

Thinker-IM: An Energy-Efficient Mixed Signal RNN Engine with Computing-in-Memory Techniques and Predictive Execution

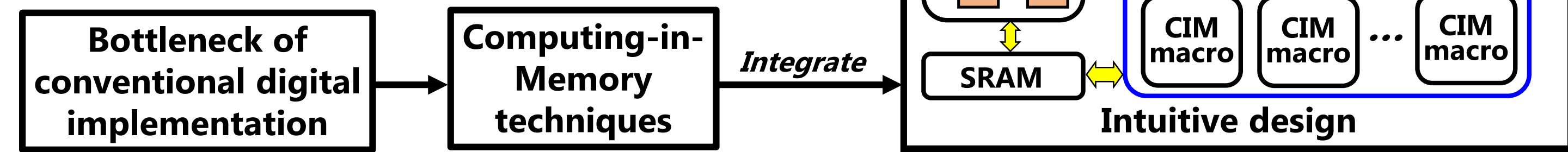
Ruiqi Guo¹, Yonggang Liu¹, Shixuan Zheng¹, Ssu-Yen Wu², Peng Ouyang³, Win-San Khwa², Xi Chen¹, Jia-Jing Chen², Xiudong Li³, Leibo Liu¹, Meng-Fan Chang², Shaojun Wei¹, Shouyi Yin^{1*}

¹Tsinghua University, Beijing; ²National Tsing Hua University, Hsinchu; ³TsingMicro Tech, Beijing; *yinsy@tinghua.edu.cn

Motivation

To achieve the ultra-low energy:

- 1: Binary neural networks;
- 2: Computing-in-memory techniques.

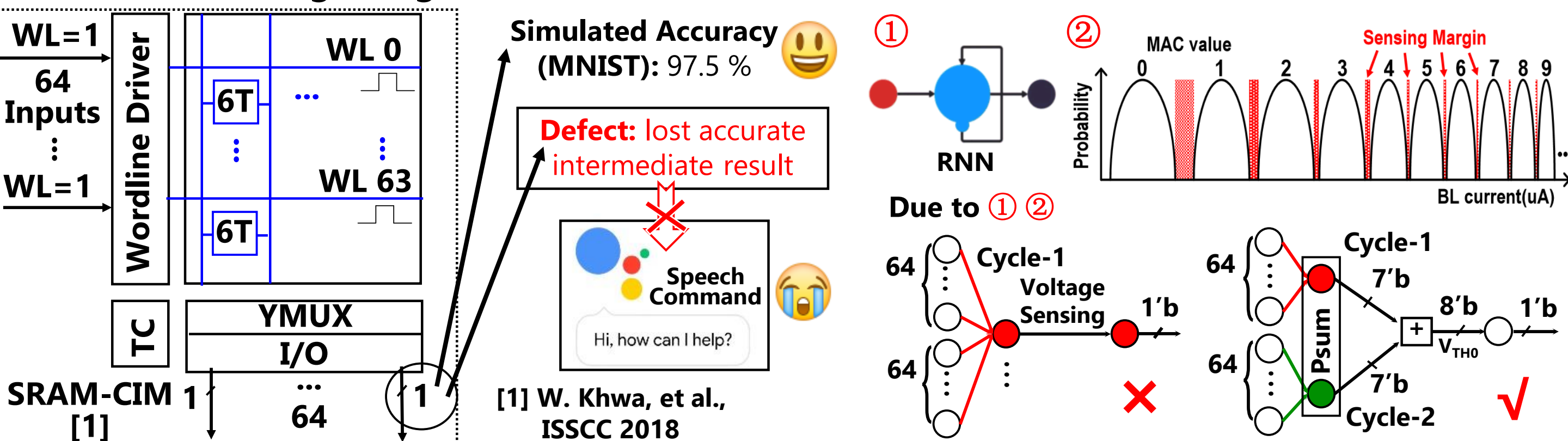


Challenge 1: Requirement for multi-CIM architecture and dataflow design

- A. Sufficient memory capacity is needed to fix all weights.
 - B. A big SRAM-CIM macro will mismatch the network shapes and computing resources. (Weights stored in the same WL need to be multiplied by the same input; Accumulated currents through the same BL contribute to the same output.) [1]
- SRAM-CIM unit-macro is needed to constitute multi-CIM architecture.
 - Dedicated computing flow is needed to schedule multiple SRAM-CIM macros.

