会津大学
THE UNIVERSITY OF AIZU

# Prospects of the Nonvolatile FPGA and Its application to Edge-AI Accelerators

## Daisuke Suzuki

Adaptive Systems Laboratory
E-mail: daisuke@u-aizu.ac.jp

# Short Biography

Japanese Name: 鈴木 大輔

English Name: Daisuke Suzuki

Birthplace: Koriyama city, Fukushima

Ph. D: Engineering in Tohoku Univ.

Academic society: IEEE, IEICE, IPSJ

Research interests:
Nonvolatile logic circuit, nonvolatile FPGA, and their application to AI accelerators

Working experience: (Actually, I was at Hanyu & Natsui laboratory, Tohoku Univ.)

Research Associate - Center for Spintronics Integrated Systems, Tohoku University (2010 - 2014).

Assistant Professor - Center for Innovative Integrated Electronic Systems, Tohoku University (2014-2015).

Assistant Professor - Frontier Research Institute for Interdisciplinary Sciences, Tohoku University (2015-2020).

**Associate Professor - Computer Engineering Division, the University of Aizu (2020-).**
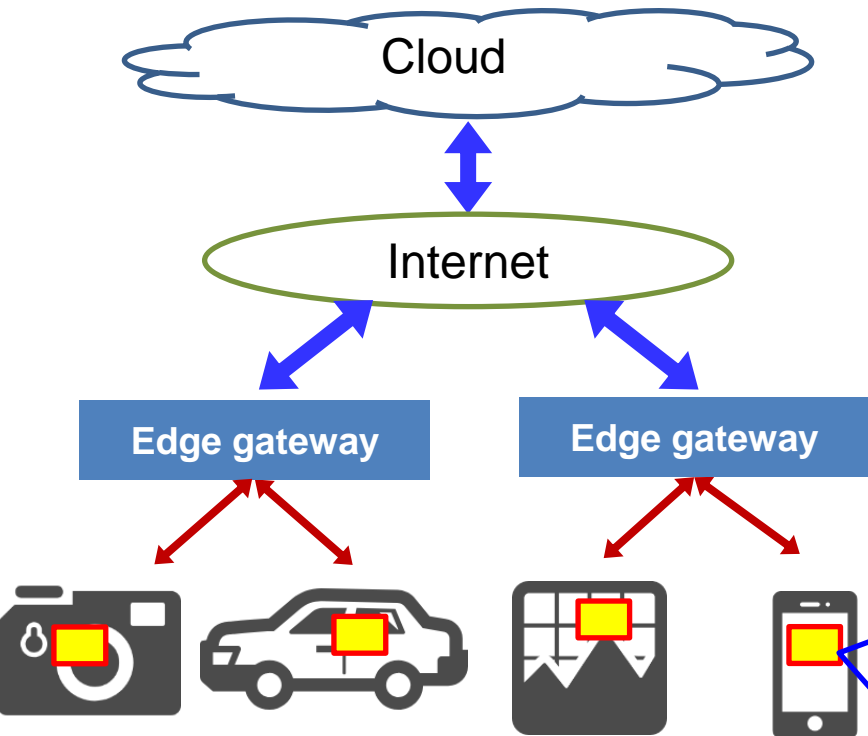
# Contents

**1. Introduction**

2. NV-FPGA and NV-LUT Circuit

3. Research Plan at UoA

4. Conclusion

## ■ Internet of Things (IoT)

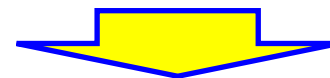Applications: Healthcare, sports, agriculture, smart house/city, automated driving, etc…)

Cloud

Internet

**Edge gateway**    **Edge gateway**

・Not cloud-centric data processing, but distributed data processing
・Recognition (Image, sound, etc.)

