

Tutorial to NeuroSim: A Versatile Benchmark Framework for Al Hardware

Prof. Shimeng Yu

School of Electrical and Computer Engineering,

Georgia Institute of Technology

Email: shimeng.yu@ece.gatech.edu

Web: https://shimeng.ece.gatech.edu/

Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

Hardware Accelerators for AI/ML

- GPU still dominates the training in cloud, FPGA is good for inference for fast prototyping
- TPU (or similar digital ASIC) is ramping up in cloud as well as edge





TPU



GPU

Conventional computing platforms ~ 0.1 TOPS/W

Digital CMOS ASICs ~ 1-10 TOPS/W

Floating-point

Fixed-point

Compute-in-memory (CIM)

Analog CMOS (or eNVMs) ~ 10-100 TOPS/W

Low-precision \rightarrow accuracy?

- To further improve energy efficiency (TOPS/W), analog CIM (possibly with eNVMs) is promising especially in the edge inference where the model is pre-trained.
- NeuroSim is suitable for early design exploration of device candidates for CIM

CIM Basics: Mixed-Signal Compute



Multi-bit memory is preferred to serve as synaptic weights

S. Yu, H. Jiang, S. Huang, X. Peng, A. Lu, "Compute-in-memory chips for deep learning: recent trends and prospects", **IEEE Circuits and Systems Magazine**, vol. 21, no. 3, pp. 31-56, 2021, *invited review*

Key Metrics for Emerging Memories

	Mainstream Charged based Memories				Emerging Non-volatile Memories					
			FLASH		DCM		STT-	SOT-	ΕοΡΛΜ	EAEET
	JIAIVI	DIAN	NOR	NAND	P CIVI		MRAM	MRAM	I CIVAIVI	ICILI
Cell area	>150F ²	6F ²	10F ²	<4F ² (3D)	4~50F ²	4~50F ²	6~50F ²	12~100F ²	6~50F ²	6~50F ²
Multi-bit	1	1	2	3-4	2-3	2-3	1	1	1	2-3
Voltage	<1V	<1V	>10V	>10V	<3V	<3V	<1V	<1V	<2V	<3V
Read time	~1ns	~10ns	~50ns	~10µs	<10ns	<10ns	<10ns	~1ns	<100ns	<50ns
Write time	~1ns	~10ns	10µs-1ms	100µs-1ms	~50ns	<100ns	<20ns	<3ns	<100ns	<100ns
Retention	N/A	~64ms	>10y	>10y	>10y	>10y	>1y	>1y	>10y	>1y
Endurance	>1E16	>1E16	~1E5	1E3~1E4	1E6~1E9	1E3~1E9	1E6~1E14	~1E12	1E9~1E12	1E6~1E9
Write Energy (J/bit)	~fJ	~10fJ	100рЈ	~10fJ	~10pJ	لq~	۲d~	لq~	~100fJ	~fJ

F: feature size of the lithography.

The energy estimation is on the cell-level (not on the array-level).

PCM /RRAM/FeFET can potentially achieve less than 4F² through 3D integration.

The numbers of this table are representative (not the best or the worst cases).

Question: which eNVM is good for CIM?

S. Yu, Semiconductor Memory Devices and Circuits, Publisher: CRC Press/Taylor & Francis, 2022. [Link]

NeuroSim: Open-Source Simulator for AI Hardware



- 500+ citations
- >100 users including industry researchers from SK Hynix, Samsung, TSMC, and Intel
- IEEE Transactions on Computer-Aided Design (TCAD) Donald O. Pederson Best Paper Award 2023
- European Design and Automation Association (EDAA) Outstanding Dissertation Award 2023

NeuroSim History

The first appearance:

P.-Y. Chen, **S. Yu**, "NeuroSim: A circuit-level benchmark simulator for neuro-inspired architectures," Workshop on Hardware and Algorithms for Learning On-a-chip **(HALO) 2015**, Austin, TX.

Interface w/ PyTorch

DNN+NeuroSim V1 (Inference)

X. Peng, IEDM'19

Interface w/ PyTorch

DNN+NeuroSim V2 (Training)

X. Peng, TCAD'20

MLP+NeuroSim

(2-layer MLP)

P.-Y. Chen, IEDM'17, P.-Y. Chen, TCAD'18 **Contributors:**

Pai-Yu Chen, Xiaochen Peng Shanshi Huang, Yandong Luo Anni Lu, Junmo Lee, James Read, etc...

Impact of NeuroSim Family

Resulted in publications 400+



Research center (7%)





DNN+NeuroSim Methodologies (Python & C++)

- End-to-end frameworks for: ^{1*}neural network inference engine [1], ^{2*}online neural network training system [2], ^{3*}3D-CIM design [3]
- ➢ Input parameters: devices → circuits → architecture hierarchy and data-flow → algorithms
- Output results: inference/training accuracy, TOPS, TOPS/W and TOPS/mm²



[1] X. Peng, et al., IEDM'2019; [2] X. Peng, et al., TCAD'2020; [3] X. Peng, et al., IEDM'2020.

DNN+NeuroSim Methodologies (Python & C++)

Key Features

- Hardware-aware quantization [4] for weight, activation, gradient, error, as well as partial sum quantization based on ADC precision
- Built-in optimized weight/activation mapping flow [5] for maximum memory utilization
- Support various large-scale network models/datasets (ResNet, VGG / CIFAR-10, 100 and ImageNet)
- Compact PPA estimation model exists for every digital/analog circuit component within a chip, e.g.
 ADC, level-shifter, D-flipflop, adder tree, etc.
- PPA estimation model is **calibrated with SPICE simulation** at module-level
- Massive options (technology node from 130nm to 1nm, memory device type, R_{ON}, ADC option...)
- Interconnect modules (e.g. H-tree) estimated with parasitic RC delay and power
- Off-chip memory access overhead included (e.g. LPDDR access energy per bit)
- Non-ideal device/circuit effects-constrained accuracy estimation (e.g. device variations/noises, ADC quantization loss/ADC offset, etc.).
- Extendibility to variants of CIM hardware and various algorithms

[4] S. Wu, et al. ICLR'2018 [5] X. Peng, et al., ISCAS'2019;

Area Estimation



Latency Estimation

Read Write

1: RC delay (digital block)

Tech-file provides temperature-dependent $I_{ON} \& I_{OFF} \rightarrow$ get R (based on transistor size)

Similarly, get C (loading cap can be wire cap)

1*: Cadence fitting function (analog ADC)

Column resistance R_{COL} is calculated based on real-trace (mapped conductance and input vector)

R_{COL} as input of latency fitting function

2: Consider operation scheme

Each block needs to consider number of operations (i.e. T=N*t) Total latency needs to be "∑" or "MAX" according to the operation scheme (depends on working in sequence or in parallel)

Dynamic Energy Estimation



2: CV² dynamic energy (digital block)

Assume critical operation scheme, count number of transistors that need to be charged-up; Sum-up CV² for all chargedup transistors;

Also times #op (number of operation).

