# Research Introduction in Ceres Lab

Cerebral and Reliable NoC-based AI Accelerator Design and Applications for Anomaly Detection in Smart Motor Systems

## Kun-Chih (Jimmy) Chen

Associate Professor/ Electric Junior Chair Professor,
Dep. Electronics and Electrical Engineering/ Institute of Electronics,
National Yang Ming Chiao Tung University (NYCU)
Email: kcchen@nycu.edu.tw
URL:  https://sites.google.com/site/cereslaben

**NYCU CERES LAB**

**Kun-Chih (Jimmy) Chen**,
Associate Professor/ Electric Junior Chair Professor
Institute of Electronics, National Yang Ming Chiao Tung University
Email: kcchen@nycu.edu.tw
Website: https://sites.google.com/site/cereslaben/advisor

## ❖ Specialty

- ❖ Multi-core System on Chip (MPSoC) design
- ❖ Neural network model and accelerator design
- ❖ Reliable system design
- ❖ VLSI CAD design
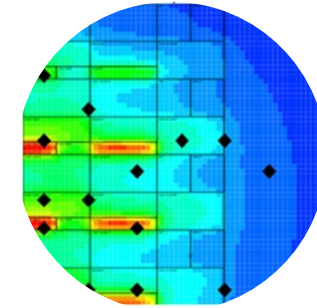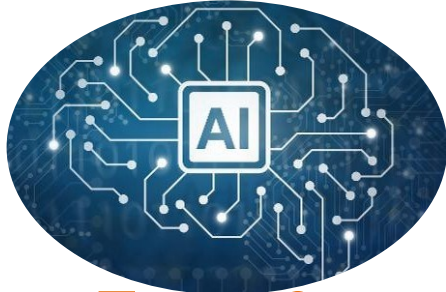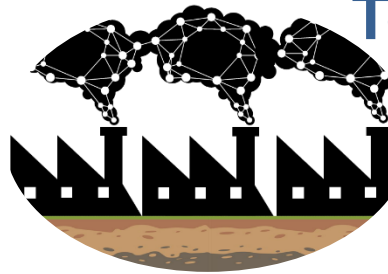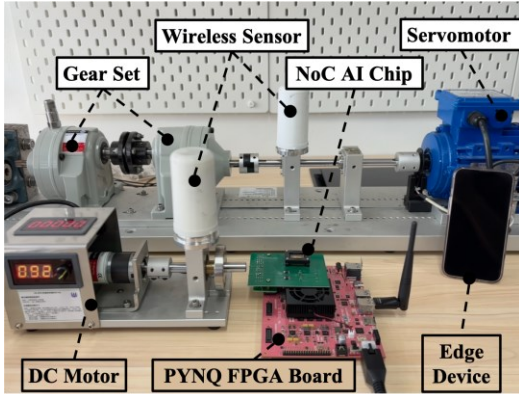- ❖ Smart Manufacturing

## ❖ Feature Honors

- ❖ Dr. Da-You Wu Memorial Award of NSTC
- ❖ IEEE TVLSI Best Paper Award
- ❖ IEEE CASS Continuing Education Featuring Selected Conference Tutorial
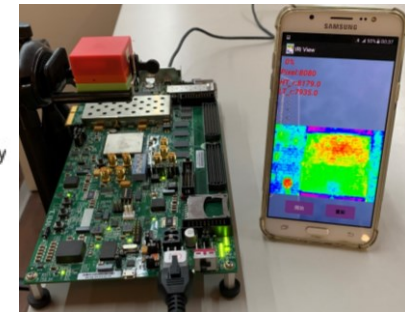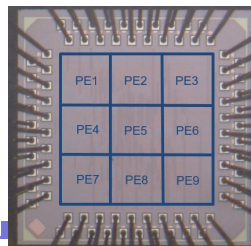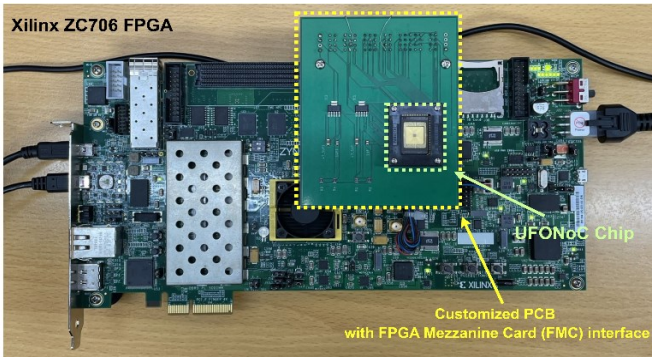- ❖ Taiwan IC Design Society Outstanding Young Scholar Award



IEEE Circuits and Systems Society takes pleasure in presenting the 2024 IEEE Transactions on Very Large-Scale Integration Systems Best Paper Award to Kun-Chih (Jimmy) Chen. For the paper "Adaptive Machine Learning-Based Proactive Thermal Management for NoC Systems" IEEE Transactions on Very Large Scale Integration (VLSI) Systems (Volume 31, Issue 8, August 2023)



IEEE Circuits and Systems Society takes pleasure in presenting the Certificate of Recognition for Contributions to CASS Continuing Education to Kun-Chih Chen. For the tutorial on Network on Chip (NoC)-based Deep Neural Network Design Framework: From Algorithms to Architectures (MWSCAS 2024)

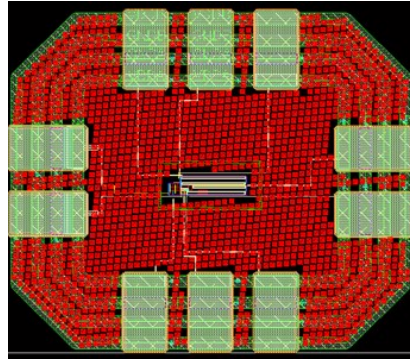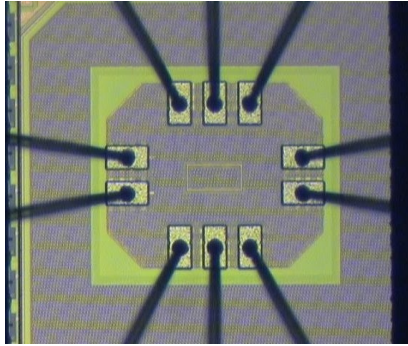# Topic 1: Smart Manufactoring

