# Heterogeneous Chiplet-based Architecture for In-Memory Acceleration of DNNs

**Gokul Krishnan**<sup>1\*</sup>, Sumit K. Mandal<sup>2</sup>, Chaitali Chakrabarti<sup>1</sup> ,Jae-sun Seo<sup>1</sup>, Umit Y. Ogras<sup>2</sup>, Yu Cao<sup>1</sup>

<sup>1</sup>Arizona State University, USA <sup>2</sup>University of Wisconsin-Madison, USA **\*Email: gkrish19@asu.edu, Yu.cao@asu.edu** 









# Agenda

- Introduction
- Motivation
- SIAM Benchmarking Tool
- Big-Little Heterogeneous Architecture
- Experiments and Results
- Conclusion

# **Everything Goes Up**

- From data volume to information processing algorithms
- Further stress on the hardware execution platform



[Micro Focus; A. Gholami, 2020]

2018

2019

12.0

6S & 6S

6 & 6+

MoCo ResNet50

2020

2021

### Hardware Cost Analysis



32 bit floatine point had

160<sup>it floating point Multiply</sup>

8-bit Multiph

32.bit stam Read 18401

32-bit DRAM Read

160<sup>it floating point Add</sup>

0.01

8-bit Add

16-bit Add

Operation	Area(um <sup>2</sup> )*
8-bit Add	36
16-bit Add	67
32-bit Add	137
16-bit Floating-point Add	1360
32-bit Floating-point Add	4184
8-bit Multiply	282
32-bit Multiply	3495
16-bit Floating-point Multiply	1640
16-bit Floating-point Multiply	7700
32-bit SRAM Read (8Kb)	-
32-bit DRAM Read	-

\*45nm Technology Node

160troains point within Note: All the above data is adopted from A. Gholami 2021

# In-Memory Computing

- In-memory computing (IMC) provides a realistic solution to mitigate the von-Neumann bottleneck
- Combines both memory access and the computation into a single unit through analog domain computation
- Crossbar-based architecture provides a good platform for MVM computations in DNNs





# Why not Monolithic IMC?

- IMC accelerators utilize a weight stationary architecture with all weights on-chip
- Monolithic IMC chips result in high fabrication cost with increasing area. Larger area -> more defects, lesser yield, and higher cost



### 2.5D/Chiplet IMC – An Alternative

- 2.5D or chiplet architecture combine many small chips using an interposer to form a large system
- Each chiplet with smaller size improves the design effort, yield, reduces defect ratio, and reduces fabrication cost
- To achieve similar or better performance as monolithic architectures, careful design of the chiplet architecture and dataflow are necessary



We propose **SIAM** !!

### What **SIAM** can do and How it Helps !

- Platform for architectural exploration for chiplet-based IMC architectures (RRAM and SRAM)
- Designers get a wide range of parameters to tune and adjust for different architectural choices (e.g. mapping, partition schemes, IMC cell etc.)

Simulator	Architecture	Circuit	Interconnect	NoP Interconnect	DRAM
GenieX	Monolithic	SPICE-based	No	No	No
RxNN	Monolithic	SPICE-based	No	No	No
NeuroSim	Monolithic	SPICE-based	P2P (H-Tree)	No	No
MNSIM	Monolithic	thic Behavior model NoC-mesh		No	No
SIAM	Monolithic & Chiplet SPICE and Behavioral Model		NoC-mesh, NoC-tree, and H-Tree	Supported (driver and interconnect)	Supported

SIAM has been open-sourced and is available at the LINK

# **SIAM Performance Benchmarking Tool**

- SIAM Block Diagram
- SIAM Architecture
- Benchmarking Engines within SIAM
- Dataflow

# **SIAM Block Diagram**

 In-memory computing (IMC) hardware performance benchmarking tool that combines device, circuits, architecture, network-on-chip (NoC), network-on-package (NoP), and DRAM evaluation



### Inputs to SIAM

User Input	Description	User Input	Description		
	DNN Algorithm	Device and Technology			
Network Structure Data Precision Sparsity	DNN network structure information Weights and activation precision DNN layer-wise sparsity	Tech Node Memory Cell Bits/Cell	Technology node for fabrication RRAM or SRAM Number of levels in RRAM		
Intra	a-Chiplet Architecture	Inter-Chiplet Architecture			
Crossbar Size Buffer Type ADC Resolution Read-out Method NoC Topology NoC Width	IMC crossbar array size SRAM or Register File Bit-precision of flash ADC Sequential or Parallel Mesh or Tree Number of channels in the NoC	Chip Mode Chiplet Structure Chiplet Size Total Chiplet Count Global Accumulator Size NoP Frequency	Monolithic or chiplet-based IMC architecture Homogeneous or custom chiplet structure Number of IMC tiles within each chiplet Fixed count or DNN specific custom count Size of global accumulator Frequency of the NoP driver and interconnect		