3: Cadence fitting function (analog block)

Column resistance R_{COL} is calculated based on real-trace (mapped conductance and input vector) R_{COL} as input of power fitting function

Energy is calculated based on latency and power fitting function

Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

DNN+NeuroSim V1 for Inference Engine



X. Peng, et al., IEDM'2019.

DNN+NeuroSim V1 for Inference Engine



ADC Design Considerations



For partial sum precision larger than 4bit, SAR ADC provides better tradeoffs than Flash ADC

Vin - Contract of the same of

-SAEN

Vout

VDD

SAEN -



17

Georgia School of Electrical and Tech Computer Engineering

Weight Mapping and Duplication for Pipelining



Georgia School of Electrical and Tech Computer Engineering

NeuroSim Validation with Real Chip Implementation

CIM RRAM at TSMC 40nm (W. Li, et al. CICC'2021)



0.5mm

mo	dulo	area	(um²)	latonov	energy (pJ)		
	uule	NeuroSim	Real chip	latency	NeuroSim	Real chip	
	DFF	719	681		1.07 × γ = 0.67	0.68	
shift-add	adder	662	663	1 cycle = 10ns	1.36 × δ = 0.11	0.10	
	inverter	2,133 INV = 1,336	1,334		1.99 × ∈ = 0.18	0.19	
accum	DFF	4,968	4,706		7.42 × γ = 4.67	4.67	
	adder	3,089	3,291	1 cycle = 10ns	6.28 × δ = 0.50	0.50	
	inverter	10,869 INV = 6,808	6,797		10.15 × ∈ = 0.91	0.93	
control	DFF	7,743	7,334		11.56 × γ = 7.28	7.28	
CONTROL	inverter	10,485 INV = 6,569	6,558		9.79 × € = 0.88	0.91	
total		31,893	31,386		15.20	15.26	

Adjustment factors introduced to calibrate:

- transistor sizing
- wiring area
- gate switching activity
- post-layout performance drop



Energy for whole array	NeuroSim	Real chip		
pre-layout	3,171.29pJ	10.4TOPS/W → 3,150.8pJ		
post-layout	3,171.29рЈ× <mark>ζ</mark>	8.48TOPS/W → 3,864.2pJ		

After calibration, the prediction with post-layout simulation error rate is less than 2%.

A. Lu, et al. AICAS 2021

DNN+NeuroSim V1 Example Results for Benchmarking

VGG-8 (8-bit activation; 8-bit weight) on CIFAR10, evaluated with DNN+NeuroSim V1.3									
Technology node (LP)	7nm		22nm						
Device	8T-SRAM	TPU-like (Google)	8T-SRAM	RRAM (Intel)	STT-MRAM (Intel)	Si:HfO2 FeFET (GF)	TPU-like (Google)		
Flash-ADC precision	4-bit	8-bit digital	4-bit	4-bit	4-bit	4-bit	8-bit digital		
Memory Cell Precision	1-bit	MAC	1-bit	1-bit	1-bit	1-bit	MAC		
Ron (Ω)	١	\	١	6k	1.41k	240k	١		
Cell area (F ²)	1080	\	280	60	100	40	١		
On/Off Ratio	١	\	۸	17	2.8	100	١		
Area (mm ²)	12.52	15.71	61.92	73.58	57.96	70.34	107.05		
L-by-L Leakage power (mW)	2.71	75.96	1.73	1.83	1.08	1.08	5.27		
Energy Efficiency (TOPS/W)	<mark>23.05</mark>	<mark>2.51</mark>	<mark>14.91</mark>	<mark>14.53</mark>	7.20	<mark>23.06</mark>	<mark>0.69</mark>		
Compute Efficiency (GOPS/mm ²)	<mark>47.26</mark>	<mark>75.00</mark>	<mark>5.54</mark>	<mark>5.48</mark>	<mark>0.62</mark>	<mark>10.43</mark>	<mark>3.60</mark>		

Memory R_{ON} is the critical factor to determine energy efficiency, desirable target: 100k ~ 1MΩ
 Compute-in-SRAM is a viable solution for high-performance applications, with outstanding energy efficiency and compute density, thanks to the logic scaling.

DNN+NeuroSim V1 Reveals Key Factors in CIM





- 7nm SRAM CIM TOPS/W is high, but suffers from leakage when the standby is frequent at edge.
- eNVM is the best option for edge device!
- 22nm FeFET CIM shows good TOPS/W, but not competitive TOPS/mm^2, how to improve?

21

Georgia School of Electrical and Tech Computer Engineering

Roadmap for Improving FeFET CIM



Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

DNN+NeuroSim V2 for Online Training

- Extends online training hardware and software analysis
- Non-ideal weight-update: asymmetry and non-linearity, device & cycle variation



X. Peng, et al., TCAD'2020.

DNN+NeuroSim V2 Dataflow and Training HW Support

Transposable synaptic array to support feed-forward & back-propagation (error calculation)
 Add weight gradient calculation (by SRAM-CIM), need frequent data reload (from DRAM)



For batch size=200, VGG-8 creates 2.6 GB, VGG-16 creates 30.4 GB, ResNet-18 creates 3.1 GB intermediate data \rightarrow Too large intermediate data for on-chip cache, thus DRAM access dominates the energy consumption!!

Asymmetry and Nonlinearity in Analog Synapse



- Asymmetry makes devices statistically easier to converge to the middle range of the conductance than approaching Gmax or Gmin.
- Frequent sign flipping causes a large unwanted conductance change towards the middle range.



Online Training: Nonlinearity & Asymmetry



Fig. 1 Analysis of nonlinearity and asymmetry (8-bit VGG-8 for CIFAR-10), w/ and w/o momentum optimization. Results show at +5/-5, the accuracy is still >80% (w/ momentum optimization).