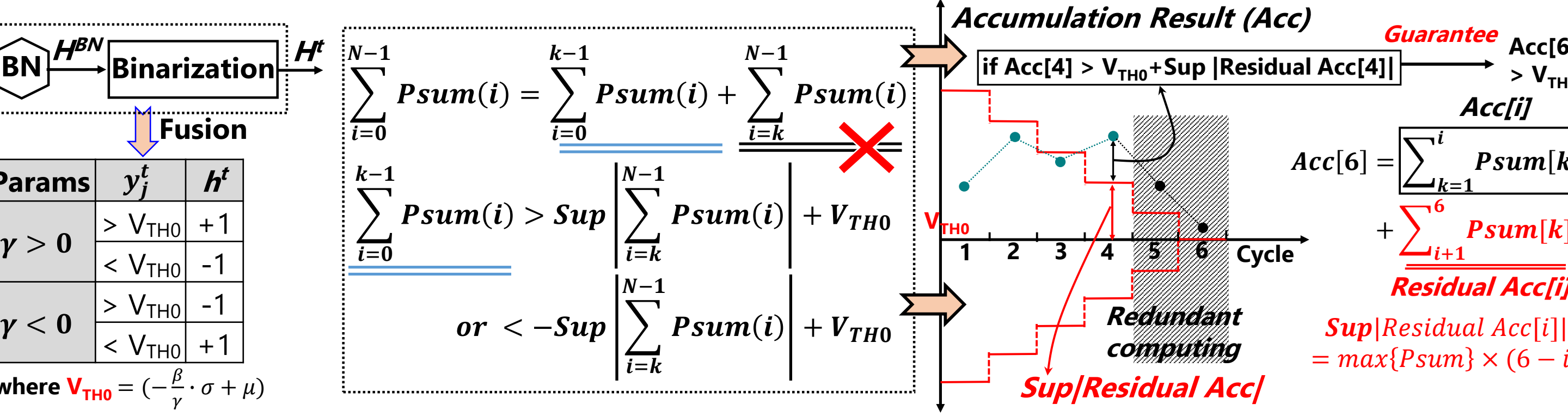
Challenge 2: Requirement for multi-bit output SRAM-CIM design

- A. RNNs are sensitive to bit-precision of intermediate results. (the current iteration is determined by the current input and previous hidden state.)
- B. RNNs should be partitioned to generate partial sums (Psums). (Small sensing margins across different MAC values of intermediate results.)

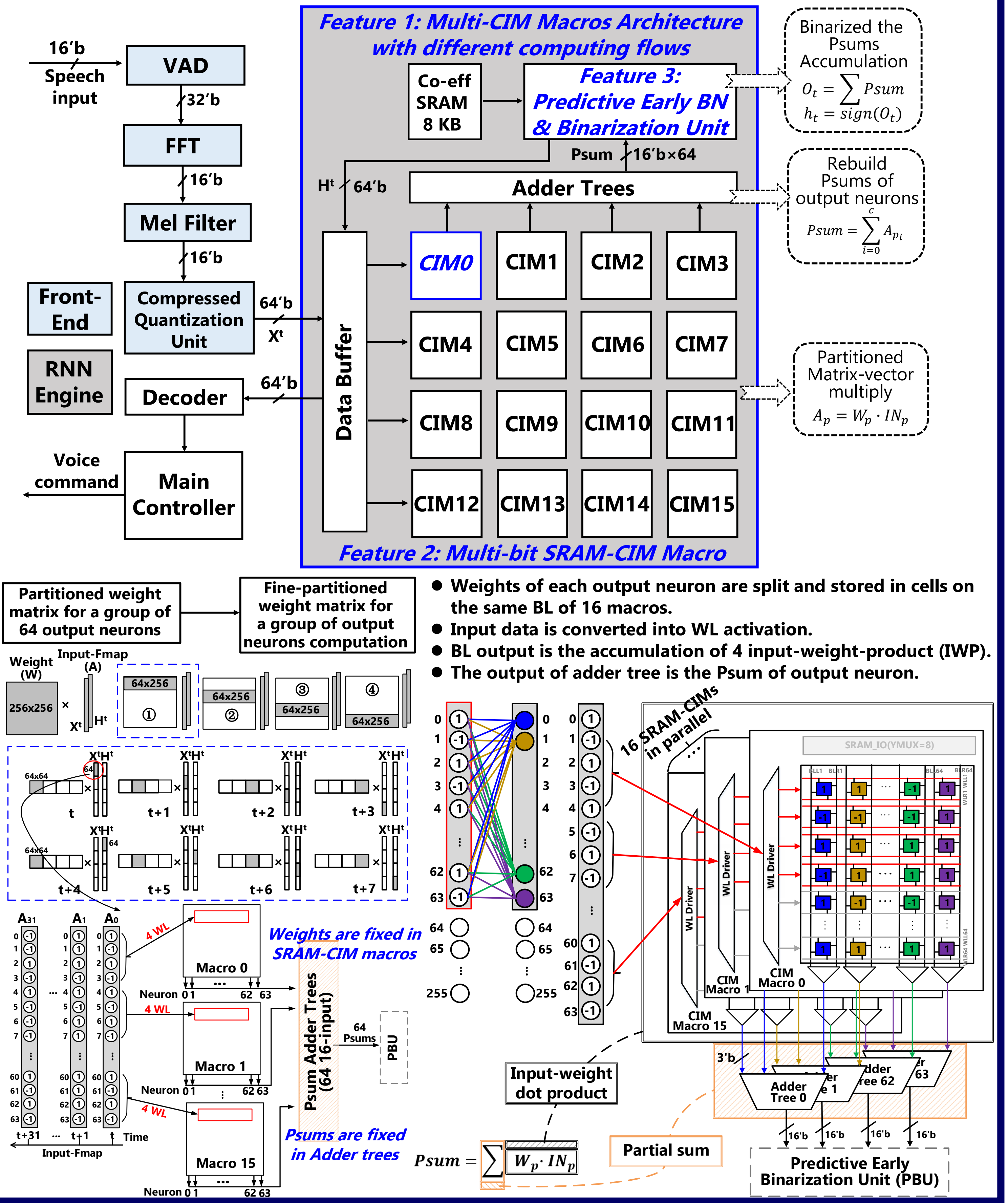


Challenge 3: binarized RNN has numerous redundant computations

- BN & binarization can be fused as a comparison, which can be early completed.