**Artificial intelligence (AI) on edge**

@ 1 Tops/W
> 30mJ/frame

ID 802564

ID 915342

ID 853468

Image recognition

*B.Moons et al. ISSCC 2017

1200mAh-1.5V Battery
< 2 Hour Operation*

[3] J. Yang et al., ISSCC 2019

・ >90B connections in 2025 [1]
・ > $2.7 trillion economic impact [2]

[1] https://www.idcjapan.co.jp/Press/Current/20180813Apr.html
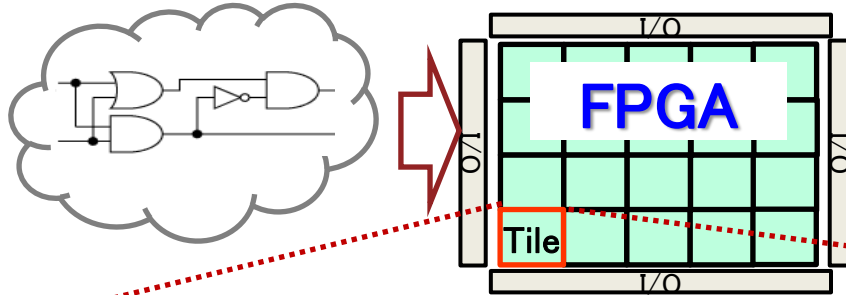[2] M. Mohammadi, et al., IEEE Communication Surveys & Tutorials, vol. 20, no. 4, pp. 2923-2960, 2018.
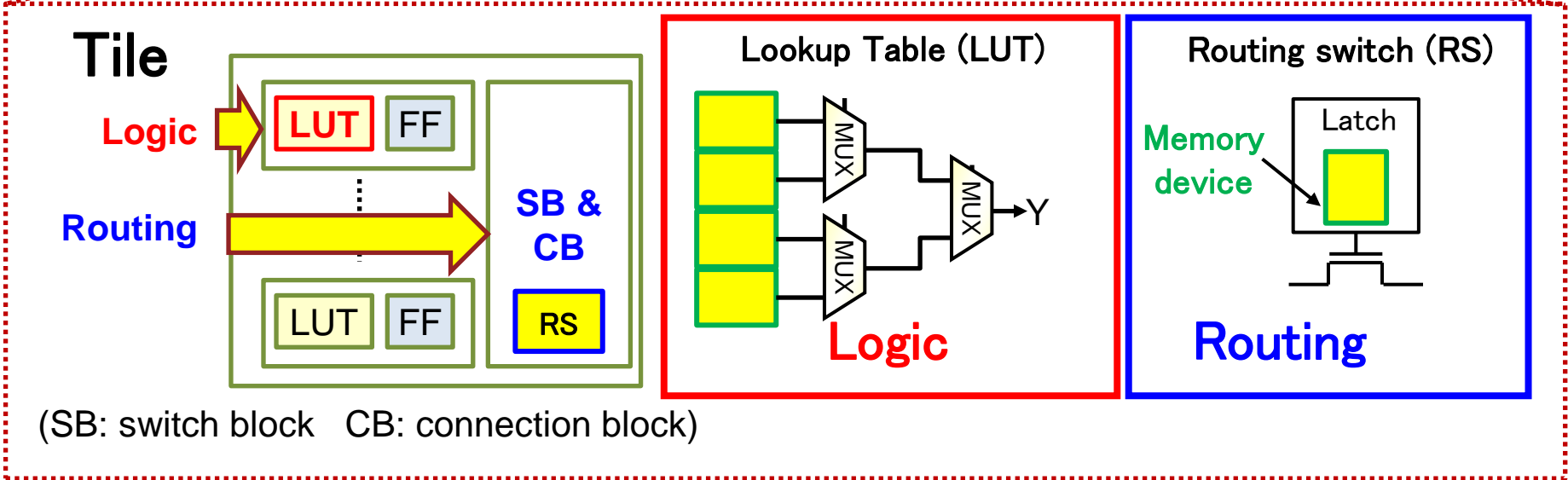
☹ **Strongly limited power supply**

## Low-power, energy-efficient edge-AI hardware is required.

**Circuit information**
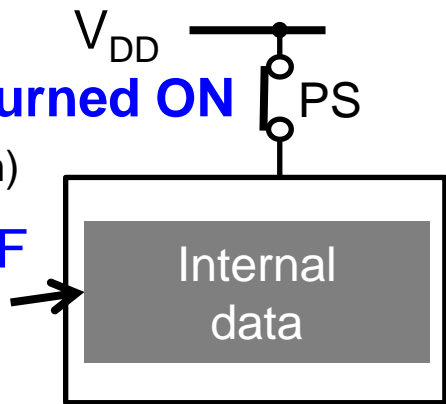
FPGA

Users can implement various digital circuits

**Tile**

Logic → LUT  FF

Routing → SB & CB  RS

(SB: switch block   CB: connection block)