Mommentum optimization:
$$\Delta W(t) = \beta \Delta W(t-1) + (1-\beta) \cdot (-\frac{\partial L}{\partial W})$$

Device & Cycle Variation Behavior Model



Online Training: Device & Cycle Variation



50

NL=+1/-1. std=0.05

50

Training *acy (%)

Accl

¹⁰⁰ # Epoch¹⁵⁰

¹⁰⁰ # Epoch¹⁵⁰

NL=+3/-3, std=0.05

200

200

NL=+5/-5, std=0.05

250

250

ining **∂**60

ā

Acc

50

50

¹⁰⁰ # Epoch¹⁵⁰

¹⁰⁰ # Epoch¹⁵⁰

Device-to-device variation

could be tolerated to certain degree, but accuracy will drop gradually with increasing variation.

L=+1/-1. std=0.05

NL=+3/-3, std=0.05

NL=+5/-5, std=0.05

200

200

200

90

80

60

50

40

250

250

250

Cycle-to-cycle variation could be tolerated to certain degree, but accuracy will drop quickly beyond the threshold.

Best Accuracy (%)

store.03 540005

storo.01

- NL=+3/-3

- NL=+5/-5

Online Training: System Performance





Fig. The data shows the <u>100th epoch</u> of FeFETbased CIM online training accelerator. (a) area breakdown by main components; (b) latency and (c) energy breakdown by main components; (d) latency and (e) energy breakdown by operations; (f) peak latency and (g) peak energy breakdown by operations.

- ADC (6-bit) & weight gradient units (SRAM-based) dominant in total area
- Buffer latency & DRAM energy is the bottleneck of performance
- Weight gradient computation is the bottleneck in the entire learning (require frequent DRAM access)



NeuroSim Extension



Reconfigurable NeuroSim [p1]



Monolithic 3D partition between memory at legacy node with **BEOL** oxide transistors and peripheral logic at advanced 7nm node

Monolithic 3D NeuroSim [p2]



Heterogeneous **3D NeuroSim** [p3]

77K and 4K cryogenic transistor technology files calibrated for data-center computing and quantum peripheral control

Cryogenic

0.2

10[.] 10-5

10 ٤

10

10

10⁻¹

10

_o 10⁻

NeuroSim [p4]

Reduce Vth until

w. 300 K NMOS

—— 300 K — 4 K, original

off becomes same

— 4 K. modified

0.4 0.6 0.8 1.0

3D NAND based architecture for GBlevel model for language, graph, genome, and recommendation system.





[p1] A. Lu et al., DATE 2020, Y. Luo et al. IEDM 2021; [p2] X. Peng et al., IEDM, 2020 [p3] X. Peng et al., TED, 2021; [p4] P. Wang el al., ISCAS, 2021; [p5] W. Shim et al., EDL 2021

Reconfigurable CIM Design

If chip area is constrained, off-chip weight reloading has two options: Option 1: reload weights (with sequential processing) Option 2: reload inputs (with batch processing)





Monolithic 3D

Specs

Architecture

Performance



Partition the design into M3D:

- Si FEOL for logic and ADC modules at 7nm
- Oxide BEOL for memory and its periphery at 45nm

X. Peng et al., IEDM, 2020



Integration	2D			3D			
[LSTP] Tech node (nm)	45 ^(a)	7	7 ^(b)	45Si+45 ^(c) 45Ox+45 ^(d) 45Ox+7		450x+7 ^(e)	7Si+7 ળ
Device	RRAM		SRAM	RRAM			SRAM
ADC precision	5-bit		4-bit	5-bit			4-bit
Cell Precision	2-bit		1-bit	2-bit			1-bit
Ron (Ω)	6 K		-	6 K			-
On/Off Ratio	150		-	150			-
numColMuxed	8						
Area (mm ²)	163.05	3.61	8.36	269.27	269.27	14.50	13.06
Energy Efficiency (TOPS/W)	8.58	77.98	30.30	14.46	15.59	81.34	36.48
Throughput (TOPS)	0.56	1.97	1.95	2.43	2.54	5.50	5.56
Power Density (W/mm ²)	4.00E-04	7.00E-03	7.72E-03	6.23E-04	6.05E-04	4.66E-03	1.17E-02
Via Density (#/mm²)	-	=	-	3.85E+03	3.85E+03	1.07E+05	1.56E+05

Heterogeneous 3D

- > TSV and hybrid bonding for 5-tiers
- Two Scheme: "Layer-by-Layer" vs. "Pipeline"
- SRAM: both logic and memory tier @ 7nm
- RRAM: 2-bit/cell, RON is 6kΩ (on/off=150), memory tier @ 22nm, logic tier @ 7nm

(b) Pipeline

Global Bu

ADC &

Shift-Add

S

Tiers Logici

Tier-2 Memory

Tier-1 Logici

TSV diameter's sweet spot: 1~3um



Shift-Add

X. Peng et al., TED, 2021

25

19

SwitchMatrix

& Memory

(a) Layer-by-Layer

Tiers Menory)

Tier-2 Memory)

Tier ILogici



Georgia School of Electrical and Tech Computer Engineering

3D NAND based GB model




Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

NeuroSim Extension Roadmap



Logic technology is progressing towards <u>angstrom era</u> Will extend the support from 7nm to beyond 1nm node Key features to be captured in NeuroSim technology libraries:

- Fin depopulation and standard cell height reduction
- DTCO and backside power delivery
- Transition from FinFET to GAA stacked nanosheet
- Further support 3D CFET (NMOS on top of PMOS)





Georgia School of Electrical and Tech Computer Engineering

Standard Cell Dimension Trend

		Technology Node (nm)								
	14	10	7	5	3	2	1			
# of fins per standard cell (NMOS + PMOS)	4+4	3+3	2+2	2+2	2+2	1+1	1+1			
Cell height (nm)	576	330	240	180	144	114	80			
Contacted poly pitch (nm)	78	64	57	51	48	45	40			
PN separation length (nm)	136	100	83	64	45	40	15			
Fin pitch (nm)	48	36	30	28	24	26	24			
Barrier thickness for M0-M2 (nm)	2.5	2.5	2.5	2.0	1.5	0.5	0.5			



https://irds.ieee.org/images/files/pdf/2021/2021IRDS_MM.pdf

- The standard cell assumptions in NeuroSim capture the actual area scaling trend in industry.

Gate pitch (CPP) scaling has significantly slowed down in recent nodes, being limited by lithography/device performance requirement (short channel effect).

- Cell height scaling is driven by fin depopulation, PN separation scaling in recent nodes.

- Moore's law is extended by Backside Power Rail/3D integration techniques beyond 1 nm node.