Architecture

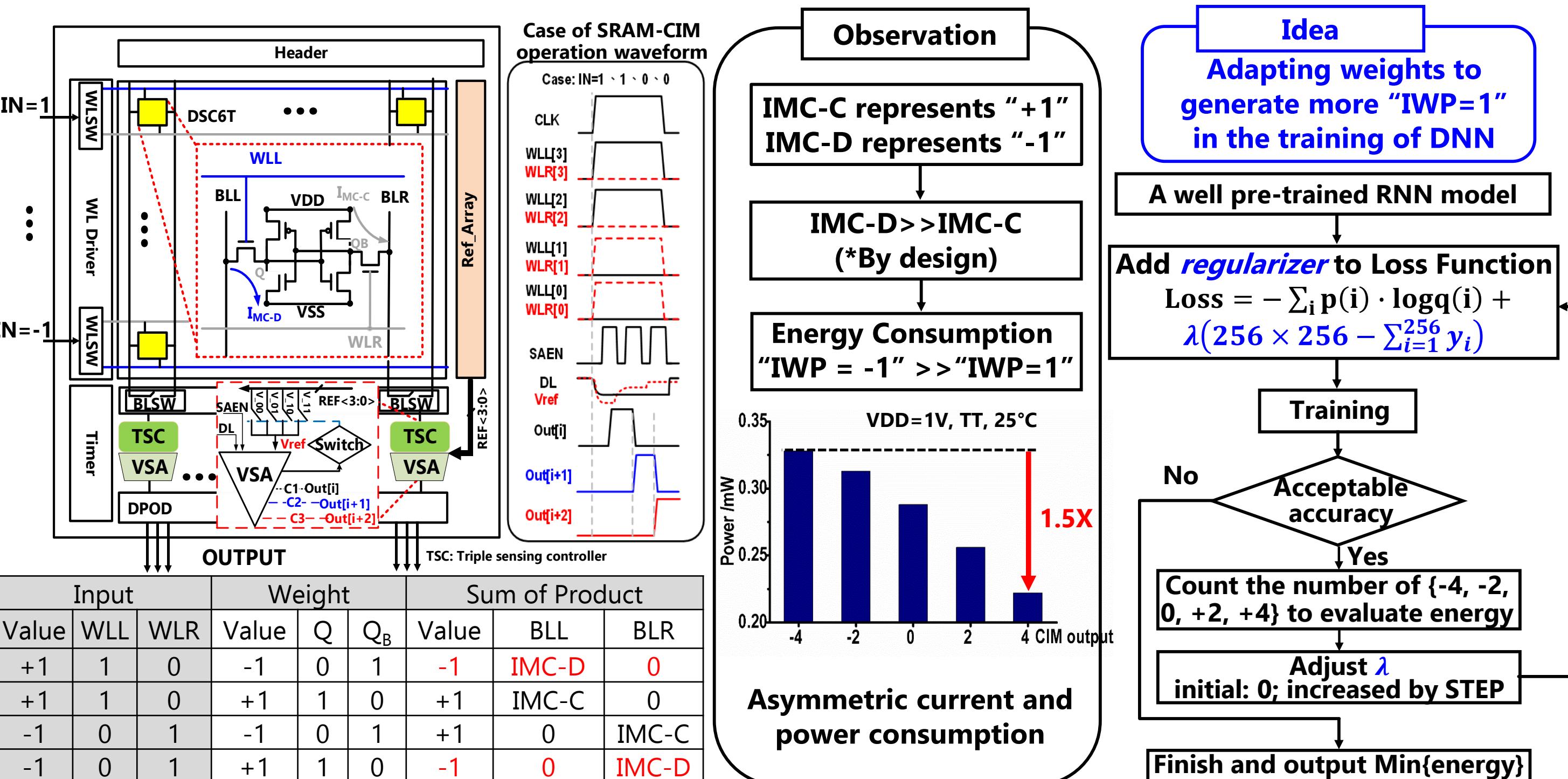


Low Power Techniques

1. SRAM-CIM Macro Design and CIM-aware Weights Adaptation

Key features of SRAM-CIM macro:

- A. Dual-split-control 6T memory cell to achieve XNOR; B. Serial-phase triple sensing controller to support 3-b output