# Topic 3: Reconfigurable Neural Network

# Topic 2: Smart Thermal Magagement on MPSoC

# Chip Gallery

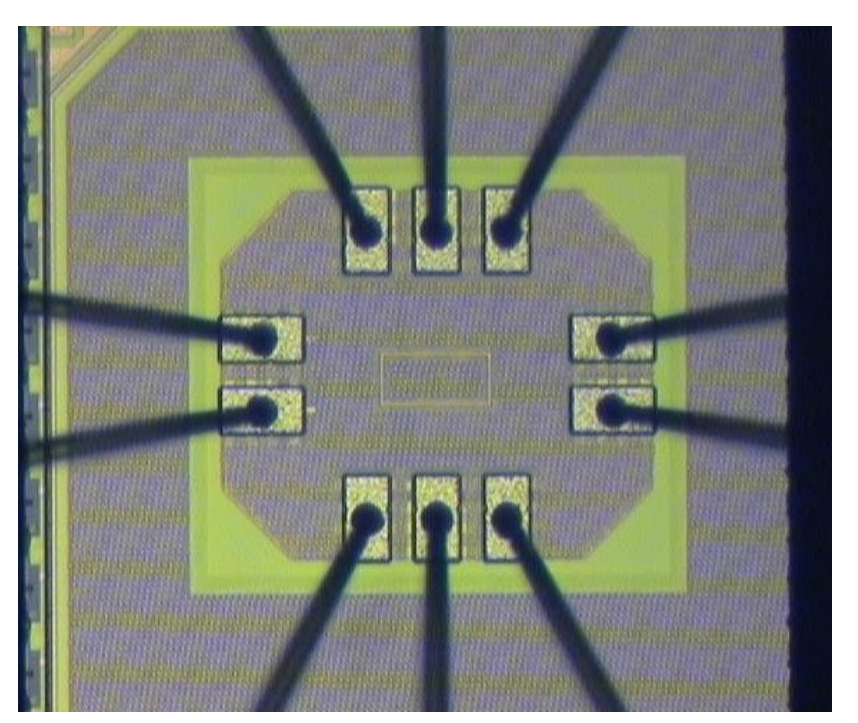

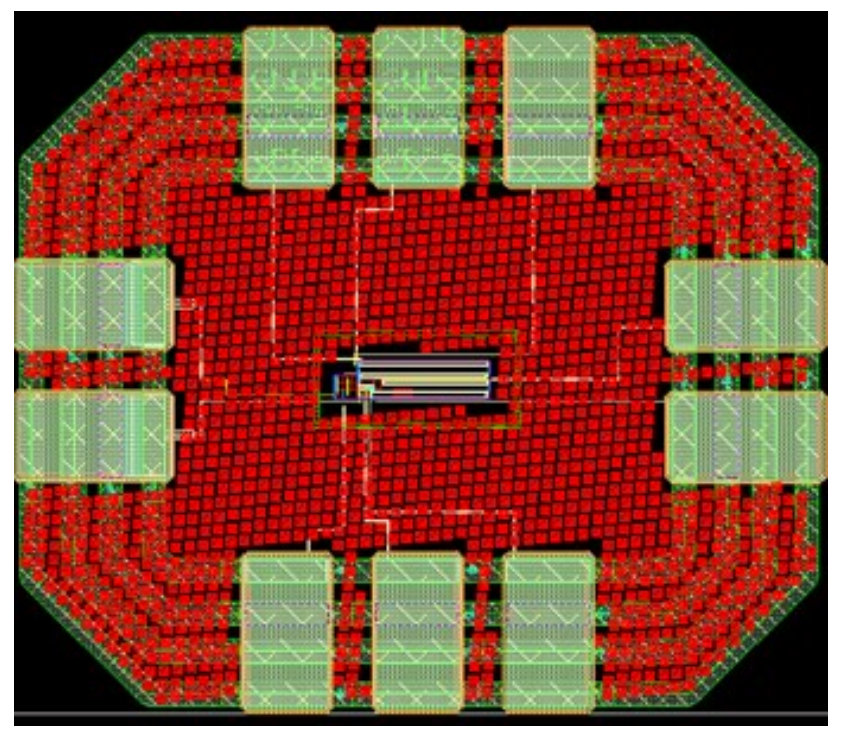

## All-digital temperature sensor

- Technology: TSMC 90nm
- Size: 0.0016 mm$^2$
- Power: 798mW
- Clock frequency: 5MHz

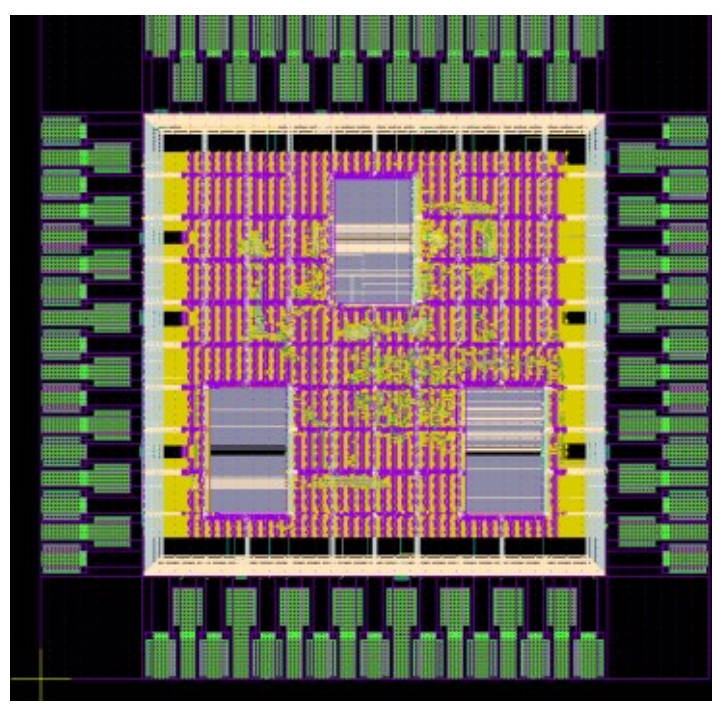

## AI-sonar for geological analysis

- Technology: TSMC 40nm
- Size: 1.56 mm$^2$
- Power: 25.5mW
- Clock frequency: 100MHz

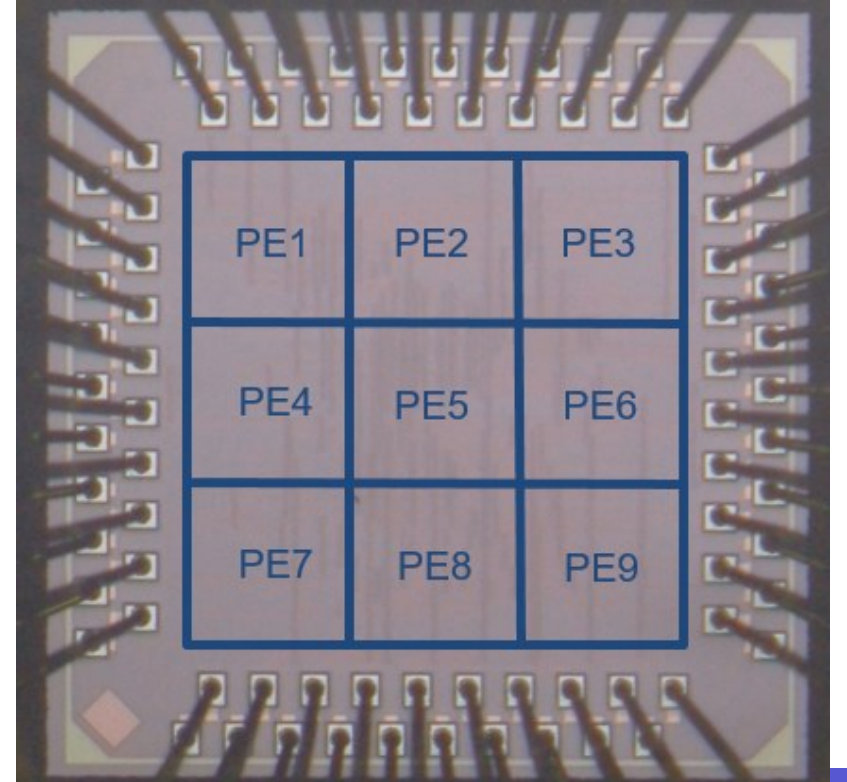


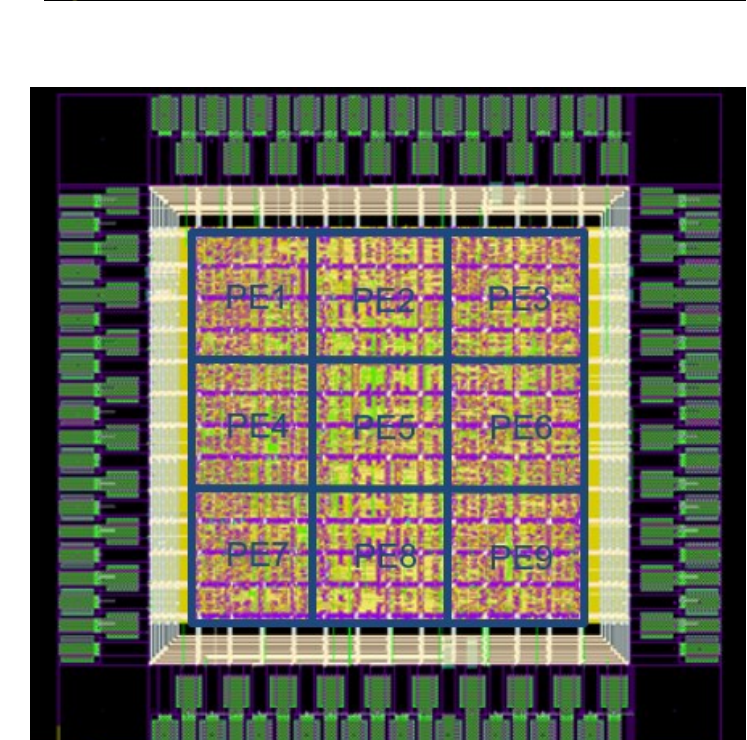

## NoC-based reconfigurable DNN

- Technology: TSMC 40nm
- Size: 0.84 mm$^2$
- Power: 10.37mW
- Clock frequency: 105MHz

# About my school
# National Yang Ming Chiao Tung University (1/2)

❖ NCTU was established in Shanghai in 1896 originally and re-established in Taiwan in 1958.

❖ Electrics Institute is the first institute when NCTU re-established in Taiwan in 1958, and the Taiwan's semiconductor and space industry were born from NCTU.