39

TCAD Simulation Assumptions

			Techn	ology Node	e (nm)					
	14	10	7	5	3	2	1	Set target I _{on} from IRDS 2016 - 2022		
Fin height [FH] (nm)	42	45	50	50	48	-	-	Design FinFET/stacked NS	GHOF	
Fin width [FW] (nm)	8	8	7	6	5	-	-	corresponding to the target technology node based on the device parameters	СРР	
NS thickness [NT] (nm)	-	-	-	-	-	6	6	Adjust S/D extension longth	FH FP	
NS width [NW] (nm)	-	-	-	-	-	15	10	(to tune S/D doping profile) to achieve	FW	
NS vertical pitch [VP] (nm)	-	-	-	-	-	14	12		Stacked NS	
# of stacked NSs	-	-	_	-	-	3.0	4.0	Calibrate work function and saturation velocity (1E7 ~	WAC	
Gate length (nm)	26	22	22	20	18	14	12	3E7 cm/s) until target I _{on} and I _{off} are satisfied	VP	
Equivalent oxide thickness [EOT] (nm)	0.9	0.8	0.7	0.65	0.6	0.55	0.5	Extract $C_{}/C_{+}/g_{}$ to set	Stacked NSs	
Gate height over fin [GHOF] (nm)	40.0	35.0	30.0	25.0	20.0	10.0 IRDS: 15	5.0 IRDS: 10	the gate capacitance, junction capacitance, and transconductance for		
Spacer width [SW] (nm)	12.0	8.0	8.0	7.0	6.0	6.0	5.0	NeuroSim simulation		

NeuroSim Input Device Parameters (from TCAD simulation)

	Technology Node (nm)								
	14	10	7	5	3	2	1		
On current/fin (µA)	54.744	58.725	60.139	61.320	64.788	66.385	59.005		
On current density (µA/um)	595.045	599.237	562.048	578.495	641.463	526.868	460.980		
Off current/fin (pA)	9.856	12.516	15.752	14.676	16.006	9.242	21.747		
g _m /fin (mS)	0.130	0.177	0.191	0.193	0.204	0.248	0.307		
Supply voltage (V)	0.800	0.750	0.700	0.700	0.700	0.650	0.600		
Gate capacitance (nF/m)	1.128	0.995	0.939	0.772	0.719	0.633	0.523		
Junction capacitance (F/m ²)	0.012	0.013	0.014	0.012	0.013	0.009	0.010		

- We allow some exceptions in IRDS standard: For 7 nm Current/fin, we decide the values based on our own projection to maintain current/fin improvement with technology scaling
- Effective width = fin height * 2 + fin width

BEOL Trend

- Exponential increase in resistivity at advanced nodes (sub 20 nm) due to surface/grain boundary scattering.
- Assume copper interconnect and introduce modified **FS-MS model** to estimate the copper resistivity as a function of aspect ratio, metal width, barrier thickness and grain size.
- Set 0.2fF/um for parasitic capacitance across all technology nodes.

 $\rho = \rho_0 \frac{3}{8} C (1 - p) \left(\frac{1}{h} + \frac{h}{A} \right) \lambda + \rho_0 \left[1 - \frac{3}{2} \alpha + 3\alpha^2 - 3\alpha^3 \ln \left(1 + \frac{1}{\alpha} \right) \right]^{-1}$ (1) (FS-MS model)



Fig. 12 Grain structure is revealed by TEM of (a) low, (b) medium, and (c) high AR wires as viewed normal to direction of current flow



values of R and p are applied to each data set.

Neurosim Metal Width Assumptions

		Technology Node (nm)								
		14	10	7	5	3	2	1		
Metal track in M0 (T)		9	7.5	6 6		6	5.7	5		
Metal pitch (nm)	M2	64	44	40	36	32	24	16		
	M0	64	44	40	30	24	20	16		
	M1	78	64	57	34	32	23	20		
Barrier thickness for M0-M2 (nm)		2.5	2.5	2.5	2.0	1.5	0.5	0.5		

Pyzyna, et al. Resistivity of copper interconnects beyond the 7 nm node

42

DCIM Scaling Trend



Georgia School of Electrical and Tech Computer Engineering

DCIM Scaling Trend



- Circuit-level analysis of DCIM down to 1 nm node is enabled by NeuroSim 1.4 framework. FOM=TOPS/WxTOPS/mm²
- Macro-level and chip-level DCIM vs ACIM benchmark at state-of-the-art technology node is performed using the updated NeuroSim 1.4 framework.
- Using the nominal chip design parameters for ResNet18-ImageNet workload, 2nm DCIM is found to outperform stateof-the-art 7nm SRAM-based ACIM.

Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

NeuroSim Extension Roadmap



Timeloop: simulator for systolic array-based DNN accelerator (TPU-like workloads)

NeuroSim: provide macro-level memory read/write latency/energy Integrate NeuroSim and Timeloop to get system-level performance with more diverse memory technologies as TPU global buffer



NeuroSim Extension Roadmap



GEM5: a common simulator for computer architecture research (CPU/GPU workloads)

NeuroSim: provide macro-level memory read/write latency/energy (similarly as CACTI and NVSim, but supporting more advanced and emerging technologies)

-> Build an interface between NeuroSim and GEM5 to run the IPC (instructions per cycle) with more diverse and accurate memory model

Emerging memories as TPU global buffer



SRAM suffers from low density and high leakage Emerging memory criteria to replace SRAM:

- high speed access (<10ns)
- high endurance (>10¹²)

Emerging memory candidates:

- 2T gain cell

RWL

WWL

- Ferroelectric memories (e.g., FeFET, FeRAM)
- Magnetic memories (e.g., STT-MRAM, SOT-MRAM)



TPU global buffer benchmarking settings

	Edge TPU	Cloud TPU				
Technology node	22nm	3nm				
Clock frequency	200MHz	700MHz				
8-bit MAC PE size	16 × 16	256 × 256				
Register file size	512-k	pit/PE				
Global buffer size	2MB	24MB				
DRAM interface	LPDDR4, 1.2pJ/bit					
	22nm SRAM	3nm SRAM				
Global buffor	22nm BEOL 2T gain cell	3D projection: 22nm BEOL 2T gain cell & 3nm FEOL peripherals				
memory device type	22nm BEOL FeFET	3D projection: 22nm BEOL FeFET & 3nm FEOL peripherals				
	22nm STT-MRAM	3nm STT-MRAM projection				
	22nm SOT-MRAM	3nm SOT-MRAM projection				

3D projection



TPU global buffer benchmark flow



TPU global buffer benchmarking results



- SRAM is still competitive for high energy efficiency

- High write energy of ferroelectric and magnetic memories should be improved

- Emerging memories with BEOL integration or scaling down with logic technology are promising