**Lookup Table (LUT)**

MUX  MUX  MUX → Y

**Logic**

**Routing switch (RS)**

Memory device → Latch

**Routing**

**Merit: Short design time, flexibility, low design cost**
**Demerit: Large amount of standby power consumption**
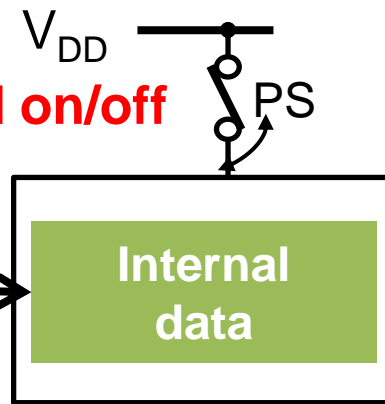
# Nonvolatile FPGA (NV-FPGA)

会津大学
THE UNIVERSITY OF AIZU

$V_{DD}$ — PS

**Always turned ON**

(PS: Power switch)

SRAM & CMOS FF
(**Volatile**)

→ Internal data

$V_{DD}$ — PS

**Quickly turned on/off**

**Nonvolatile memory & FF**

→ **Internal data**

(A: Active I: Idle)

I  A  I  A  I

Power | Dynamic

**Static**

Time

CMOS-only FPGA(volatile)

☹ **Always on to keep data**
**→ High standby power**

I  A  I  A  I

Power

Time

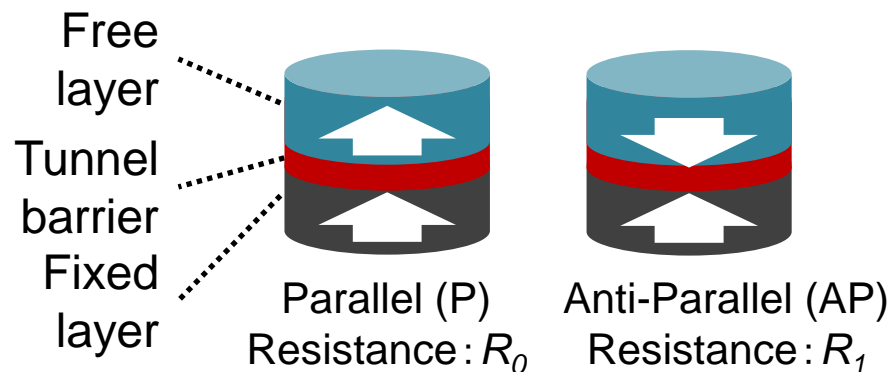**Nonvolatile FPGA**

☺ **Instant on/off w/o losing data**
**→ Ultra-low standby power**

## NV-FPGA -> Suitable for IoT/AI device

会津大学
THE UNIVERSITY OF AIZU

Free layer

Tunnel barrier

Fixed layer

Parallel (P)
Resistance：$R_0$

Anti-Parallel (AP)
Resistance：$R_1$

**Variable resistor**

D. Suzuki, et al.,
VLSI Circuits, 2009.

MTJ      MTJ

Metal

CMOS

300nm

Cross-sectional SEM image

$R_M$

M=1

$I_{W1}$

$I_{W0}$

$R_1$

M=0

$R_0$

$I_M$

R-I Characteristic

**One possible candidate for nonvolatile storage element**

# Summary of Research Roadmap

4th Meeting of the CE Division

# Contents

1. Introduction

**2. NV-FPGA and NV-LUT Circuit**
   - LIM-based LUT Circuit
   - Only-Once-Write Shifting
   - NV-FPGA-Embedded MCU

3. Research Plan at UoA

4. Conclusion

# Logic-In-Memory (LIM) Structure

Area

Buffer →

Buffer

**Overhead due to NV-storage function**

| SRAM-based | NVSRAM-based | LIM-based |
|---|---|---|
| Configuration memory | Configuration memory | SA/buffer |
| | | Configuration memory |
| MUX tree | MUX tree | MUX tree |

