### A Heuristic and Greedy Weight Remapping Scheme with Hardware Optimization for Irregular Sparse Neural Networks Implemented on CIM Accelerator in Edge AI Applications

Lizhou Wu<sup>1</sup>, Chenyang Zhao<sup>1</sup>, Jingbo Wang<sup>3</sup>, Xueru Yu<sup>2</sup>, Shoumian Chen<sup>2</sup>, Chen Li<sup>2</sup>, Jun Han<sup>1</sup>, Xiaoyong Xue<sup>1</sup>, Xiaoyang Zeng<sup>1</sup>

- <sup>1</sup> Fudan University, Shanghai, China
- <sup>2</sup> Shanghai Integrated Circuits R&D Center, China
- <sup>3</sup> Tongji University, China



- Background & Motivation
- Proposed Irregular-Sparsity-Aware CIM Accelerator
  - Heuristic & Greedy Weight Remapping
  - Hardware Optimization: <u>Zero-Skipping Circuit & Pipeline</u>
- **Experiment & Results**
- □ Conclusion

#### Background & Motivation

- Proposed Irregular-Sparsity-Aware CIM Accelerator
  - Heuristic & Greedy Weight Remapping
  - Hardware Optimization: <u>Zero-Skipping Circuit & Pipeline</u>
- **Experiment & Results**
- **Conclusion**

### **Background: Analog CIM accelerator for NN**

 $(\cdot)$ 

#### Crossbar CIM Pattern

- Map weights to the memory crossbar MAC are performed on the bitline
- Drive activations from the wordline Inputs are shared by all the columns

#### OU-based CIM Pattern

- Full crossbar CIM is too ideal:
  - ADC overhead / PVT deviatioin
- **Operating on sub-units (OU):** - Practical / Easier for sparsity

Restricted Column # 
Reduce ADC #

Restricted Row # Relieve ADC precision



### **Background: NN Sparsity & Pruning**

#### Neural Network Sparsity

- Activation (ReLU) & Weight (Pruning)
- Compact storage / Speedup / Efficiency

] Irregular Pruning

- Merits: Low accuracy drop / Simplicity
- Defects: Low utilization / High index cost

Our Target

### **Granularity Exploration**

Prune at fine-grained group









element-wise

channel-wise

Nice for Software Hard for Hardware

### **Background: Previous Work**

Weight mapping based on k-means clustering to skip all-0-crossbars



# Weight compression with input-distributing network



### **Motivation: Irregular Sparsity Utilization**

# Challenge — Contradiction between the irregular sparsity & regular CIM Motivation

Both concepts of OU & PG provide opportunites for sparsity utilization
 — Joint granularity exploration for the OU & PG

✓ Scattered zero-elements can be aggregated by weight remapping

— Apply heuristic method to efficiently search for the optima

70%-94% sparsity utilization (up to 1.6× improvement)

✓ Both the weight & actiavtion sparsity can be leveraged for acceleration

— Optimize the hardware to realize pipeline acceleration

3-7.6×speedup & 2.1-4.8×energy efficiency

#### Background & Motivation

- Proposed Irregular-Sparsity-Aware CIM Accelerator
  - Heuristic & Greedy Weight Remapping
  - Hardware Optimization: <u>Zero-Skipping Circuit & Pipeline</u>
- **Experiment & Results**
- **Conclusion**

#### □ Column shuffle remapping (CSR) for OU-based CIM



Prune the model and get the mask matrix





rows for compact storage

#### Definition of the parameters in the CSR problem

- ✓ CH<sub>in</sub> & CH<sub>out</sub>: the shape of the flattened 2D weight matrix
- $\checkmark$  **P**<sub>WL</sub> & **P**<sub>BL</sub>: the parallelism of the MVM conducted in one OU
- $\checkmark~N_{C}$  : The degree of the CSR freedom
- ✓ Pruning Factor (β):

 $\beta$  elements in one row are pruned and shuffled as a block

✓ Shuffle Window (SW):

 $\alpha$  OUs in one row are grouped as a shuffle window

$$N_{C} = \alpha \cdot P_{BL} / \beta$$

CSR is restricted to SW to relieve  $N_{C} \log N_{C} N_{C}/2$  the hardware overhead for shuffle



#### □ Heuristic CSR exploration based on Genetic Algorithm



- □ Storage compression based on Greedy Algorithm
  - $\checkmark$  The maximum deviation between OUs is less than  $P_{WL}$
  - ✓ Column shuffle can be merged with the heuristic process



#### Background & Motivation

#### Proposed Irregular-Sparsity-Aware CIM Accelerator

- Heuristic & Greedy Weight Remapping
- Hardware Optimization: <u>Zero-Skipping Circuit & Pipeline</u>
- **Experiment & Results**

Conclusion

- Zero-Skipping Circuit Loading (flattened act.):
  - Discard the activations corresponding to the compressed rows.
  - Record the zero flag for each activation bit.