 - in terms of standby leakage, eNVMs are more attractive when the compute activity is low (<1%) for edge intelligence



Georgia School of Electrical and Tech Computer Engineering

NeuroSim Extension Roadmap



Georgia School of Electrical and Tech Computer Engineering

Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

Motivation: large language model (transformer)

- Large model size: ~ hundreds GB parameters
- Various types of matrix multiplication (MM) workloads
 - Dense MM (e. g. W_QX), sparse input (e.g. AV), dense input & sparse output (e.g. QK^T)





in MHSA module

Technology	Digital MAC (7nm)	SRAM CIM (7nm)	2-bit FeFET CIM (22nm)
Energy/8-bit MAC	~1pJ	~0.5pJ	~0.2pJ
SpMM	Yes	No	No
Write endurance	Unlimited	Unlimited	~10 ⁵
Write speed	1 cycle	1 cycle	> 100ns



- 3D Heterogeneous computing platform
- Approximate computing for attention probability

Approximate computing for attention prob.

- Use low precision (4-bit) CIM computing to determine the small attention scores
- Use 8-bit digital MAC for the attention scores with high magnitude



Pseudo code: approximate computing # 4-bit approx. computing with CIM $S_{low} = Q_{4b} \times K_{4b}^{T}$ (CIM) Atten_prob_{low} = Softmax(S_{low}) Mask = where(Atten_prob_{low} > threshold, 1, 0)

8-bit digital MAC for large scores $S_{high} = Q_{8b} \times K_{8b}^{T}$ for Mask > 0 $S = where(Mask > 0, S_{high}, S_{low})$

do softmax

S = S/sqrt(d_h) Atten_prob = Softmax(S) Atten_prob = max(Atten_prob, threshold)

linear combination with V
Out = SpMM(Atten_prob, V)

Fig. An illustration of the approximate computing based MHSA

3D Dataflow for MHSA computing: Q, K

- Mapping: Q, K in the same cube
- IFM delivery: broadcast to each heads
 - By group (input channel)
- Q, K: send to the bottom SRAM CIM for storage
 - Q, K of the same head will be in the same tile
 - Send to the other cube if it can not store all the heads (e.g. 4 heads in this case)



Fig. the dataflow for IFM delivery (left) and partial sum accumulation (right) for Q, K multiplication

3D Dataflow for MHSA computing: QK^T

- SRAM CIM for approximate computing
 - Send the mask of attention prob of the digital MAC
- 8-bit MAC for attention prob of high magnitude
 - Output stationary dataflow





Fig. the approximate computing scheme using CIM and digital systolic array

2.5D/3D Package-Level Benchmark Framework



Design Parameters

TABLE 1. The configurations of the proposed design and baselines												
Design name		Configuration of a cube (16 tiers)										
H3D_Hybrid (Optimal)	7nm systolic + 7nm SRAM CIM + 22nm FeFET CIM × 14											
H3D_MAC	7ni	7nm systolic + 7nm SRAM buffer + 22nm FeFET buffer× 14										
H3D_FeFET	2	2nm 1-bi	t FeFE	T CIN	VI × 2 + 2	22nr	n 2-bit	FeFET (CIM × 14			
H3D_SRAM	7nm SRAM CIM× 16											
TABLE	2. The device	and circu	uit par	ame	ters in t	he s	imulat	ion				
Device	Tech. node Cell bit		it C	Cell size		ADC-bit		allel ad	Write			
8T SRAM	F=7nm	1-bit		1080F ²		8-bit	8 r	ows	0.8V/1cycle			
FeFET	F=22nm	1 or 2-l	oit	54F ²		3-bit 8		ows	3V/50ns (1-bit)			
TABLE 3. The interconnect parameters for simulation												
Name	Pitch	Diamete		ter Heig		t C _{TSV}		Band width	Energy			
TSV	20µm	10µ	เท	100µr		~360fF		256- bit	0.35pJ/bit			
Name	Length	Pitch	Wid	Width H		(wire	Band width	Energy			
Si interposer	500µm	10µm	2μr	2µm		~	73fF	256- bit	0.17pJ/bit			
TABLE 4.	The projected p	performa	ince o	f a 7ı	nm syst	olic	array (32×32)				
Name	MAC energy	y T	hroug	hput	(dense)	1	Throughput (sparse)		Area			
Sparse engine	0.34pJ/MA0	2	1024 MACs/cycle				6 MAC	75 /cycle	0.89mm ²			

Interconnect Modeling and Validation



RC parasitics modeled for TSV, given the height/diameter and distance of keep-out-zone

2.5D links modeled as RC network

modeled

Reference

TSV-link1 TSV-link2

Georgia School of Electrical and Tech Computer Engineering

Impact of 2.5D/3D Technology Parameters



Fig. (a) The interposer area (form factor) and energy efficiency with different number of tiers stacked in a 3D cube. Energy efficiency reduces in more stacked tiers due to parasitics; (b) The interposer area and latency with different TSV bandwidth. The latency is reduced by 21.5% with 256-bit TSV bandwidth.

Results: approximate computing accuracy

- 4-bit weight and activation is needed for QK multiplication
- Negligible accuracy loss using Bert-base (~110M parameters) or GPT-2 (~1.7G parameters)



Fig. (a) The F1 score and sparsity of attention score with different threshold values to zero out small attention scores.(b) The accuracy w. and w/o. hardware effects of SRAM CIM array when calculating the attention scores for Bert-base and GPT-2

Microsoft Research Paraphrase Corpus (MRPC) Corpus of Linguistic Acceptability (CoLA)

Results: Hardware PPA Breakdown



A layer-by-layer computing scheme is adopted to achieve low power consumption for edge device. With such as low total power (~30mW), the temperature rise is small.

Energy breakdown for different workloads



Technology/Packaging Scaling Perspective



- Two scaling schemes are considered, one with packaging scaling only while the other with both packaging and transistor technology scaling.
- It is observed that both the energy efficiency and the latency could be bounded by computing if there is only packaging technology scaling.
- Therefore, it is important to synchronize the scaling of the packaging technology and transistor technology.

Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

Bioinformatics Challenges

• Motivation:

• Mass spectrometry (MS) and genome analysis are two foundational workloads for proteomics, drug discovery, and genetic diagnosis



- Challenges:
 - MS and genome data are booming, e.g. Over 521 TB MS data in UCSD MassIVE database (as of 05/27/2022)
 - MS and genome analysis are memory-intensive workloads that take long time (a few hours to days)



1. Karimi, M. R., Karimi, A. H., Abolmaali, S., Sadeghi, M., & Schmitz, U. (2022). Prospects and challenges of cancer systems medicine: From genes to disease networks. Briefings in Bioinformatics, 23(1), bbab343.