**NVSRAM-based:**

Buffer → D

X → MUX tree
6
64

| SA | SA | | SA |
|---|---|---|---|
| MTJ0 | MTJ1 | .... | MTJ63 |
| WT | WT | | WT |

Configuration memory

**LIM-based:**

SA/buffer → D

**Logic & memory**

X → MUX tree
6
64

| MTJ0 | MTJ1 | .... | MTJ63 | | WT |
|---|---|---|---|---|---|
| WT | WT | | WT | | |

Configuration memory

☺ Logic and memory functions are compactly merged.

SA: Sense amplifier   WT: Write transistor

## Compact circuitry by sharing circuit components.

☹ Process variation affects current levels
$I_L$ (logic 0), $I_H$ (logic 1), and $I_{REF}$ (reference current).



Redundant MTJ device

Variation compensation

Reference tree

MTJ layer

CMOS layer

Circuit information

$R_{r0}$

$R_{r1}$   $R_{r2}$   $R_{r3}$

Probability

$I_F=I_L$   $I_{REF}$   $I_F=I_H$

Error   Error

Current Level

**w/o redundant MTJ device**

Probability

$I_F=I_L$   $I_{REF}$   $I_F=I_H$

Current Level

**w/ redundant MTJ device**

## Variation resilient circuitry without area overhead

D. Suzuki, et al., VLSI Circuits, 2015.

## Area [$\mu m^2$]

SRAM: 1088
NVSRAM: 1344
-56%
Proposed: 480

## Standby power [nW]

SRAM: 400.0
NVSRAM: 22.0
-99%
Proposed: 3.0

(NVSRAM: Nonvolatile SRAM)

**Area and standby power reduction by LIM structure.**

[1] D. Suzuki et al., Jpn. J. Appl. Phys., **57,** 04FE09 (2018).

Data-shit function -> Key function of the **LUT circuit**



## Step 1

Output data stream $D_0$

OUT

**Memory cell (Active)**

**(1) Read**

Input data stream $D_6$, $D_5$, $D_4$

**(2) Write**

| M[3] $=D_3$ | M[2] $=D_2$ | M[1] $=D_1$ | M[0] $=D_0$ |

$A = 0$  2

Decoder (read/write)

(A: Address)

## Step 2

Output data stream $D_1$, $D_0$

OUT

**(1) Read**

**Active**

Input data stream $D_6$, $D_5$, $D_4$

**(2) Write**

| M[3] $=D_3$ | M[2] $=D_2$ | M[1] $=D_1$ | M[0] $=D_4$ |

$A = 1$  2

Decoder (read/write)

**Updated**

☺ **Number of write access per cycle is minimized to one.**

## **Write power reduction by minimizing # of write access**

# Conventional Data-Shift Function

■ SRAM-based LUT circuit



| # of write access per cycle (K-input LUT) | Conventional | Proposed |
|---|---|---|
| | $2^K$ | 1 |

**Proposed method is further fewer write access.
-> Low-write-power consumption**

会津大学
THE UNIVERSITY OF AIZU

Source: Renesas



**Power gap**

**Requirement for IoT sensor nodes powered by harvested energy**

- Simple terminal control
- Digital signal processing of sensor signals
- Trend analysis
- Recognition and authentication

Power consumption (μW): 1000000, 100000, 10000, 1000, 100, 10, 1

Operation frequency (MHz): 0, 100, 200

**Ultra-low-power/high-performance MCU is required for intelligent sensor node app.**

[Concept]

Replace sequential processing by CPU with **parallel processing by FPGA**

→ reduce processing time and increase the amount of standby state

→ further improve energy efficiency

Sequential processing with NV-CPU

**Parallel processing with FPGA-based accelerator**

Power overhead due to accelerator

**Power time reduction by accelerator**

**Total power reduction by PG & FPGA-based acceleration**

47.14μW Operation at 200MHz is achieved.

