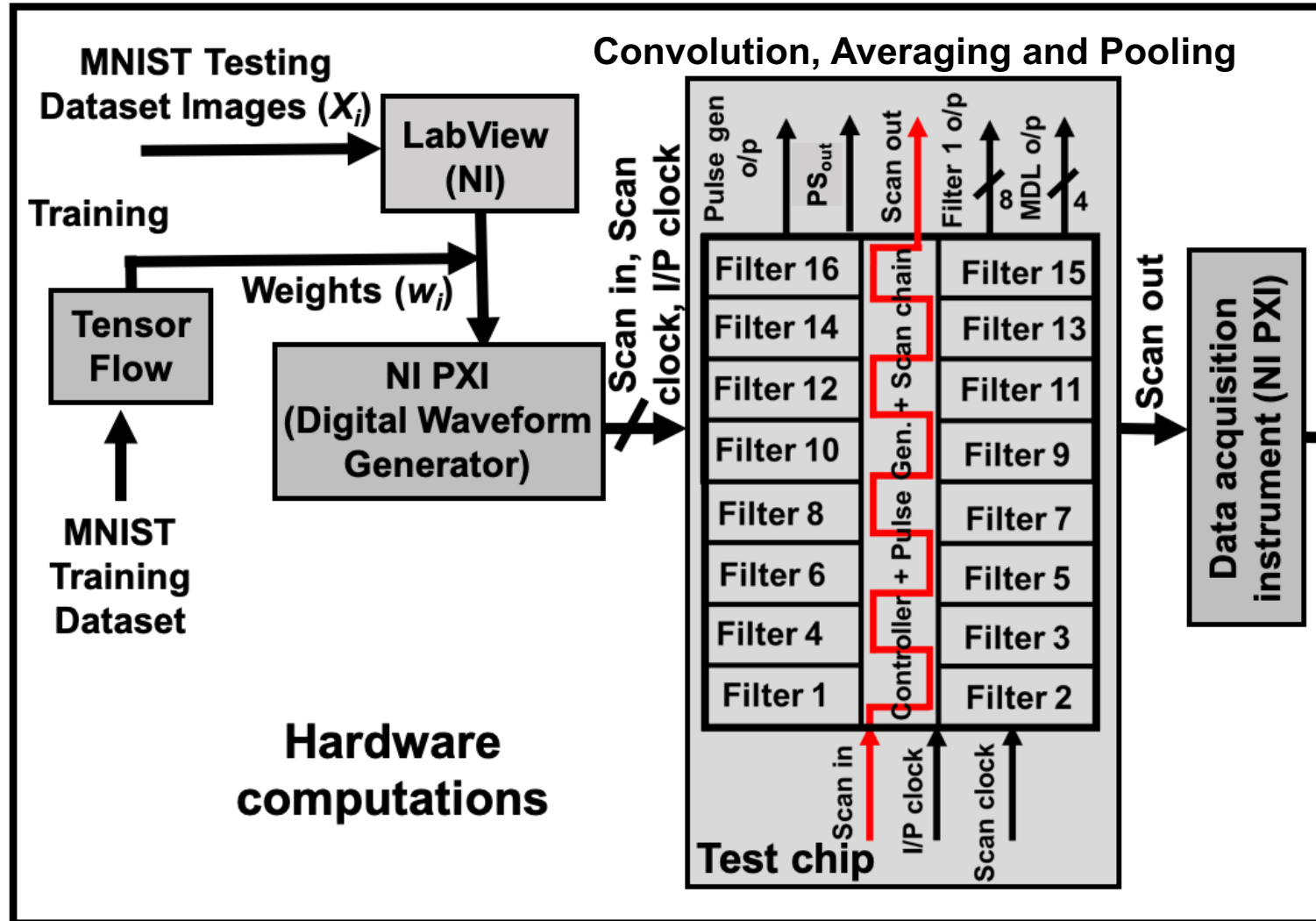- Comparison with Prior Work

- Summary

# Overall Dataflow Diagram



Training

Weights ($w_i$)

Tensor Flow

MNIST Training Dataset

# Overall Dataflow Diagram

# Overall Dataflow Diagram

# LeNet-5 Architecture & Parameters



Ref: LeCun et al., Gradient-based learning applied to document recognition. Proc. Of IEEE, 1998

| Parameters for LeNet-5 | C1 | C3 |
|---|---|---|
| Filter Size | 5*5*1*6 | 5*5*6*16 |
| Input/Filter bit width | 8bits/ 1bit | 8bits/ 1bit |
| Input Size | 32*32*1 | 14*14*6 |
| Output Size | 14*14*6 | 5*5*16 |
| #Filters | 6 | 16 |
| #Operations/convolution* | (25*4*6)*2 | (150*4* 16)*2 |