2. Predictive early BN & Binarization Mechanism and Processing Unit

Exact Prediction

$$\sum_{i=0}^{k-1} Psum(i) > Sup \sum_{i=k}^{N-1} Psum(i) + V_{TH0}$$

$$\text{or } < -Sup \sum_{i=k}^{N-1} Psum(i) + V_{TH0}$$

Aggressive Prediction

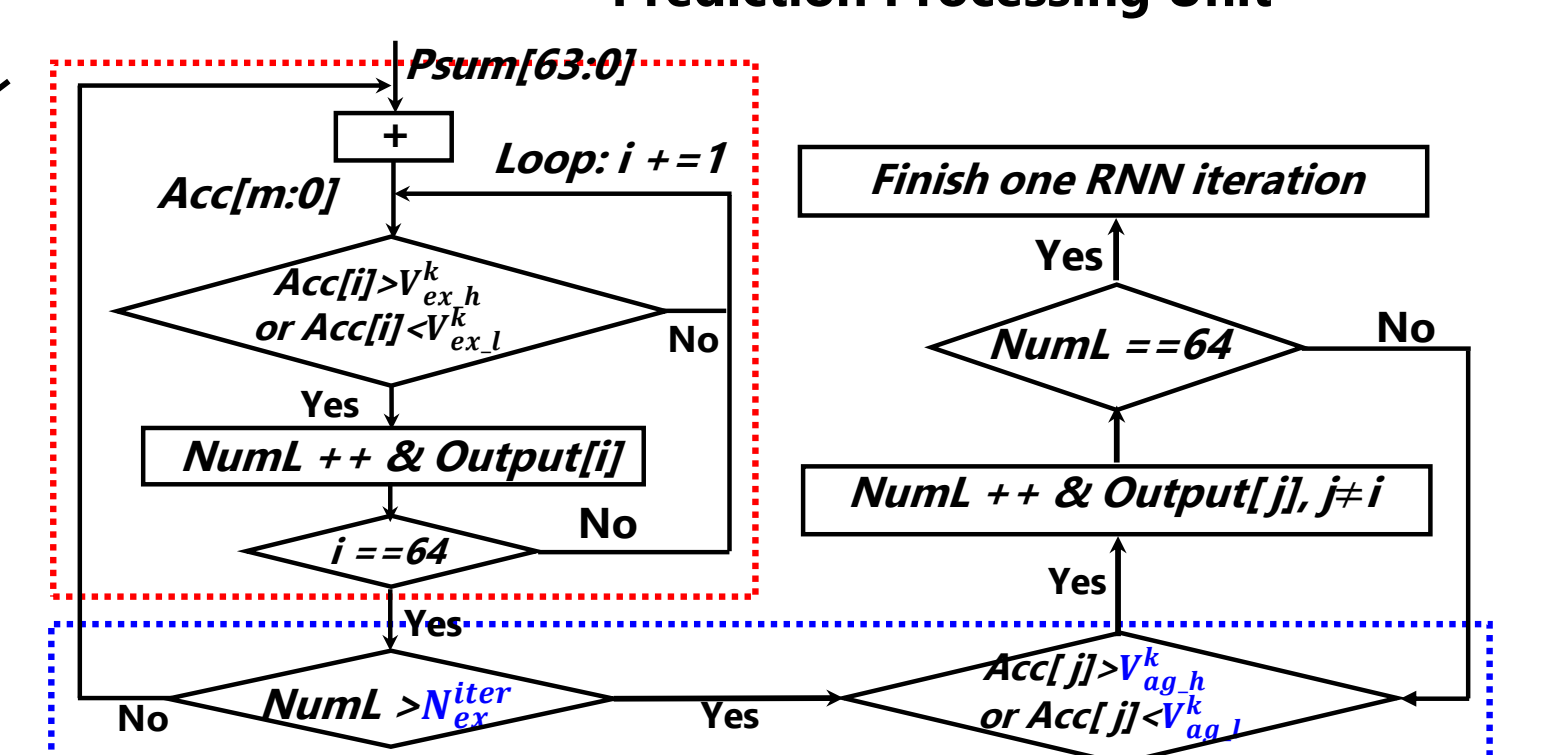
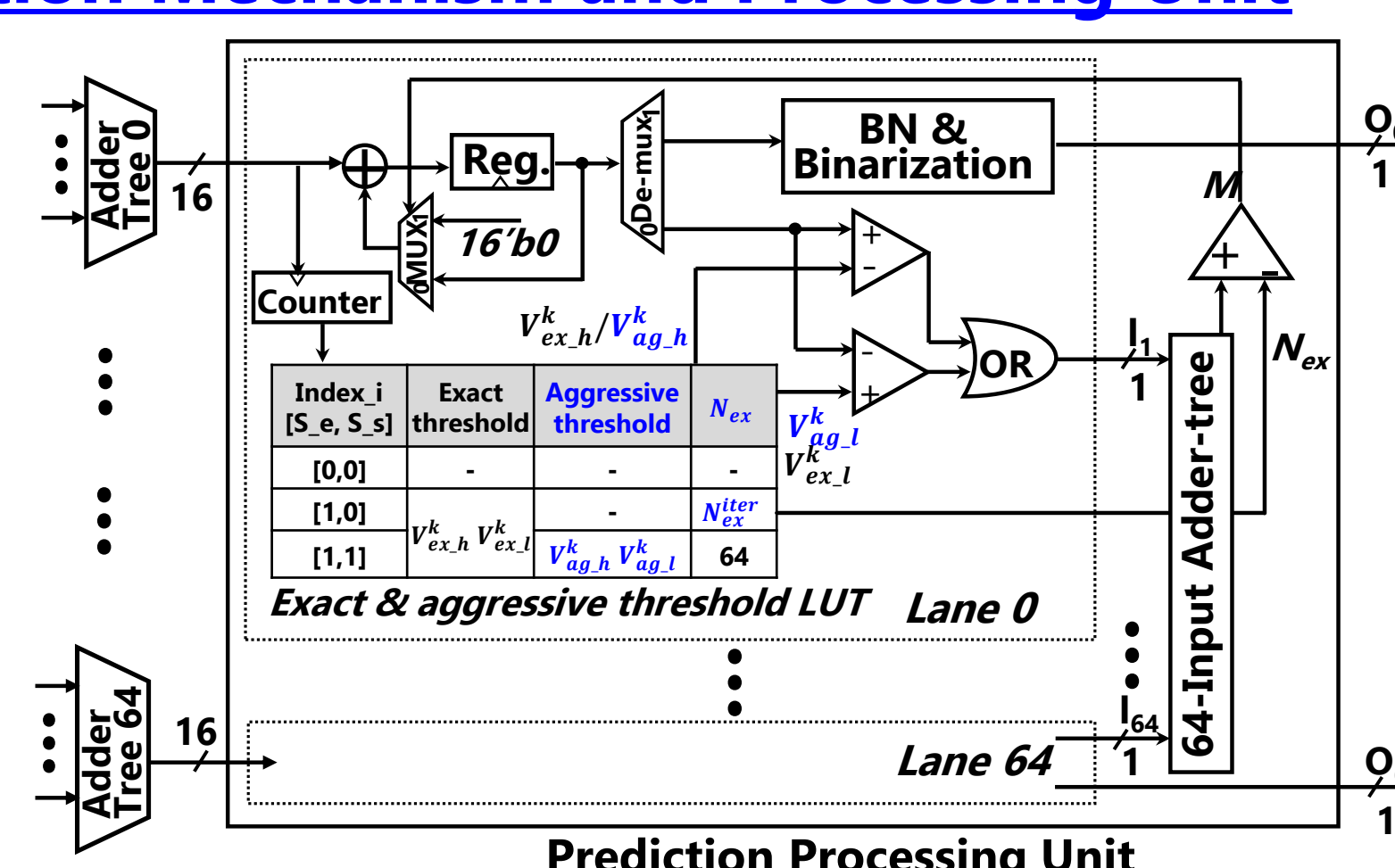
- A. Due to the fault-tolerate nature in RNN, a portion of Neurons (Nex) which are exactly predicted is enough to guarantee the final recognition accuracy.
- B. Vex can be relaxed to an aggressive threshold (Vag)