Motivation & Goals

• State of the art:

- Mass spectrometry: Ann-Solo¹ and HOMS-TC² for peptide identification and HyperSpec³ for clustering
- Genome analysis: BioHD⁴ for pattern matching



• Metrics:

- Algorithm: Accuracy and identification rate for identification/matching
- Hardware: Power, performance, area (PPA) for proposed designs
- 1. W. Bittremieux. et al. Extremely fast and accurate open modification spectral library searching of high-resolution mass spectra using feature hashing and graphics processing units. Journal of Proteome Research 18, 3792–3799 (2019)
- 2. Kang, Jaeyoung, et al. "Accelerating open modification spectral library searching on tensor core in high-dimensional space." Bioinformatics 39.7 (2023): btad404.
- 3. Xu, Weihong, et al. "HyperSpec: Ultrafast Mass Spectra Clustering in Hyperdimensional Space." Journal of Proteome Research (2023).
- 4. Zou, Zhuowen, et al. "Biohd: an efficient genome sequence search platform using hyperdimensional memorization." ISCA, 2022.

Motivation & Goals



• Proposed solution:

- Algorithm: Use hyperdimensional computing (HDC) to present/process bio workloads
- Hardware: Develop AIMS architecture based on 3D NAND Flash incorporated with 2.5D/3D heterogeneous integration to reduce the latency and improve the energy efficiency of current Mass Spectrometry (MS) analysis on CPU/GPU by 10-100x
- Success criteria:
 - Demonstrate accuracy of HDC-aided bio algorithms on CPU/GPU
 - Present 3D NAND-based accelerator or memory system
 - Evaluate the hardware feasibility using electrothermal simulation. Compare the PPA metrics with SOA GPU or in-memory solutions

Hyperdimensional Computing (HDC)

• Hyperdimensional Computing consists of two phases: training(store) phase and inference phase.



E.g., Genome Sequencing

Proposed 3D NAND HDC Engine

 The peripheral circuits are adapted from common 3D NAND periphery and adding additional ADC and comparator for hyperdimensional computing



3D NAND Associative Memory

• The query hypervector is sent in as BL voltage and times the conductance representing the corresponding class hypervector. The summed current is the dot product used for similarity check.



Heterogeneous Integration

• Two heterogeneous integration techniques for 3D NAND, Cu-Cu hybrid bonding and CMOS under array (CuA), are incorporated to achieve a compact form factor.


Simulation Parameters

- Dataset
 - Training dataset [1]:10,575 microorganism genome sequences with an average length of 3 million nitrogen bases per sequence.
 - Testing dataset: Artificially generate mutated sequences with the mutation rate of 10⁻³ per generation for classification.
- 3D NAND Parameters from industry [2].
- Digital circuits are synthesized using ASAP 7 [3].
- 3D NAND Peripherals are extracted from NeuroSim [4].

[1] "WoL: Reference Phylogeny for Microbes" https://biocore.github.io/wol/
[2] Hang-Ting Lue et al. 2019. Optimal Design Methods to Transform 3D NAND Flash into a High-Density, High-Bandwidth and Low-Power Nonvolatile Computing in Memory (nvCIM) Accelerator for Deep-Learning Neural Networks (DNN). In IEEE International Electron Devices Meeting (IEDM). 38.1.1–38.1.4.

[3] L. T. Clark et al., "ASAP7: A 7-nm FinFET predictive process design kit," Microelectronics Journal, vol. 53, pp. 105-115, July 2016.

[4] Xiaochen Peng et al. 2019. DNN+NeuroSim: An End-to-End Benchmarking Framework for Compute-in-Memory Accelerators with Versatile Device Technologies. In 2019 IEEE International Electron Devices Meeting (IEDM). 32.5.1–32.5.4.

	Parameters	Values	
Advanced	Technology	7nm FinFET Process	
CMOS Tier	VDD1	0.7 V	
	ADC Type	4-bit Flash CSA	
	Encoder Dimensionality	10000	
3D NAND	Device Dimensions	N/A	
Physical Parameter	SSL Pitch	220 nm	
	BL Pitch	100 nm	
	No. of WL/SSL/BL	32/15/10000	
	No. of Block	64	
	Array Size	$960~\text{SSL}\times10000~\text{BL}$	
	WL Staircase pitch	500 nm	
3D NAND Electrical	WL Read Voltage	1 V/4.5 V (V_select/V_pass)	
Parameter	SSL Read Voltage	4.5 V(activated)	
	BL Read Voltage	0.2 V	
	I_{on}/I_{off}	2 nA/1 pA	
CMOS	Technology	130nm Process	
Under Array	VDD2	1.3 V	

ADC Precision Requirement

- The simulation result indicates that the minimum requirement of ADC bit is 4-bit.
- The summed current distribution of each SL before sending to ADC with ON current of 2nA [1] is also simulated. The range of summed currents is in few μA which is reasonable for sensing.



[1] Hang-Ting Lue et al. 2019. Optimal Design Methods to Transform 3D NAND Flash into a High-Density, High-Bandwidth and Low-Power Nonvolatile Computing in Memory (nvCIM) Accelerator for Deep-Learning Neural Networks (DNN). In IEEE International Electron Devices Meeting (IEDM). 38.1.1–38.1.4.

3D NAND Device Non-ideal Effects

- The 3D NAND parameters [1] are ON current = 2nA, OFF current = 1pA, ADC bit precision = 10-bit, VBL = 0.2V and Vpass,r = 4.5V.
- From the silicon data, the current variation (σ/μ) is 0.1 and current shift (δ/μ) is ±0.05 which is tolerable from the simulation results.



[1] Hang-Ting Lue et al. 2019. Optimal Design Methods to Transform 3D NAND Flash into a High-Density, High-Bandwidth and Low-Power Nonvolatile Computing in Memory (nvCIM) Accelerator for Deep-Learning Neural Networks (DNN). In IEEE International Electron Devices Meeting (IEDM). 38.1.1–38.1.4.

Performance

- Peak performance and average performance with different lengths of the genome sequences are calculated.
- Throughput is estimated to produce more than 100 genome sequences classification per second and the average energy consumption for a single genome is
- 16.86 mJ • Comparing With commercial 3D NAND which requires >16 mJ to transfer the class hypervectors, the 3D NAND read energy of 11.7 μ J in the proposed system is >1000x smaller.