*Assuming 1 Multiply-Accumulate-Average, and 1 Pooling = 2 operations

**Hardware computations**

**Convolution, Averaging and Pooling (Fixed Point)**

**Software computations**

**FCN (16 bit Floating Point)**

# 40nm test-chip die micrograph



| Technology | 40nm CMOS |
|---|---|
| Area | 0.124mm$^2$ |
| # standard cells | ~ 40,000 |
| # of filters supported | 16 |
| Data type supported | Signed/Unsigned |
| Voltage Range | 0.375V - 1.10 V |

© 2019 IEEE
International Solid-State Circuits Conference

*14.4:* **All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing**

*50 of 67*

# MDL Functionality Demonstration



> **EN == 1 →** MDL acts as a delay line

> **EN == 0 →** MDL acts as a memory line

# Pulse Generator Functionality Demonstration

**Measured waveforms in 16x speedup mode**

© 2019 IEEE
International Solid-State Circuits Conference

*14.4:* **All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing**

*52 of 67*

# Speedup (1x-16x) Modes Demonstration

**Measured waveforms for image pixel value 214**

© 2019 IEEE
International Solid-State Circuits Conference

**14.4:** All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

*53 of 67*

# Speedup (1x-16x) Modes Demonstration

**Measured waveforms for image pixel value 214**

# Speedup (1x-16x) Modes Demonstration

**Measured waveforms for image pixel value 214**

14.4: All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

# Measured Accuracy @LeNet-5

*Hardware results for 100 test images*

**LeNet-5 on MNIST dataset**



V=650mV, Frequency=25MHz, T=25$^0$C
**Test chip measurements**

Classification Accuracy % (y-axis: 95, 96, 97, 98, 99, 100)

S/W (software)   1x   4x   8x   16x
Speedup Modes
**Weights: +1/-1**

- **Convolution Layers:** 8-bit fixed point inputs (Hardware/Test-chip)

- **FCN and Software implementation:** 16-bit floating point inputs and weights

**Weights (+1/-1):**   **Classification Accuracy**
- Software: 98.81%
- Test-chip:
  - ✓ 98% for speedup modes 1x-8x
  - ✓ 97% for speedup mode 16x

# Measured Accuracy @LeNet-5

*Hardware results for 100 test images*



**LeNet-5 on MNIST dataset**

V=650mV, Frequency=25MHz, T=25$^0$C
**Test chip measurements**

Classification Accuracy %

S/W (software) | 1x | 4x | 8x | 16x — **Speedup Modes**
**Weights: +1/-1**

S/W | 1x | 4x | 8x | 16x — **Speedup Modes**
**Weights: 0/1**

> **Convolution Layers:** 8-bit fixed point inputs (Hardware/Test-chip)

> **FCN and Software implementation:** 16-bit floating point inputs and weights

**Weights (+1/-1):**
> Software: 98.81%
> Test-chip:
> ✓ 98% for speedup modes 1x-8x
> ✓ 97% for speedup mode 16x

**Weights (0/+1):**
> Software: 97.62%
> Test-chip:
> ✓ 97% for speedup modes 1x-8x
> ✓ 96% for speedup mode 16x

**Classification Accuracy**

# Measured Accuracy *vs.* Voltage



> Operational down to **375mV** with more than **90%** accuracy in **1x** mode

> Classification accuracy of **97%** observed at voltages down till **537mV** in **16x** mode

# Measured Throughput *vs.* Voltage



- Throughput increases with voltage and speedup modes

- Peak Throughput for **C1 = 0.128 GOPS** @585mV

- Peak Throughput for **C3 = 0.38 GOPS** @585mV

© 2019 IEEE
International Solid-State Circuits Conference

**14.4:** All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

*59 of 67*

# Measured Energy Efficiency *vs.* Voltage



**Peak Energy Efficiency for C1 =
4.61 TOPS/W @496mV**

**Peak Energy Efficiency for C3 =
13.46 TOPS/W @496mV**

*14.4:* All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

# Simulation Results for AlexNet



- ➤ Dataset: Subset of ImageNet dataset (Classes - Cats and Dogs)

- ➤ Convolution and Pooling – 8 bit fixed point inputs and weights

- ➤ Fully Connected Network Layers – 16 bit floating point inputs and weights

- ➤ 13% accuracy drop observed in simulation
  - Multiple MDL lines are used for multi-bit weights dot product
  - Residue from all MDLs ➔ higher accuracy loss