    ❖ The first wafer in Taiwan (1964)

    ❖ The first Bipolar Transistors in Taiwan (1965)

    ❖ The first IBM computing system in Taiwan (1968)

    ❖ The first hybrid rocket in Taiwan (2010)

    ❖ The first sounding rocket in Taiwan (2014)

# About my school
# National Yang Ming Chiao Tung University (2/2)

❖ In 2021, NCTU merged with a prestigious medical university, National Yang Ming University, and rebranded as National Yang Ming Chiao Tung University (NYCU).

Medical, Bio-technology     Engineering, Science, Administration

# NYCU at a Glance

**nycu**

**21,703** Students

(1,300 Overseas Students)

8,612  Undergraduates

13,091  Graduates

**2,454**  Faculties

(135 International Faculties)

1,154   Full-time faculties

958   Part-time faculties

89   Research Staff

253   Staff

# Ecosystem Neighboring NYCU

NYCU

Surrounded by
**Science Park and National Institutes**

【光復校區 Guangfu Campus 】

the first transistor in Taiwan

National Measurement Laboratory (NML)

National Center for High-performance Computing (NCHC)

Taiwan Semiconductor Research Institute (TSRI) (Consolidation of CIC & NDL)

National Synchrotron Radiation Research Center (NSRRC)

Taiwan Instrument Research Institute (TIRI)

Taiwan Space Agency (TASA)

**Hsinchu Science Park**

National Chiao Tung University

# NYCU-EE is the largest EE department in Taiwan

❖ The best EE department in Taiwan

   ❖ World ranking: 39; Taiwan ranking: 1

❖ IEEE Fellows: 20; IET Fellows: 2; NAE Academician: 3; NAI Fellows: 2

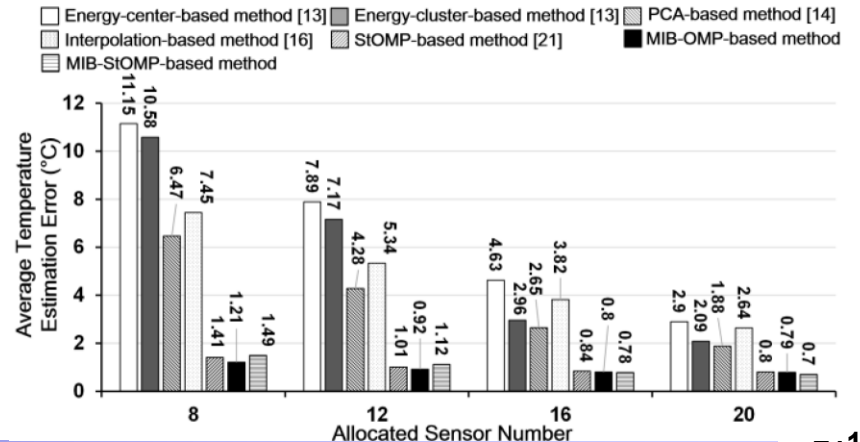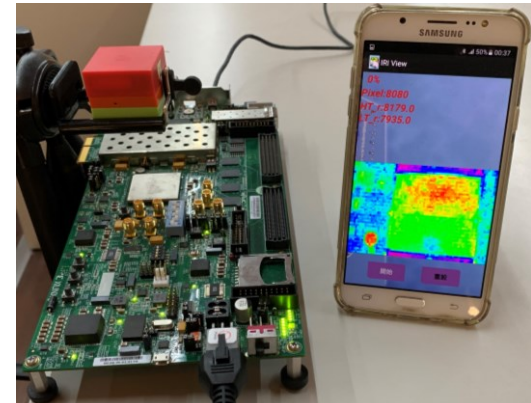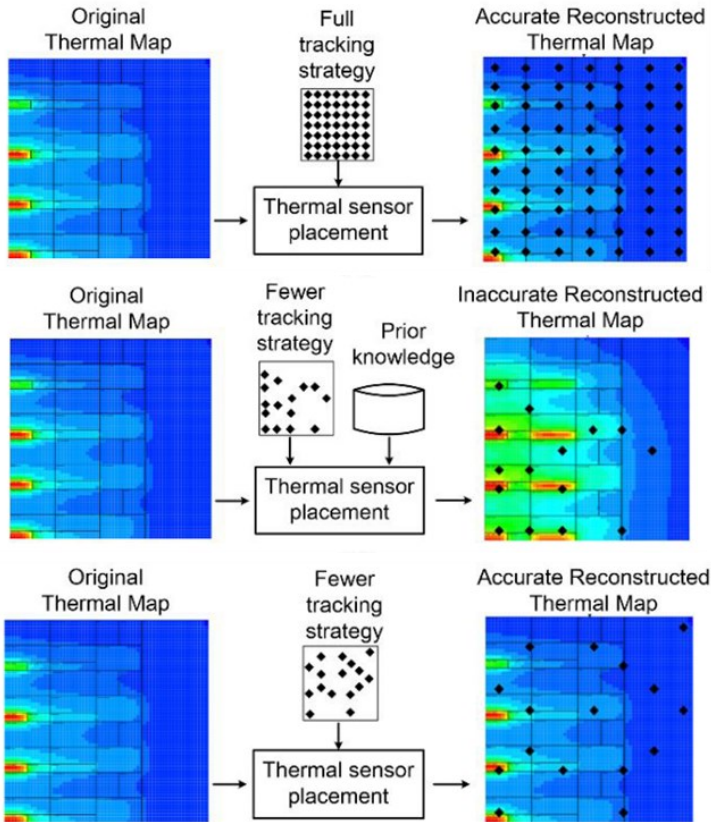# Research Pillar 1:

# Smart Thermal Management on MPSoC

# CS-based Low-cost Thermal Sensor Placement

❖ **<u>Features</u>** (IEEE TCAD 2022)

- Adopting Compressive Sensing (CS) to achieve fast sensor placement
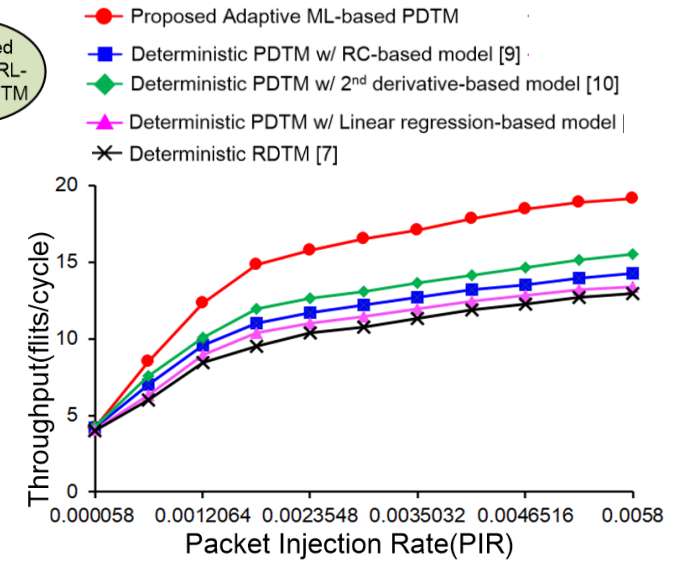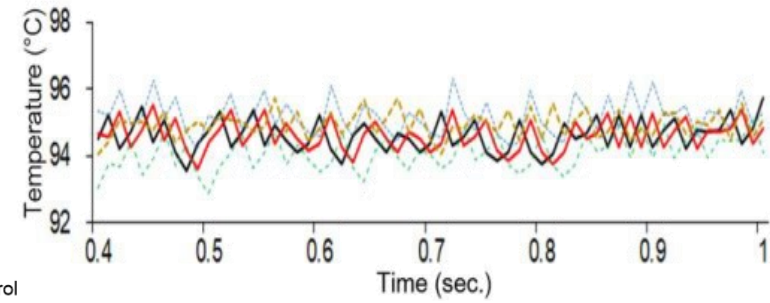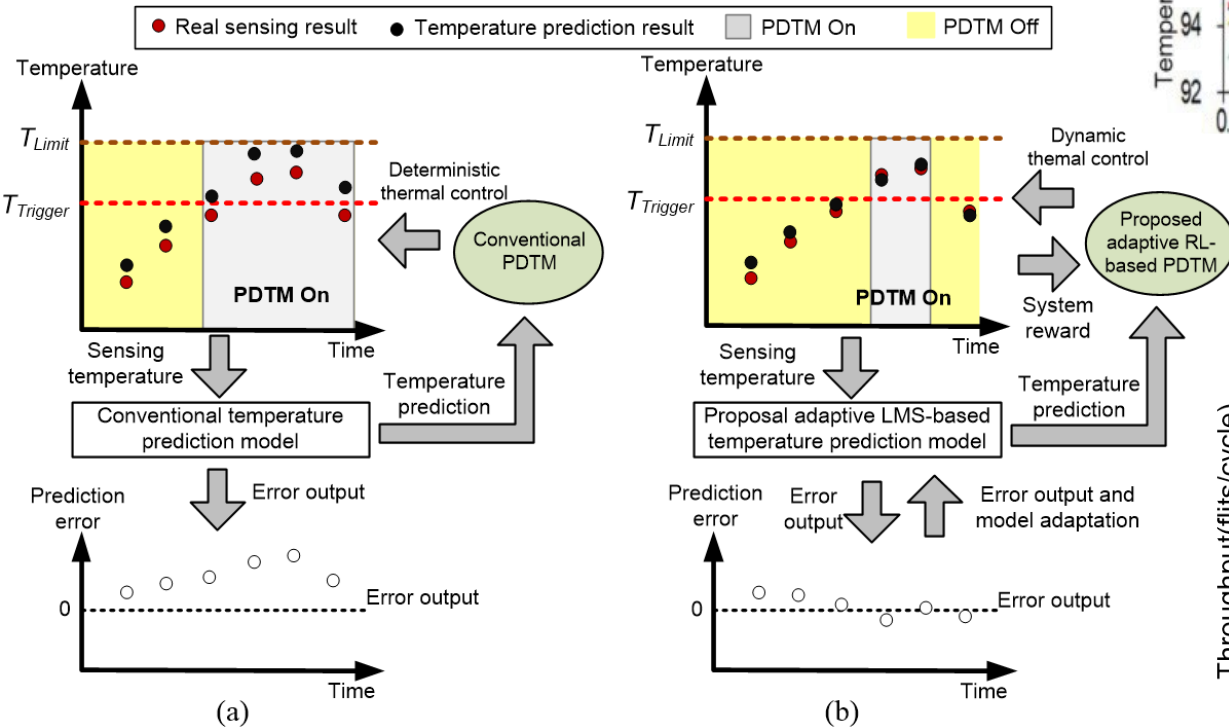- Proposing a novel temperature reconstruction method to build the temperature distribution with low computing cost