(Nex and Vag are determined by training offline)

Predictive early BN & Binarization Mechanism

Exact Prediction

Aggressive Prediction



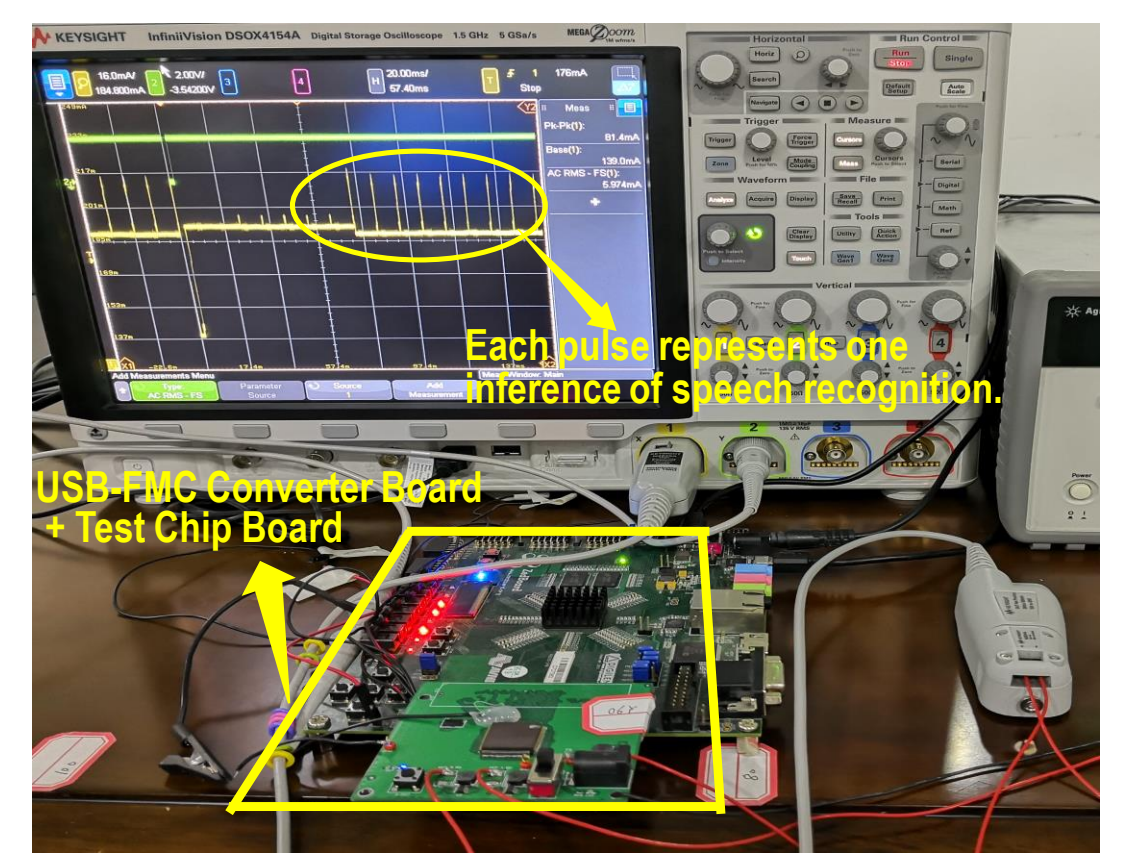
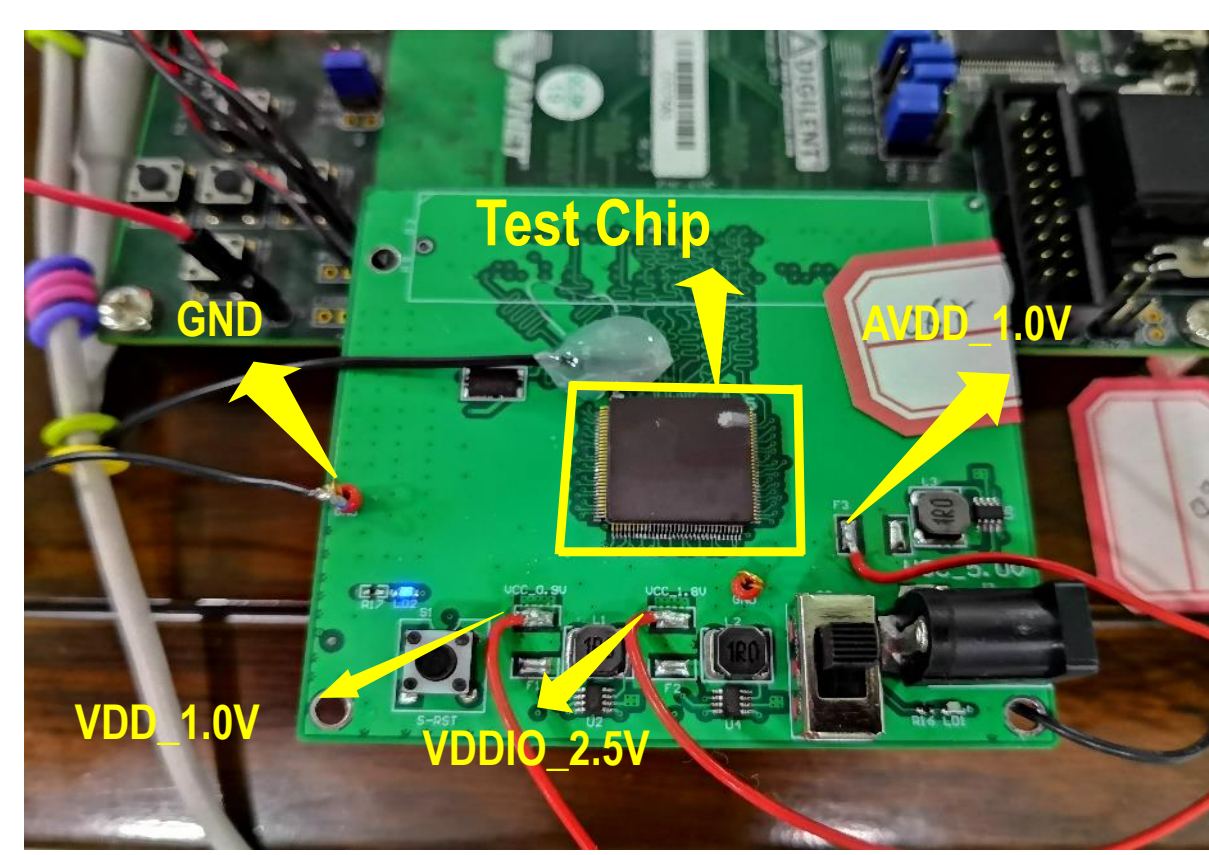
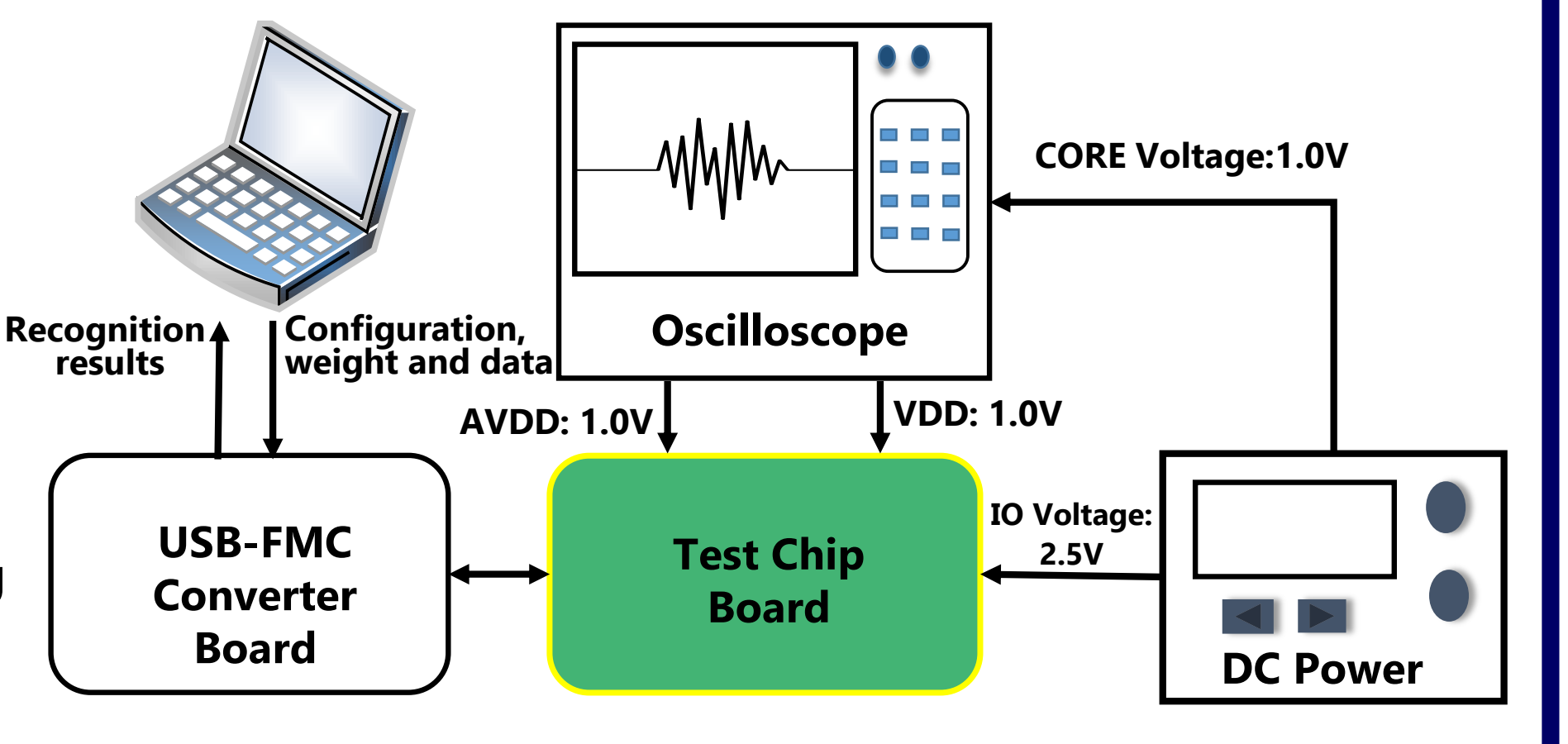
Verification