Peak Performance	280663	1.568 mJ	11.7 uJ	1.58 mJ
Avg. Performance	3107507	16.85 mJ	11.7 uJ	16.86 mJ
Latency				
Peak Performance	280663	0.7 ms	0.121 ms	0.82 ms
Avg. Performance	3107507	7.7 ms	0.121 ms	7.89 ms

Performance Breakdown

- Area breakdown shows a balanced design between two wafers.
- Latency and energy are dominated by the encoders. Future improvement is needed.





Outline

- Introduction to AI hardware and CIM paradigm
- NeuroSim V1 inference methodologies: area, latency, and energy estimation
- NeuroSim V2 training methodologies: supporting quantization and other non-ideal device effects.
- NeuroSim extension 1: technology updates to 1 nm node and DCIM support
- NeuroSim extension 2: TPU-like architecture benchmark with novel global buffer memory designs
- NeuroSim extension 3: chiplet based integration for ultra-large-scale transformer model
- NeuroSim extension 4: 3D NAND based CIM for hyperdimensional computing
- Demo of running inference engine benchmarking (DNN+NeuroSim V1.4)

- NeuroSim core: consists of multiple .cpp/.h files
 - .cpp files in the Neurosim core can be classified into 4 categories:
 - (1) Parameter files: ex) Param.cpp, Technology.cpp
 - (2) Circuit module files: ex) Adder.cpp, Buffer.cpp, Comparator.cpp ...
 - (3) Subarray/ProcessingUnit/Tile/Chip.cpp: Describes the circuit operation in each level of chip hierarchy
 - (4) main.cpp: description of the overall chip computation flow, final performance metrics calculation
 - Parameter files should be accessed to set the device/hardware parameters
 - Subarray.cpp, ProcessingUnit.cpp, Tile.cpp, and Chip.cpp include parameter files, circuit module files to estimate the power/area/latency of CIM hardware.

NeuroSim core: device parameters input in param.cpp

```
memcelltype = 2;
                                           // 1: cell.memCellType = Type::SRAM
                                           // 2: cell.memCellType = Type::RRAM
• Type:
                                           // 3: cell.memCellType = Type::FeFET
                                                   // precision of memory device
              cellBit = 2;
              /*** parameters for SRAM ***/
              // due the scaling, suggested SRAM cell size above 22nm: 160F^2
              // SRAM cell size at 14nm: 300F^2
              // SRAM cell size at 10nm: 400F^2
              // SRAM cell size at 7nm: 600F^2
              heightInFeatureSizeSRAM = 10;
                                                   // SRAM Cell height in feature size
  SRAM:
۲
              widthInFeatureSizeSRAM = 28:
                                                  // SRAM Cell width in feature size
              widthSRAMCellNMOS = 2;
              widthSRAMCellPMOS = 1;
              widthAccessCMOS = 1:
              minSenseVoltage = 0.1;
              /*** parameters for analog synaptic devices ***/
              heightInFeatureSize1T1R = 4; // 1T1R Cell height in feature size
              widthInFeatureSizelT1R = 12;
                                                // 1T1R Cell width in feature size
              heightInFeatureSizeCrossbar = 2;
                                                // Crossbar Cell height in feature size
              widthInFeatureSizeCrossbar = 2:
                                                 // Crossbar Cell width in feature size
              resistanceOn = 6e3;
                                               // Ron resistance at Vr in the reported measurement data (need to
              resistanceOff = 6e3*150;
                                                // Roff resistance at Vr in the reported measurement dat (need to
  eNVM:
              maxConductance = (double) 1/resistanceOn;
              minConductance = (double) 1/resistanceOff;
              readVoltage = 0.5;
                                                 // On-chip read voltage for memory cell
              readPulseWidth = 10e-9;
                                                 // read pulse width in sec
                                                 // Gate voltage for the transistor in 1T1R
              accessVoltage = 1.1;
              resistanceAccess = resistanceOn*IR DROP TOLERANCE;
                                                                          // resistance of access CMOS in 1T1R
                                                 // Enable level shifer if writeVoltage > 1.5V
              writeVoltage = 2;
```

• NeuroSim core: V1.4 update - SRAM scaling to 1nm node

// 1.4 update : SRAM size update

```
else if (technode==5){ // IRDS
if (technode>14){
                                                                                        widthSRAMCellNMOS = 1;
widthSRAMCellNMOS = 1;
                                                                                        widthSRAMCellPMOS = 1;
widthSRAMCellPMOS = 1;
                                                                                        widthAccessCMOS = 1;
widthAccessCMOS = 1;
                                                                                        heightInFeatureSizeSRAM = 19.2;
                                                                                                                                // SRAM Cell height in feature size
heightInFeatureSizeSRAM = 10;
                                    // SRAM Cell height in feature size
                                                                                                                                // SRAM Cell width in feature size
                                                                                        widthInFeatureSizeSRAM = 43.75;
widthInFeatureSizeSRAM = 28;
                                   // SRAM Cell width in feature size
}
                                                                                        else if (technode==3){ // IRDS
else if (technode==14){ // Samsung 14 nm
                                                                                        widthSRAMCellNMOS = 1;
widthSRAMCellNMOS = 1;
                                                                                        widthSRAMCellPMOS = 1;
widthSRAMCellPMOS = 1:
                                                                                        widthAccessCMOS = 1;
widthAccessCMOS = 1;
                                                                                        heightInFeatureSizeSRAM = 30;
                                                                                                                              // SRAM Cell height in feature size
heightInFeatureSizeSRAM = 10.6;
                                      // SRAM Cell height in feature size
                                                                                        widthInFeatureSizeSRAM = 68.26;
                                                                                                                                // SRAM Cell width in feature size
widthInFeatureSizeSRAM = 30.8;
                                     // SRAM Cell width in feature size
}
                                                                                        else if (technode==2){ // IRDS
else if (technode==10){ // TSMC 10 nm
                                                                                        widthSRAMCellNMOS = 1;
widthSRAMCellNMOS = 1;
                                                                                        widthSRAMCellPMOS = 1;
widthSRAMCellPMOS = 1;
                                                                                        widthAccessCMOS = 1;
widthAccessCMOS = 1;
                                      // SRAM Cell height in feature size
                                                                                        heightInFeatureSizeSRAM = 42;
                                                                                                                              // SRAM Cell height in feature size
heightInFeatureSizeSRAM = 12.8;
                                                                                        widthInFeatureSizeSRAM = 120;// 111.42;
                                                                                                                                        // SRAM Cell width in feature size
widthInFeatureSizeSRAM = 31.25;
                                      // SRAM Cell width in feature size
3
                                                                                        else if (technode==1){ // IRDS
else if (technode==7){ // TSMC IEDM 2016
widthSRAMCellNMOS = 1;
                                                                                        widthSRAMCellNMOS = 1;
widthSRAMCellPMOS = 1;
                                                                                        widthSRAMCellPMOS = 1;
widthAccessCMOS = 1;
                                                                                        widthAccessCMOS = 1;
heightInFeatureSizeSRAM = 16;
                                    // SRAM Cell height in feature size
                                                                                        heightInFeatureSizeSRAM = 80;
                                                                                                                              // SRAM Cell height in feature size
widthInFeatureSizeSRAM = 34.43;
                                      // SRAM Cell width in feature size
                                                                                                                              // SRAM Cell width in feature size
                                                                                        widthInFeatureSizeSRAM = 144;
```