# Adaptive ML Method for Proactive Thermal Management

❖ **Features** (IEEE TVLSI 2024 Best Paper; ISCAS 2020 Best Student Paper)
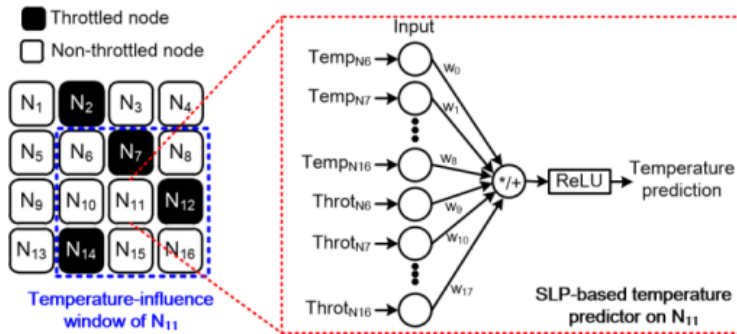
− Adopting online learning to predict the on-chip temperature precisely

− Adopting adaptive reinforcement learning to fine-grained control the system temperature

# Adaptive Single Layer Perception (ASLP)-based Temperature Prediction

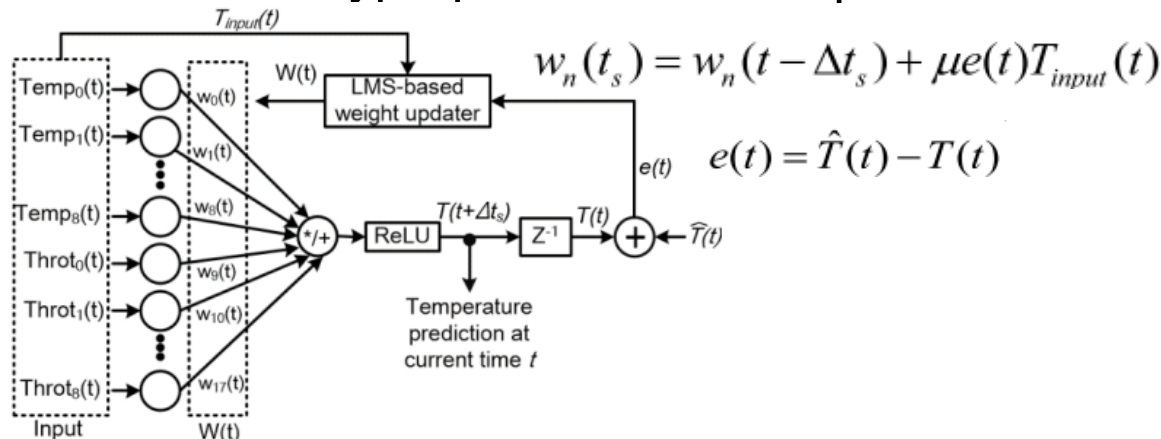❖ Adopt an SLP to re-model the temperature prediction equation



$$T(k + N \mid k) = B^N T(k \mid k) + \sum_{i=k}^{q} B^{q-i} Du(i)$$

$$T(k + N \mid k) = \sum_{n=1}^{k-q+2} w_n \cdot X_n$$
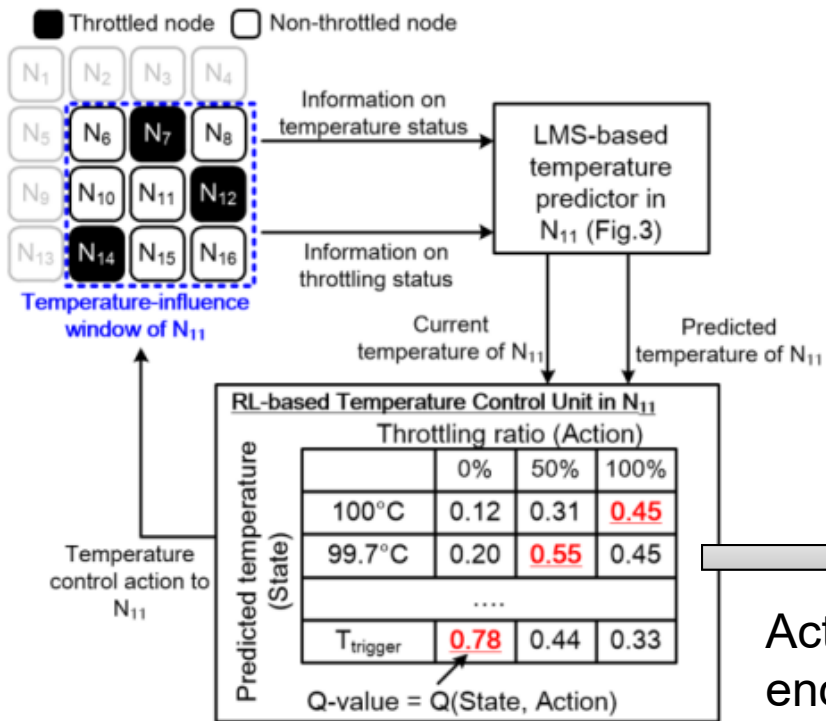
❖ Adopt the LMS-based adaptive filter to update the parameters ($w_n$) at runtime to fit the hyperplane of the temperature behavior



$$w_n(t_s) = w_n(t - \Delta t_s) + \mu e(t) T_{input}(t)$$
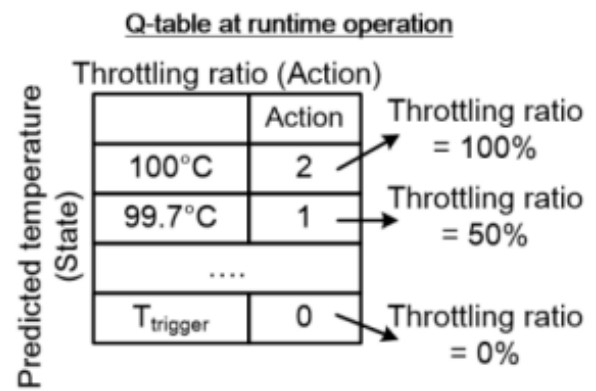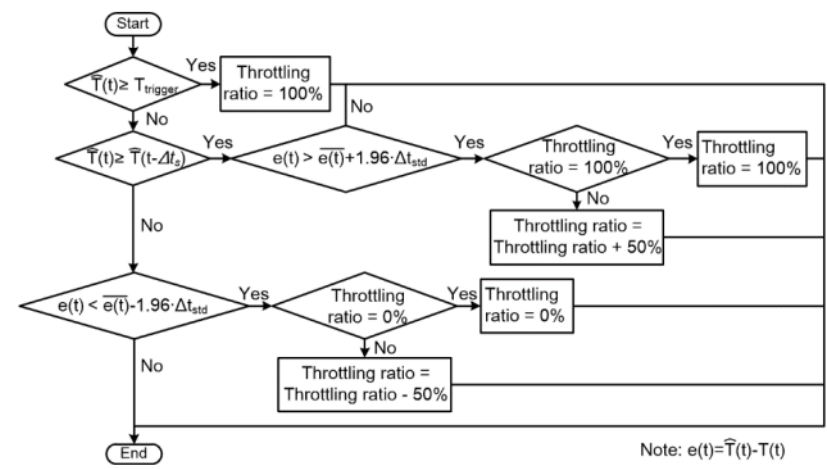
$$e(t) = \hat{T}(t) - T(t)$$

# Adaptive Reinforcement Learning-based Temperature Control Mechanism

❖ Adopt the Q-learning the select the proper action to throttle the thermal-emergency NoC nodes.