#### Skipping (per act. bit):

- Skip all-0-bits according to the zero flag.
- Fire the effective bits to the word-line driver.

Joint the activations' bit-wise sparsity with weights' row-wise sparsity for acceleration



#### **The overall architecture and dataflow**

#### Cache (per OU row):

 Reuse the act. for the OUs in the same row

#### Load & Skip (per OU):

 Drive the effective bits of the act. to the OU

#### CIM (per OU):

 Traverse the OUs in the row-first order



#### ☐ The overall architecture and dataflow

#### ADC (per OU):

 Accumulate the bitwise multiplication results

#### **BENES (per SW):**

 Reorder the results to the normal sequence

#### S&A (per SW):

 Shift & add for weighted psum accumulation



#### □ The timing diagram of the zero-skipped pipeline

- The bit-serial pattern benefits directly from the proposed ZSC
- ✓ The Cache & Load
   stages are emitted by
   the handshake signal
- The skip of the all-0 OU can be hidden by
   the bit-serial process

		1			<u></u>	<u> </u>							<u></u>
ache	SWO			1					SW1		İ		1
Load		000					0U1	0U2	•			-	
Skip		AX					• •		AX				-
СІМ				A1 W0	A1 W1	A1 W2	A3 W0	A3 W1	A3 W2	A2 W0	İ		
ADC	Z	ZFV Value* W0					A1 W2	A3 W0	A3 W1	A3 W2	A2 W0		
enes	0	C1: 1010 C2: 0101				A1 W0	A1 W1	A1 W2	A3 W0	A3 W1	A3 W2	A2 W0	
S&A	A C3: 0001 C7: 0010 C6: 0000 C8: 0001						A1 W0	A1 W1	A1 W2	A3 W0	A3 W1	A3 W2	A2 W0

- Background & Motivation
- Proposed Irregular-Sparsity-Aware CIM Accelerator
  - Heuristic & Greedy Weight Remapping
  - Hardware Optimization: <u>Zero-Skipping Circuit & Pipeline</u>
- **Experiment & Results**
- **Conclusion**

### **Experiment: Case Study on Granularity**



### **Experiment: Setup**

### □ Hardware Configuration

- □ Memory: 128×128 RRAM crossbar
- □ OU: 16×8 sub-crossbar with 2-bits cell percision
- □ Peripheral: 5-bits ADC×16 / 16to16 Benes Network

#### Pruning & Quantization

- **D** Prune: Magnitude method based on L2 norm ( $\beta$ =4)
- Quant: LSQ method W@int8, A@uint8

### Behavior-Level Simulator

- □ Record the per-layer activations during inference (on ImageNet1K)
- □ Model the timing of the zero-skipped pipeline and the energy of the macro
- □ Analyze the performance according to the configuration and runtime trace

### **Experiment: Compression Rate**

#### □ 70%-94% sparsity utilization — up to 1.6× improv.



### **Experiment: Speedup & Energy Efficiency**

□ 3-7.6×speedup & 2.1-4.8×energy efficiency



- Background & Motivation
- Proposed Irregular-Sparsity-Aware CIM Accelerator
  - Heuristic & Greedy Weight Remapping
  - Hardware Optimization: <u>Zero-Skipping Circuit & Pipeline</u>
- **Experiment & Results**
- □ Conclusion

### Conclusion

The Co-optimzied Irregular-Sparsity-Aware CIM Accelerator

Heuristic & Greedy Weight Remapping

✓ Exploit the irregular weight sparsity for the regular CIM

• Zero-Skipping Circuit & Pipeline

✓ Utilize the exploited sparsity for speedup & energy saving

Joint exploration of the OU & prunning granularity
 ✓ Strike a balance between the compression rate, index overhead and degree of parallelism

## Thank You !