1. Introduction

2. NV-FPGA and NV-LUT Circuit

**3. Research Plan at UoA**

4. Conclusion

☹ In current situation, almost of NV-FPGA design is manual.
-> Establish design automation flow of the NV-FPGA

Cell library for the synthesizable NV-FPGA components (e. g. NV-FF)

Standard cell library

HDL code of the synthesizable NV-FPGA

Logic simulation

Synthesis

Gate-level simulation

Place & route

Post-layout simulation

**Design automation of NV-FPGA (CRF2020)**

Standard EDA tool flow for implementing the synthesizable NV-FPGA

Digital circuit RTL →

**Tool for implementing digital circuit on the NV-FPGA**
JSPS: KAKENHI
20K11725

Configuration bitstream

I/O

NV−FPGA chip

I/O I/O

Tile

I/O

会津大学
THE UNIVERSITY OF AIZU

**Example: Binary Convolutional Neural Network (BCNN) [1]**

[1] M. Courbariaux et al., arXiv:1602.02830, 2016.



-> **dog**
-> cat
-> lion
-> bird

Convolution    Pooling    Convolution    Pooling    Fully-connected



$\otimes$    K    ...    ...    =    R    C    M    N

Input feature map (ifm)    M weights    Output feature map (ofm)

```
for (to=0; to<M; to++)
  for (row=0; row<R; row++)
    for (col=0; col<C; col++)
      for (ti=0; ti<N; ti++)
        for (i=0; j<K; i++)
          for (j=0; j<K; j++)
            ofm[to][row][col] += weight[to][ti][row+i][col+j] * ifm[ti][row+i][col+j];
```
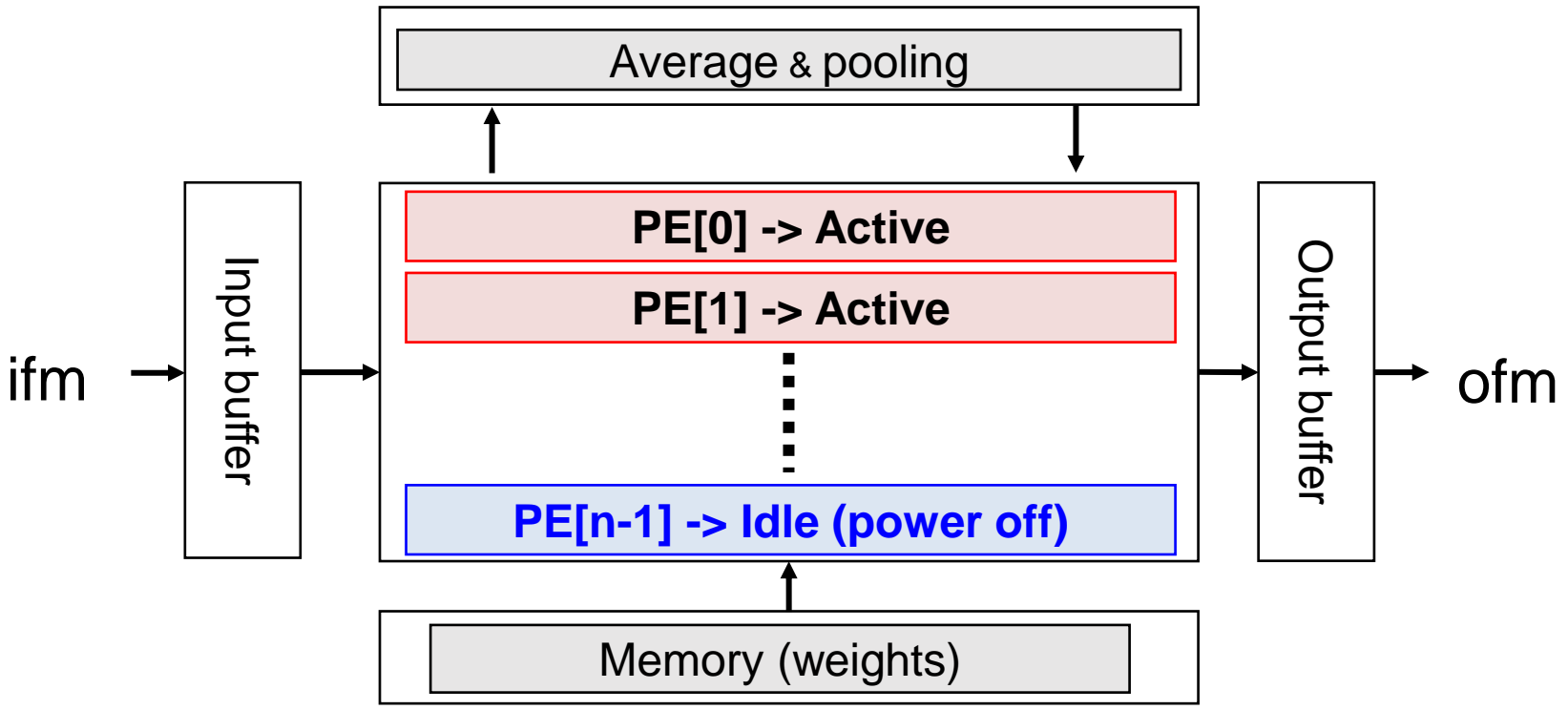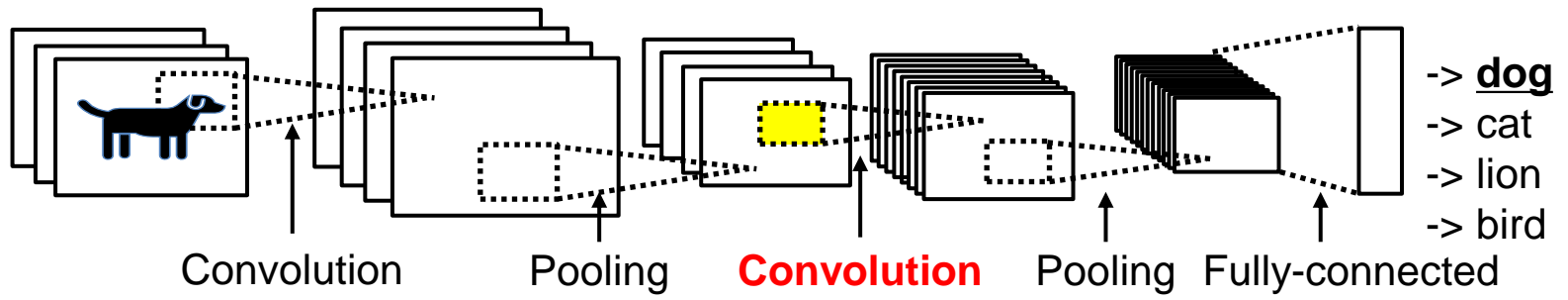
**Binary information**