# Experimental Results

# Research Pillar 2:

# Reconfigurable Neural Network design

# Lego-based DNNoC Design Paradigm

❖ **<u>Features</u>** (IEEE JETCAS 2021)

− Adopting the Neu-Lego design to mitigate the analysis complexity

− Adopting flexible Network on Chip (NoC) interconnection to reduce interconnection complexity and reduce time-to-market

# *DNNoC* Construction (1/2) : Model Analysis

❖ Obtains the required number of *NeuLego* PEs and model information from the given DNN model.

❖ Each *NeuLego* PE is used to process data from the same data dimension.

❖ Facilitates data exchange.

# *DNNoC* Construction (1/2) : Model Analysis

❖ Obtains the required number of *NeuLego* PEs and model information from the given DNN model.

❖ Each *NeuLego* PE is used to process data from the same data dimension.

  ❖ Facilitates data exchange.



**Convolution layer**
Channel: 3
Kernel size: 5x5

**Max Pooling layer**
Channel: 3
Kernel size: 2x2

**Convolution layer**
Channel: 3
Kernel size: 4x4

**Max Pooling layer**
Channel: 3
Kernel size: 2x2

**Dense layer**
Neurons: (1) 128 (2) 32 (3) 10

**Model analysis**

**The required NeuLego PEs**

# *DNNoC* Construction (2/2) : *DNNoC* Construction Flow and Lego Placement

❖ We propose to share the computing resources and find the proper number of *NeuLego* PEs for a given DNN model.

    ❖ Improve hardware efficiency



**The required NeuLego PEs**

**The constructed *DNNoC***

**The *NeuLego* PE pool**

**Mesh-NoC size : 3**

**Row-based PE alignment**

# *DNNoC* Construction (2/2) :  *DNNoC* Construction Flow and Lego Placement

❖ We propose to share the computing resources and find the proper number of *NeuLego* PEs for a given DNN model.

   ❖ Improve hardware efficiency



**The required NeuLego PEs**

**The constructed *DNNoC***

**The *NeuLego* PE pool**

**Mesh-NoC size : 3**

**Row-based PE alignment**

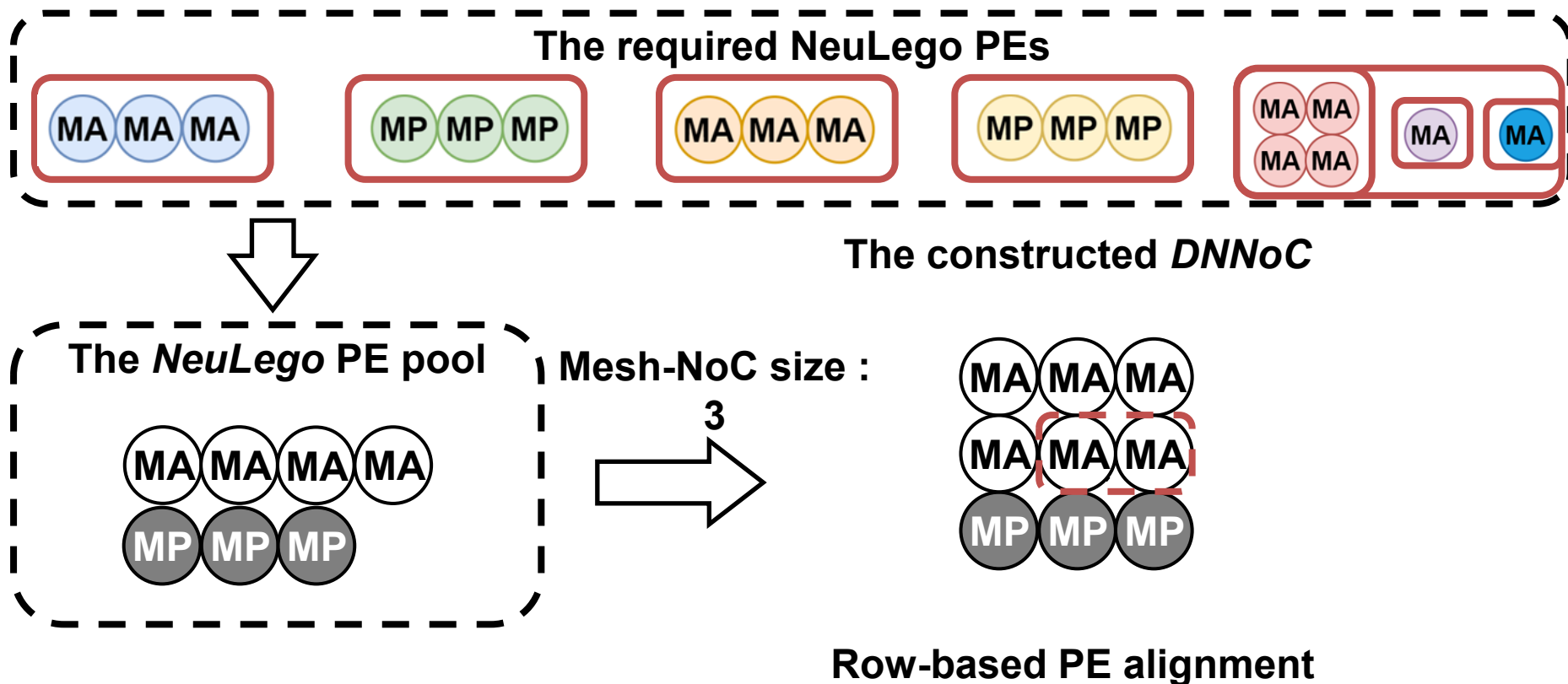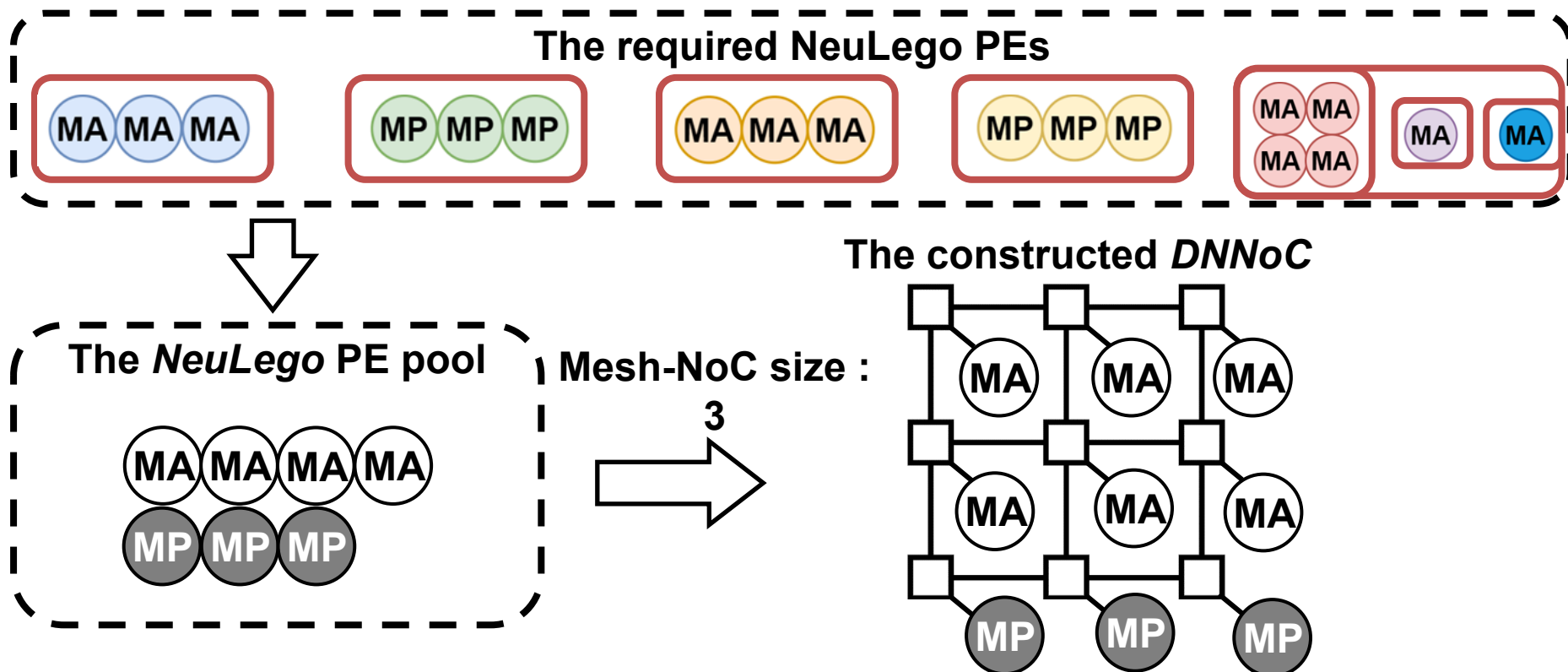# *DNNoC* Construction (2/2) : *DNNoC* Construction Flow and Lego Placement