Demonstration System

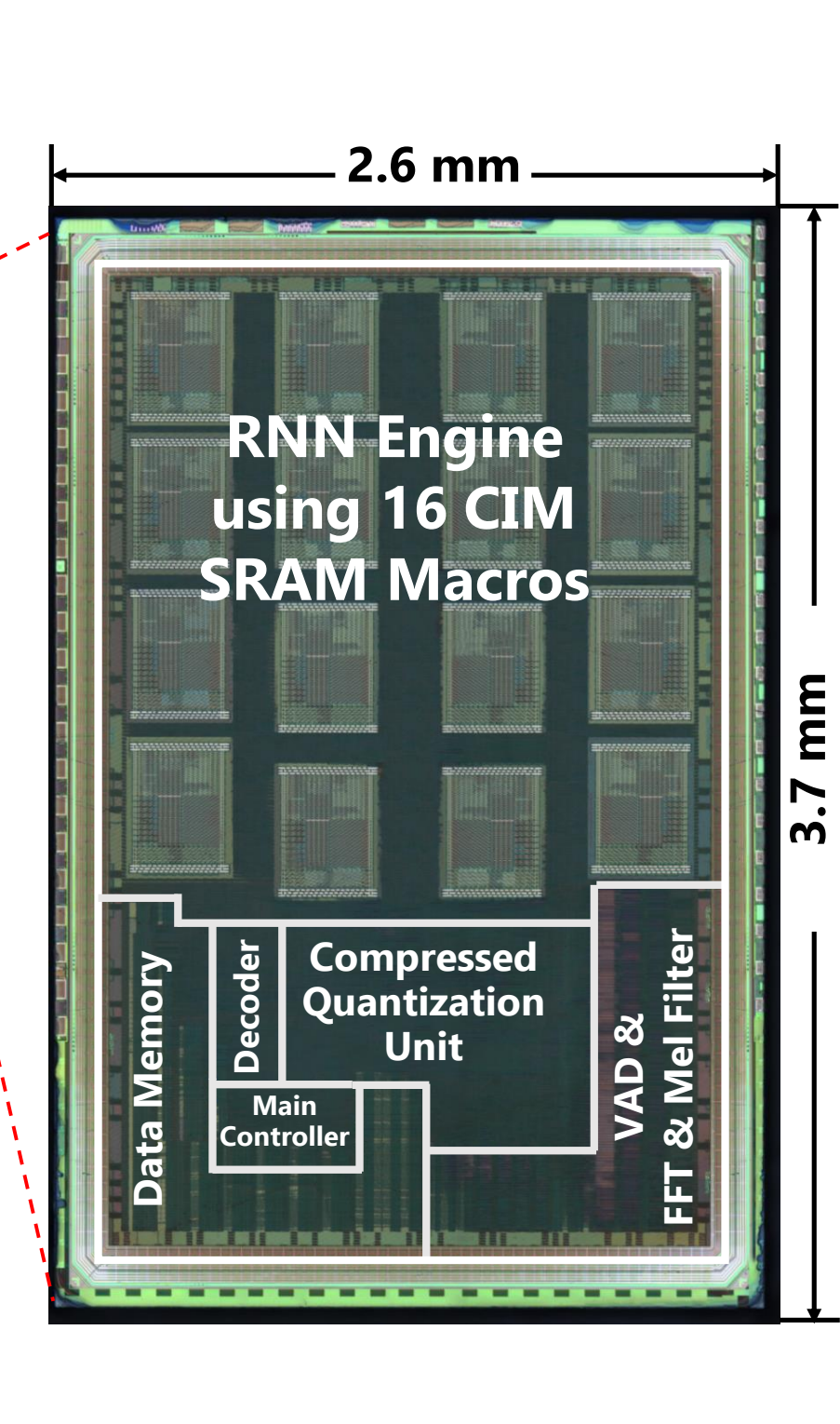
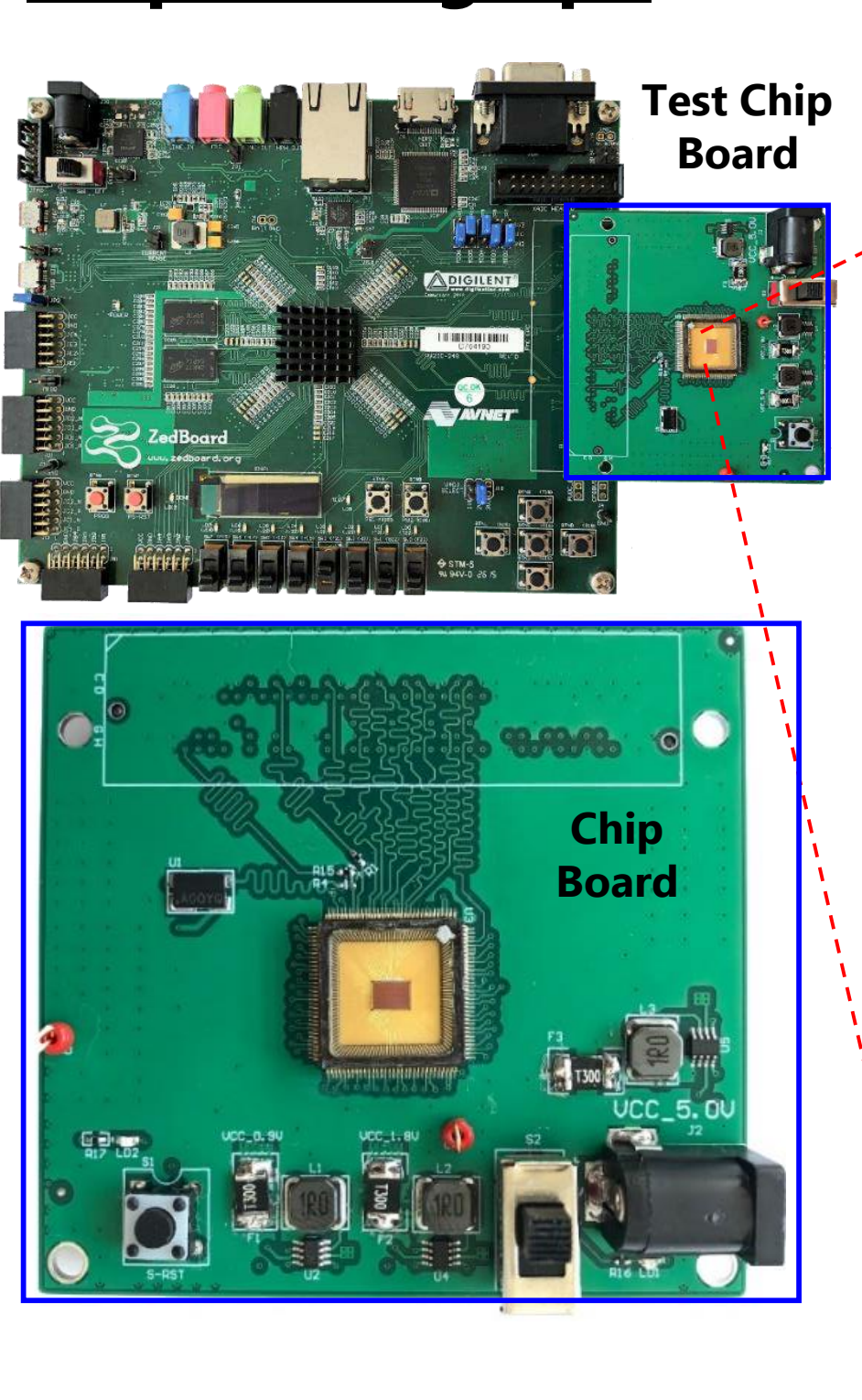
Test chip communicates with PC

- neural network weights, BN parameters, etc;
- testing data;
- configurations;
- recognition result.

Oscilloscope measures the working current



Chip Photograph



Chip Summary:

- Process: 65 nm CMOS
- Supply Voltage: 0.9 – 1.1 V
- Frequency: 5 – 75 MHz
- Core Size: 3.1 × 2 mm²
- Die Size: 3.7 × 2.6 mm²
- Neural Energy Efficiency: 5.1 pJ/Neuron @ 0.9 V, 75 MHz
- Arithmetic Energy Efficiency: 11.7 TOPS/W @ 0.9 V, 75 MHz

Key Features:

- A. Multiple SRAM-CIM architecture
- B. Multi-bit output SRAM-CIM
- C. Low-current training flow for SRAM-CIM architecture
- D. Predictive early BN and binarization method