**SUM -> Bit-count**    **Multiplication -> XNOR**

**Massively parallel computing is required.**

Convolution    Pooling    **Convolution**    Pooling    Fully-connected

-> **dog**
-> cat
-> lion
-> bird

| Average & pooling |
| --- |

Input buffer

ifm →

PE[0] -> Active

PE[1] -> Active

PE[n-1] -> Idle (power off)

Output buffer

→ ofm

Memory (weights)

**Massively parallel architecture with no wasted standby power consumption**

[with Prof. Ben]
Competitive Research Fund 2020 in UoA,
`` Development of an Energy-efficient
Heterogeneous Spiking Neuro-inspired System
for Deep Neural Networks."

[with Prof. Saito]
IoT/AI Device Cluster

# Contents

1. Introduction

2. NV-FPGA and NV-LUT Circuit

3. Research Plan at UoA

**4. Conclusion**

# Conclusion

## NV-FPGA and NV-LUT Circuit

➢ Compact & variation resilient circuity by using LIM structure

➢ Low-power data shifting by using only-once-write shifting

➢ Energy-efficient NV-MCU chip by embedding NV-FPGA


## Research Plan

➢ Establish EDA Tool Flow for NV-FPGA

➢ Design NV-FPGA-based AI Accelerators

➢ Collaborations with UoA members (Prof. Ben, Prof. Saito, etc.)