❖ We propose to share the computing resources and find the proper number of *NeuLego* PEs for a given DNN model.

    ❖ Improve hardware efficiency

# *DNNoC* Execution

❖ **Layer-wise dynamic mapping algorithm**
  ❖ The available computing resources of the constructed *DNNoC* would be sufficient to fit the largest layer of the target DNN model.
  ❖ Maps a large-scale DNN model to the source-limited *DNNoC* platform.
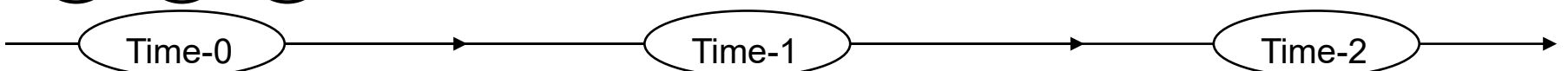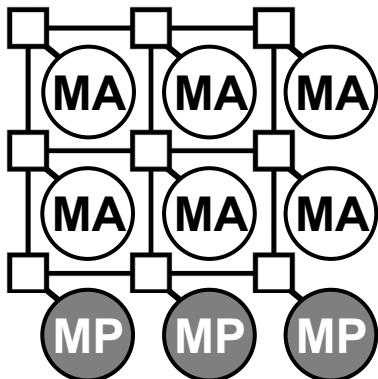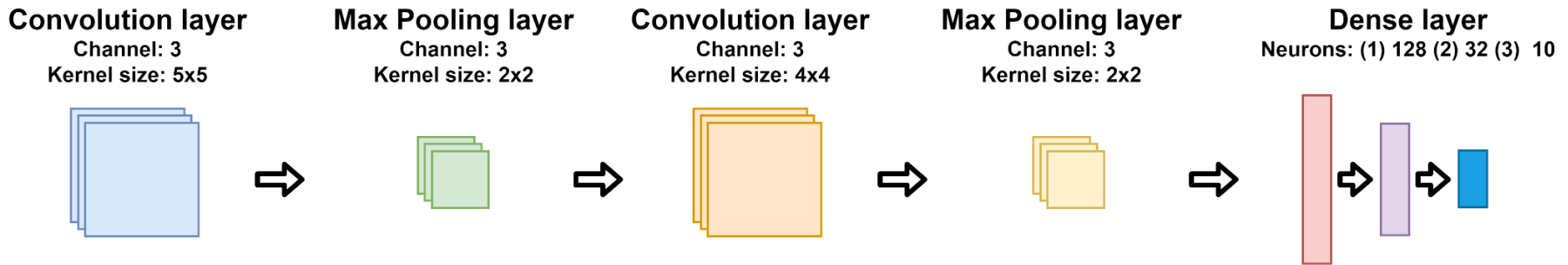
# *DNNoC* Execution

❖ **Layer-wise dynamic mapping algorithm**

   ❖ The available computing resources of the constructed *DNNoC* would be sufficient to fit the largest layer of the target DNN model.

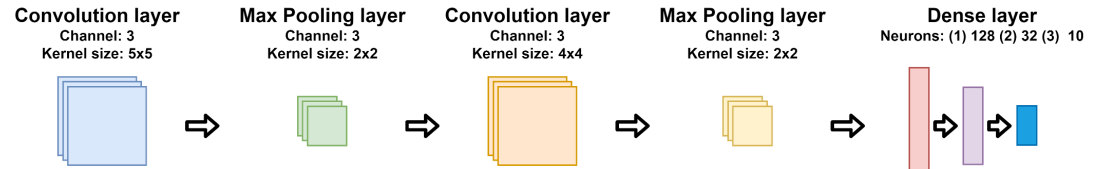   ❖ Maps a large-scale DNN model to the source-limited *DNNoC* platform.

# The first NoC-based Reconfigurable DNN Accelerator with AXI communication protocol in Taiwan

❖ **Features** (VLSI CAD Symposium Best Paper Award)

- Support arbitrary kernel size and shape to compute the convolution operations

- Adopting flexible Network on Chip (NoC) interconnection to reduce interconnection complexity and reduce time-to-market

- Adopt AXI4-stream communication protocol

| Technology | | TSMC 40nm |
|---|---|---|
| Area (mm$^2$) | Chip | 1.4 x 1.4 |
| | Core | 0.84 x 0.84 |
| | Gate count | 6,871k |
| IO/Core VDD (v) | | 2.5/0.9 |
| Clock freq (MHz) | | 105 |
| Power (mW) | | 10.3672 |
| Throughput (GOPS) | | 143.5 |

# Design Challenge of DNN Accelerator:
# Various kernel size

❖ The convolutional kernel sizes are usually not fixed in the DNN model.

  ❖ Worst-case design consideration

❖ The register size of processing element (PE) is usually based on the largest kernel size in the target model.

  ❖ Low utilization of PE computational capability.

  ❖ Cannot process the operation.

| DNN model | Kernel size/shape |
|---|---|
| AlexNet | 3x3, 5x5, 11x11 |
| GoogLeNet | 1x1, 3x3, 5x5, 7x7 |
| DeepSpeech2 | 21x11, 41x11 |

# Channel-wise Convolution Operation

❖ The channel-wise convolution operation.
 ❖ PE is low applicable for arbitrary kernel size.
 ❖ PE generates channel partial sum (*CPsum*).

$$CPsum_{(i,j,c)} = \sum_{m=1}^{h}\sum_{n=1}^{h}\left(I_{(i+m-1,j+n-1,c)} \times W_{(m,n,c)}\right)$$

$$OFmap_{(i,j)} = \sum_{c=1}^{d}\left(CPsum_{(i,j,c)}\right)$$

# Weight-wise NN Processing Mechanism (1/2)

❖ We can exploit the shape parameters to infer all Input Feature Map Data (*IFD)* that the weight will convolve with.

❖ PE will process one weight with corresponding inputs and accumulate operation partial sum (*OPsum*).

$$OPsum_{(i,j,m,n,c)} = I_{(i+m-1,j+n-1,c)} \times W_{(m,n,c)}$$

$$OFmap_{(i,j)} = \sum_{c=1}^{d}\sum_{m=1}^{h}\sum_{n=1}^{h}\left(OPsum_{(i,j,m,n,c)}\right)$$

# Weight-wise NN Processing Mechanism (2/2)

❖ Computing data register (*CD_REG*).