# Performance Summary

| LeNet-5 Results/Metrics | Convolution Layer – C1 | | | | Convolution Layer – C3 | | | |
|---|---|---|---|---|---|---|---|---|
| Speedup Mode | 1x | 4x | 8x | 16x | 1x | 4x | 8x | 16x |
| Input clock frequency (MHz) | 24.0 | 24.0 | 24.0 | 24.0 | 24.0 | 24.0 | 24.0 | 24.0 |
| MAC clock frequency (MHz) | 0.19 | 0.75 | 1.50 | 3.00 | 0.19 | 0.75 | 1.50 | 3.00 |
| Convolution Cycle Time (us) | 149.3 | 37.33 | 18.67 | 9.33 | 842.67 | 210.67 | 105.33 | 52.67 |
| Operating Voltage (V) | 0.537 | 0.537 | 0.537 | 0.537 | 0.537 | 0.537 | 0.537 | 0.537 |
| Power (µW) | 28.67 | 28.67 | 28.67 | 28.67 | 30.17 | 30.17 | 30.17 | 30.17 |
| Throughput (GOPS) | 0.008 | 0.032 | 0.064 | 0.128 | 0.023 | 0.091 | 0.183 | 0.365 |
| Energy Efficiency (TOPS/W) | 0.29 | 1.16 | 2.33 | 4.65 | 0.76 | 3.02 | 6.04 | 12.08 |

➢ Peak Energy Efficiency: **12.08(4.65) TOPS/W** for C3(C1) layers @537mV

➢ Peak Throughput: **0.365(0.128) GOPS** for C3(C1) layers @537mV

# Outline

- Motivation

- Concept of Time-domain MAC

- Proposed CNN Engine Architecture

- Chip implementation and measurements

- Comparison with Prior Work

- Summary

© 2019 IEEE
International Solid-State Circuits Conference

**14.4:** All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

*63 of 67*

# Comparison with prior approaches

| Reference | Tech. (nm) | Circuit Type | Input/ Weight Size | Chip Size (mm²) | Pooling | Low Vcc Op. | Cap. or ADCs | Accuracy | Throughput (GOPS) | Power (µW) | Energy Efficiency (TOPS/W) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ISSCC'18 [1] | 65 | Analog | 6/1bits | 0.067 | No | No | Yes | 96.0% | 10.70 | 380.7 | 28.10 |
| ISSCC'18 [2] | 65 | Analog | 8bits | 1.44 | No | No | Yes | 96.0% | - | - | 3.125 |
| ISSCC'16 [3] | 65 | Digital | 16bits | 16.00 | Yes | No | No | 98.3% | 64 | 4.51E+4 | 1.42 |
| VLSI'16 [4] | 40 | Digital | 6/4bits | 2.40 | No | Yes | No | 98.0% | 102 | 3.90E+4 | 2.60 |
| CICC'17 [5] | 65 | Time | 8/3bits | 0.24 | No | Yes | Yes | 91.0% | 0.396 | 2.05E+4 | 0.019 |
| ISSCC'18 [6] | 55 | Time | 6/6bits | 3.125 | No | Yes | No | - | 2.152 | 690 | 3.12 |
| **This work (MDL CNN)** | **40** | **Time** | **8bits/ 1bit*** | **0.124** | **Yes** | **Yes** | **No** | **97.0%** | **0.365** | **30.17** | **12.08** |

*Scalable to multi-bit weights                    *References mentioned in paper*

# Conclusions

➢ Proposed **Bidirectional Memory Delay Line** for energy efficient time-domain MAC computation

➢ All-digital compact and **technology scaling friendly** design (no ADCs, DACs, frequency modulators)

➢ Low power design supporting **near-threshold voltage** operation and 16x performance boost with 4 input encoding modes

➢ Configurable MDL lengths with **on-chip pooling** and averaging operations

➢ Demonstrated the proposed time-domain CNN engine in **40nm CMOS** node achieving **12.08 TOPS/W** energy efficiency and **0.365 GOPS** throughput at 537mV

# Acknowledgements

➢ Authors would like to thank TSMC University shuttle program for the test- chip fabrication support.

➢ Authors would like to thank AMD for the financial support

➢ Authors also thank Vignesh Radhakrishnan and Jacob Rohan (graduate students in ECE, UT Austin) for helping with the test-chip measurements.

© 2019 IEEE
International Solid-State Circuits Conference

**14.4:** All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing

*66 of 67*

14.4: All-Digital Time-Domain CNN engine using Bi-directional Memory Delay Lines for Energy Efficient Edge Computing