#### **SIAM Architecture**

- Array of IMC chiplets, accumulator, buffer, and DRAM connected by an NoP fabric
- Supports both RRAM and SRAM-based IMC crossbar architectures



#### Dataflow





### Why Heterogeneous Chiplet Architecture

- Inherent non-linear weights and activations distribution in DNNs
- Adverse impact on the IMC utilization resulting in higher area and energy
- Affects hardware cost of the NoP within the architecture





# **Big-Little Heterogeneous Architecture**

Bank of big and little chiplets connect by an interposed and bridge-based
NoP
Little Chiplet NoP (Interposer)



# Mapping Overview

- Mapping algorithm aims to maximize IMC utilization by utilizing the Big-Little IMC chiplet -> determine the config of the architecture
- Little bank with smaller IMC are used for the initial/smaller layers while big bank is used for the larger/deeper layers -> Map the DNN layers
- NoP designed to exploit the volume of data movement in each of the banks -> determine the NoP configurations
- Little bank servicing most of the initial layers has higher data volume movement while the big bank has a lower data movement

### **Experiments and Results**

# IMC Utilization with Big-Little

- Utilize the mapping algorithm to determine the best configuration for big and little chiplets
- 256-64 and 256-32 have similar utilization. But 256-64 provides more resources and better energy-efficiency due to reduced peripheral circuits
  - Little Chiplet: 25 in number, 25 tiles/Chiplet, and 64x64 IMC size
  - Big Chiplet: 11 in number, 36 tiles per chiplet, and 256x256 IMC size



### Performance Comparison

- Compare performance with a homogenous all little and all big architecture for VGG-19 on CIFAR-100
- Proposed big-little architecture achieves reduced area, lower energy, and reduced latency

Configuration	Area				Energy				Latency						
e e magaranten	IMC	NoP	NoC	Total	Normalized to	IMC	NoP	NoC	Total	Normalized to	IMC	NoP	NoC	Total	Normalized to
	(%)	(%)	(%)	(mm <sup>2</sup> )	big-little (×)	(%)	(%)	(%)	(mJ)	big-little (×)	(%)	(%)	(%)	(ms)	big-little ( $\times$ )
Little only	11.9	88.0	0.1	952.1	10.9	99.7	0.2	0.1	1.3	4.1	99.7	0.1	0.2	1.6	1.3
Big only	44.0	55.5	0.5	597.2	6.8	78.6	11.0	10.4	0.43	1.3	99.6	0.1	0.3	3.2	2.7
Big-Little (this work)	52.4	47.4	0.2	87.4	1.0	99.8	0.1	0.1	0.32	1.0	99.2	0.3	0.5	1.2	1.0

### **EDAP** Comparison

- We compare the energy-delay-area product of the overall big-little architecture with all little and all big configurations
- Proposed big-little architecture achieves up to 329x improvement in EDA, while consistently outperforming the all big and all little configurations



## Support for Unseen Workloads

- We design our architecture to support different workloads by utilizing a local DRAM for each chiplet
- For an unseen workload, weights are written into the IMC arrays multiple times to complete one inference operation
- Ratio of DRAM energy and compute energy for different chiplet configurations

# Chiplets	VGG-1	.6	VGG-1	19	** All weights of VGG-19 fit on-chip
	#partitions	Ratio	#partitions	Ratio	with this config
36	2	1.1	1	0.08**	
25	2	2.1	2	131	-
16	3	3.6	2	161	-

### Comparison with Other Platforms

- Compared to Nvidia V100 and T4 GPUs, the big-little architecture achieves up to 9.6x improvement in area and 99.6x improvement in energy efficiency
- Compared to state-of-the-art accelerator from Nvidia (SIMBA), the big-little architecture achieves 2.4x area improvement and 18.4x improvement in energy efficiency

Platform	Area (mm <sup>2</sup> )	Energy Efficiency (Images/s/W)
Nvidia V100 GPU*	815	8.3
Nvidia T4 GPU*	525	15.5
SIMBA [23]	215	45
Big-Little (this work)	85	827

# Key Take Away

- We motivate the need for chiplet architectures for scalable acceleration of DNNs
- We introduce a novel benchmarking simulator SIAM that can support a wide range of configurations for architectural exploration
- We propose a Big-Little IMC architecture that utilizes a heterogeneous compute and interconnect structure for DNN acceleration
- Experimental evaluation of the proposed big-little architecture shows up to 9.6x improvement in area and 99.6x improvement in energy efficiency over state-of-the art GPUs and accelerators (SIMBA)