❖ Scaling factor register (*SF_REG*).
 ❖ SF register size will not be restricted.
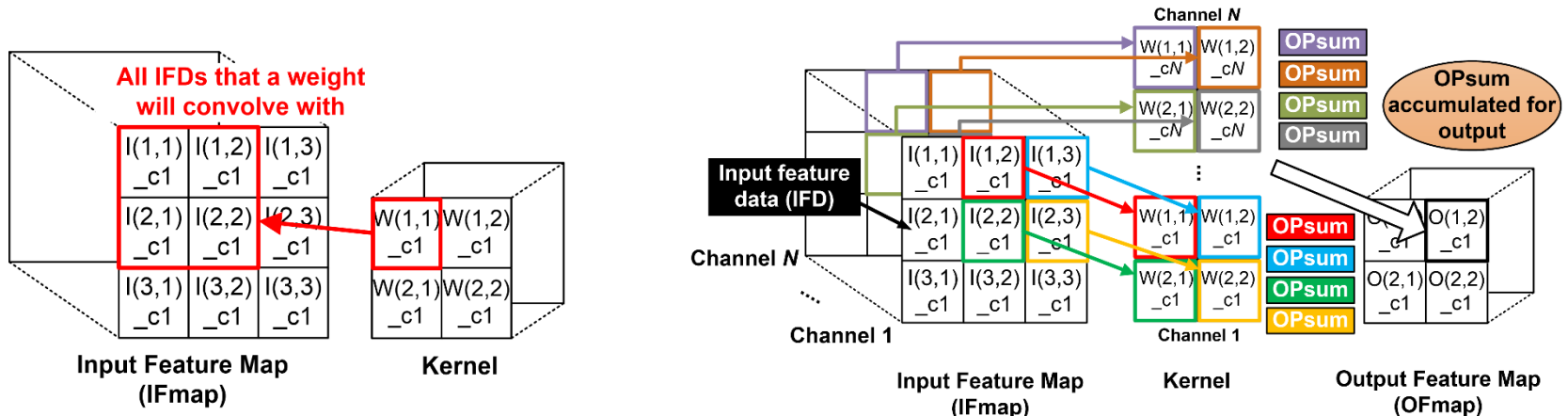
❖ Reduce memory access.

| I1 | I2 | I3 |
|----|----|----|
| I4 | I5 | I6 |
| I7 | I8 | I9 |

$*$

| W1 | W2 |
|----|----|
| W3 | W4 |

$=$

| O1 | O2 |
|----|----|
| O3 | O4 |

$$O_1 = (I_1 \times W_1) + (I_2 \times W_2) + (I_4 \times W_3) + (I_5 \times W_4)$$
$$O_2 = (I_2 \times W_1) + (I_3 \times W_2) + (I_5 \times W_3) + (I_6 \times W_4)$$
$$O_3 = (I_4 \times W_1) + (I_5 \times W_2) + (I_7 \times W_3) + (I_8 \times W_4)$$
$$O_4 = (I_5 \times W_1) + (I_6 \times W_2) + (I_8 \times W_3) + (I_9 \times W_4)$$

*OPsum*

**CD_REG**

W2

**OPsum_REG**

O1 O2
O3 O4
**Output**

X  +

I2  I3

I5  I6

**SF_REG**

# Hybrid Data Reuse Method by Using NoC

❖ After accessing the data from on-chip memory once, PE will share duplicated data through packet transmission.

  ❖ Does not need to design complicated dataflow.
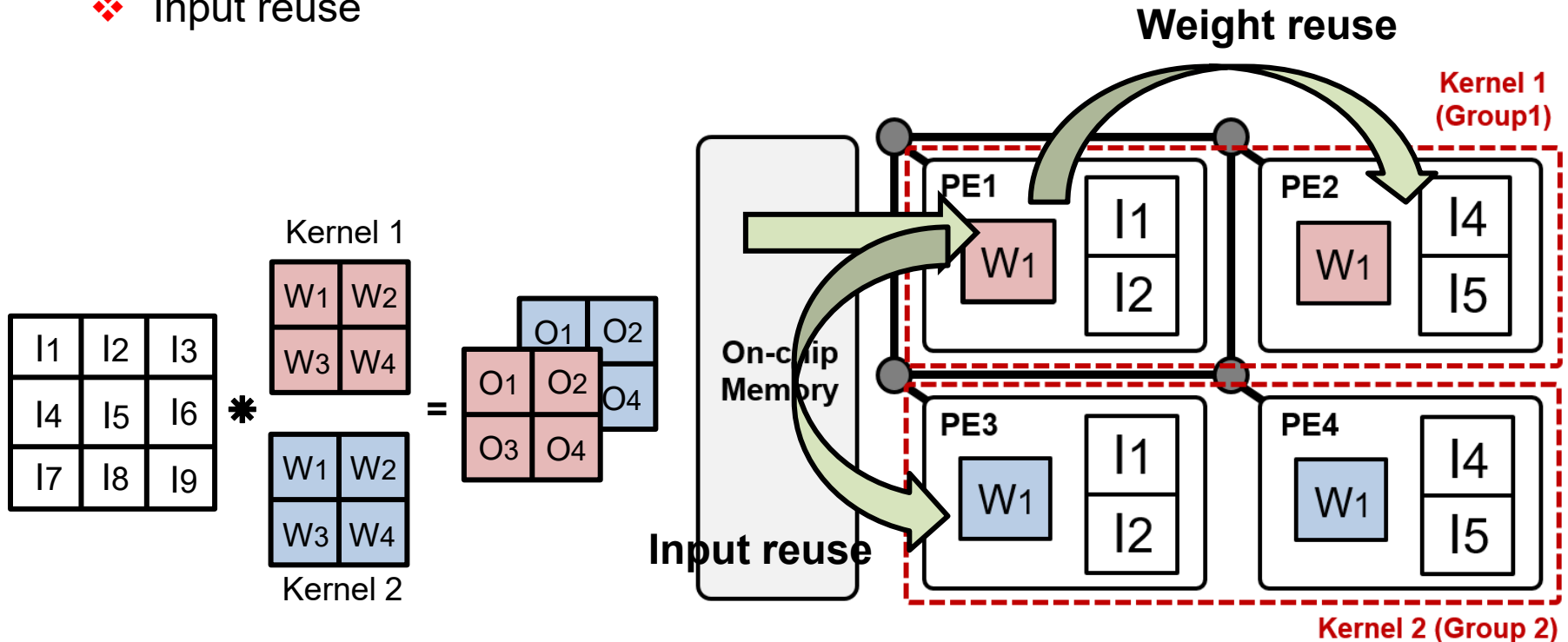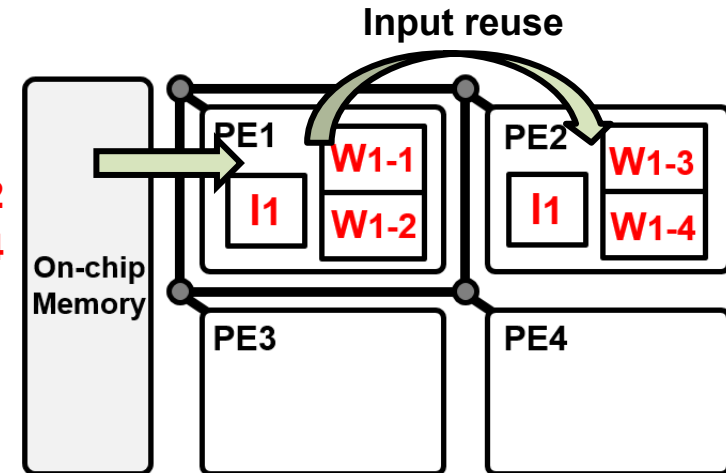
  ❖ Weight reuse

  ❖ Input reuse

# The Processing of Fully-connected Layer

❖ The proposed processing mechanism can also be applied in the fully-connected layer.

  ❖ Store the input data to the *CD_REG*, and the corresponding weights to the *SF_REG*, respectively.

  ❖ Input reuse

$$O_1 = (I_1 \times W_{1-1}) + (I_2 \times W_{2-1}) + (I_3 \times W_{3-1}) + (I_4 \times W_{4-1}) + (I_5 \times W_{5-1})$$
$$O_2 = (I_1 \times W_{1-2}) + (I_2 \times W_{2-2}) + (I_3 \times W_{3-2}) + (I_4 \times W_{4-2}) + (I_5 \times W_{5-2})$$
$$O_3 = (I_1 \times W_{1-3}) + (I_2 \times W_{2-3}) + (I_3 \times W_{3-3}) + (I_4 \times W_{4-3}) + (I_5 \times W_{5-3})$$
$$O_4 = (I_1 \times W_{1-4}) + (I_2 \times W_{2-4}) + (I_3 \times W_{3-4}) + (I_4 \times W_{4-4}) + (I_5 \times W_{5-4})$$

# The Applicability of Max-Pooing Layer

❖ The proposed mechanism can be performed in the max-pooling layer.

❖ The *SF_REG* size will be designed as a multiple of the filter size.

❖ The comparator can be reused by ReLU.

# Research Pillar 3:

## Smart Manufacturing

# Mini-factory in Ceres Lab



https://www.youtubeeducation.com/watch?v=_9scuX6REvQ

# Overview of Problem Formulation



*Reduce the computation complexity*

*Determining the number of layers for computing*

*Change in working condition*

Sensor $S_1$

Sensor $S_2$

Sensor $S_3$

Sensor $S_4$

Sensor $S_5$

Sensor $S_6$

$S_1$
$S_2$
$S_3$
$S_4$
$S_5$
$S_6$

*NN with 1 Hidden Layer*

*NN with 2 Hidden Layer*

*Transfer weights*

*1. Sensor Selection*

*2. Feature Extraction*

*3. Model Selection*

*4. Model Adaptation*

# Sensor selection and neural architecture search(NAS) methods for RUL estimation

❖ **Features** (IEEE JETCAS 2023; Highlighted IEEE JETCAS paper)

  – Adopting LASSO method to select the valuable sensors along with model training and estimate RUL precisely

  – A lightweight NAS method is proposed to find a fit neural network model during the model training phase

# Feature Dynamic Adaptive Thresholding Normalization (D-FATN) to Mitigate the Negative Transfer Learning

❖ **<u>Features</u>** (accepted by IEEE TIM 2024)

   – Adopting <span style="color:red">Dynamic Feature Adaptive Thresholded Normalization (D-FATN)</span> to enhances important features while suppressing redundant ones by utilizing mini-batch statistics for normalization

# Proposed Dynamic Feature Adaptive Thresholding Normalization (D-FATN) to Mitigate NTL

❖ Comparison of input feature distribution across various regularization techniques for different models



Input feature distribution transfer scenario during Fine-Tuning, a) Standard BN, b) Stochastic Normalization, c) Proposed Feature Adaptive Thresholded Normalization (FATN).

# Overview of the Proposed Work using Transfer Learning Model



Overview of the proposed 1D CNN model for bearing fault diagnosis adopting fine-tuned-based TL.

# Dataset Description

## CWRU Dataset description

| CWRU Datasets | Health conditions | Number of samples | Operation conditions |
|---|---|---|---|
| A | N/IRF/ORF/RF | 10 x 400 | 0 HP (1797 rpm) |
| B | N/IRF/ORF/RF | 10 x 400 | 1 HP (1772 rpm) |
| C | N/IRF/ORF/RF | 10 x 400 | 2 HP (1750 rpm) |
| D | N/IRF/ORF/RF | 10 x 400 | 3 HP (1730 rpm) |



## Paderborn Dataset description

| Paderborn Dataset | Health conditions | Number of samples | Operation conditions |
|---|---|---|---|
| PD | N/IRF/ORF | 1200/2200/2400 | 1500 rpm |



## NYCU Dataset description

| Paderborn Dataset | Health conditions | Number of samples | Operation conditions |
|---|---|---|---|
| NYCU | N/BF/CF/IRF/ORF | 4,096 | 1000 rpm |

# Experimental Results of the Proposed D-FATN in mitigating NTL

| Source → Target | L1 [26] | L2 [27] | L2-SP [22] | DELTA [23] | BN [24] | SN [25] | D-FATN |
|---|---|---|---|---|---|---|---|
| A → B | 85.1 | **96.4** | 94.4 | 95.4 | 94.0 | 95.0 | 95.2 |
| A → C | 83.8 | 94.4 | 88.7 | 93.4 | **97.2** | 94.6 | 95.2 |
| A → D | 80.8 | 95.9 | 91.0 | 93.3 | 87.5 | 94.5 | **96.8** |
| B → A | 87.7 | 96.6 | 96.0 | 95.8 | 97.2 | 97.3 | **98.7** |
| B → C | 89.0 | 96.5 | 96.1 | 94.0 | 96.4 | 95.6 | **98.2** |
| B → D | 87.5 | 96.4 | 94.7 | 96.2 | 97.3 | 96.8 | **97.6** |
| C → A | 88.5 | 92.6 | 94.2 | 95.0 | 93.4 | 95.6 | **96.5** |
| C → B | 89.0 | 96.0 | 95.0 | 94.7 | 94.0 | 95.2 | **96.9** |
| C → D | 89.8 | 95.8 | 94.0 | 93.2 | 95.5 | 96.4 | **98.5** |
| D → A | 89.7 | 95.2 | 94.5 | 94.6 | 96.0 | 92.7 | **97.5** |
| D → B | 89.1 | 96.0 | 94.8 | 94.5 | 94.1 | 95.7 | **97.3** |
| D → C | 87.6 | 96.7 | 95.1 | 94.2 | 93.9 | **97.2** | 96.6 |
| Average | 87.3 | 95.7 | 94.0 | 94.5 | 94.7 | 95.5 | **97.1** |

Same environment condition

| Source → Target | L1 [26] | L2 [27] | L2-SP [22] | DELTA [23] | BN [24] | SN [25] | D-FATN |
|---|---|---|---|---|---|---|---|
| A → PU | 86.7 | 96.1 | 91.4 | 94.3 | 95.2 | 93.8 | **97.4** |
| B → PU | 89.5 | 95.3 | 92.4 | 95.2 | 95.0 | 95.7 | **96.3** |
| C → PU | 88.9 | 96.4 | 95.1 | 94.4 | 95.5 | 97.2 | **98.4** |
| D → PU | 89.4 | 95.0 | 93.8 | 93.7 | 95.3 | 95.9 | **97.8** |
| PU → A | 88.9 | 94.7 | 96.6 | 94.5 | 89.0 | 96.2 | **97.8** |
| PU → B | 83.5 | 95.6 | 95.4 | 96.3 | 91.1 | 96.9 | **97.3** |
| PU → C | 80.7 | 90.9 | 88.6 | 94.6 | 91.0 | 92.0 | **96.8** |
| PU → D | 81.7 | 93.9 | 94.8 | 94.4 | 82.1 | 95.1 | **98.9** |
| Average | 86.1 | 94.7 | 93.5 | 94.6 | 91.7 | 95.3 | **97.5** |

Different environment condition

| Source → Target | L1 [26] | L2 [27] | L2-SP [22] | DEL-TA [23] | BN [24] | SN [25] | D-FATN |
|---|---|---|---|---|---|---|---|
| A → NYCU | 80.3 | 92.3 | 96.4 | 94.5 | 92.9 | 94.7 | **97.9** |
| B → NYCU | 90.6 | 94.4 | 95.5 | 96.3 | 90.5 | 95.7 | **97.0** |
| C → NYCU | 87.3 | 95.3 | 94.8 | 95.8 | 96.5 | 97.4 | **98.2** |
| D → NYCU | 89.5 | 96.4 | 95.4 | 92.0 | 95.2 | **97.3** | 97.1 |
| PU → NYCU | 87.3 | 96.9 | 94.8 | 95.9 | 89.4 | 97.3 | **98.2** |
| NYCU → A | 86.1 | 94.3 | 95.0 | 95.4 | 93.3 | 94.4 | **98.0** |
| NYCU → B | 86.8 | 96.9 | 95.0 | 94.2 | 94.0 | **97.8** | 97.2 |
| NYCU → C | 87.3 | 95.0 | 94.5 | 94.1 | 93.8 | 96.4 | **98.2** |
| NYCU → D | 87.6 | 96.3 | 95.1 | 95.3 | 93.5 | 96.3 | **97.9** |
| NYCU → PU | 87.1 | 96.5 | 92.2 | 95.7 | 93.9 | 96.1 | **97.3** |
| Average | 86.9 | 95.4 | 94.8 | 94.9 | 93.3 | 96.2 | **97.8** |

Different environment condition

# Conclusion

❖ **Three research pillars in Ceres Lab.**

  ❖ Smart Thermal Management on MPSoC

    ➢ Thermal sensor placement and temperature distribution reconstruction

    ➢ Temperature prediction and management

    ➢ Novel neural computing methods (not covered today)

      ➢ SNN, stochastic computing, semi-quantum computing

  ❖ Reconfigurable Neural Network design

    ➢ Lego-based DNN accelerator design flow

    ➢ NoC-based reconfigurable DNN accelerator

    ➢ Fast protocol translation for NoC-based TLM computing (not covered today)

  ❖ Smart Manufacturing

    ➢ LASSO-based NAS design and sensor selection

    ➢ Negative Transfer Learning (NTL) problem mitigation

      ➢ Dynamic Feature Adaptive Thresholded Normalization in CONV layer

      ➢ Source Free Unsupervised Domain Adaption (not covered today)

*Thank you for listening!!!*