- NeuroSim core: circuit parameters input in param.cpp
 - Technology:

// technode: 130 --> wireWidth: 175 // technode: 90 --> wireWidth: 110 // technode: 65 --> wireWidth: 105 // technode: 45 --> wireWidth: 80 // technode: 32 --> wireWidth: 56 --> wireWidth: 40 // technode: 22 // technode: 14 --> wireWidth: 25 // technode: 10, 7 \rightarrow wireWidth: 18 technode = 22;// Technology featuresize = 40e-9;// Wire width for subArray simulation wireWidth = 40;// wireWidth of the cell for Accuracy calculation

• Sub-array size: numRowSubArray = 128; numColSubArray = 128; // # of rows in single subArray
// # of columns in single subArray

ADC type/sharing/resolution:

```
      SARADC = false;
      // false: MLSA

      // true: sar ADC

      currentMode = true;
      // false: MLSA use VSA

      // true: MLSA use CSA

      numColMuxed = 8;
      // How many columns share 1 ADC (for eNVM and FeFET) or parallel SRAM

      levelOutput = 32;
      // # of levels of the multilevelSenseAmp output, should be in 2^N forms; e.g. 32 levels --> 5-bit ADC
```

- NeuroSim core: V1.4 update improved interconnect models & scaling
 - // 1.4 update: wirewidth } else if ((8 <= wireWidth) && (wireWidth < 10)){</pre> if (wireWidth >= 175) { AR = 3.00; Rho = 7.87*1e-8; AR = 1.6;} else { • Wire width: Rho = 2.01*1e-8; exit(-1); puts("Wire width out of range"); else if ((110 <= wireWidth) && (wireWidth < 175)) {</pre> AR = 1.6;Rho = 2.20*1e-8;Rho = Rho * 1 / (1- ((2*AR*wireWidth + wireWidth)*barrierthickness / (AR*pow(wireWidth,2))); // 1.4 update: Metal0 if (Metal0 >= 175) { } else if ((8 <= Metal0) && (Metal0< 10)){</pre> AR Metal0 = 3.00; Rho Metal0 = 7.87*1e-8; AR Metal0 = 1.6; } else { Metal0: • Rho Metal0 = 2.01*1e-8; exit(-1); puts("Wire width out of range"); else if ((110 <= Metal0) && (Metal0< 175)) { AR Metal0 = 1.6; Rho Metal0 = Rho Metal0 * 1 / (1- ((2*AR Metal0*Metal0 + Metal0)*barrierthickness / (AR Metal0*pow(Metal0,2))); Rho Metal0 = 2.20*1e-8; } else if ((8 <= Metal1) && (Metal1 < 10)){</pre> // 1.4 update: Metal1 AR Metal1 = 3.00; Rho Metal1 = 7.87*1e-8; if (Metal1 >= 175) { } else { AR Metal1 = 1.6; exit(-1); puts("Wire width out of range"); • Metal1: Rho Metal1 = 2.01*1e-8; else if ((110 <= Metal1) && (Metal1 < 175)) { Rho Metal1 = Rho Metal1 * 1 / (1- ((2*AR Metal1*Metal1 + Metal1)*barrierthickness / (AR Metal1*pow(Metal1,2))); AR Metal1 = 1.6; Metal0 unitwireresis = Rho Metal0 / (Metal0*1e-9 * Metal0*1e-9 * AR Metal0); Rho Metal1 = 2.20*1e-8; Metal1 unitwireresis = Rho Metal1 / (Metal1*1e-9 * Metal1*1e-9 * AR Metal1);

- NeuroSim core: technology parameters in Technology.cpp
- -> contains comprehensive device specifications for logic transistors down to 1 nm node.

```
if (transistorType == conventional) {
    if (featureSizeInNano == 130) {
        if (deviceRoadmap == HP) {
            /* PTM model: 130nm_HP.pm, from http://ptm.asu.edu/ */
            vdd = 1.3;
            vth = 128.4855e-3;
            phyGateLength = 7.5e-8;
            capIdealGate = 6.058401e-10;
            capFringe = 6.119807e-10;
            effectiveResistanceMultiplier = 1.54; /* from CACTI */
            current_gmPmos=3.94E+02;
            current_gmPmos=2.61E+02;
            currentOnNmos[0] = 0.93e3;
```

else if (featureSizeInNano == 1) {

// 1.4 update: IRDS 2022
vdd = 0.6;
vth = vth_list[6];
PitchFin= 24e-9;
phyGateLength = 1.2e-8;

// 1.4 update: IRDS 2022 - GAA specfic parameters
max_fin_per_GAA=1;
max_sheet_num=4;
thickness_sheet=6*1e-9;
width_sheet=10*1e-9;

widthFin= width_sheet; // for drain height calculation
effective_width=(thickness_sheet+width_sheet)*2;

capIdealGate = caplist[6] * 1E-18 / (effective_width*max_sheet_num); cap_draintotal = cap_draintotallist[6]/ (effective_width); capFringe = 0;

effectiveResistanceMultiplier = eff_res_mul[6]; /* from CACTI */
current_gmNmos= gm[6];
current_gmPmos= gm[6];
gm_oncurrent = gm[6]; // gm at on current

- Python wrapper that interfaces with PyTorch
 - inference.py: accuracy estimation w/ hardware effects
 - Input non-ideality effects of characterized devices + algorithmic parameters (dataset, weight/bit resolution)

Model/dataset/mode	<pre>parser.add_argument('dataset', default='cifar10', help='cifar10 cifar100 imagenet') parser.add_argument('model', default='VGG8', help='VGG8 DenseNet40 ResNet18') parser.add_argument('mode', default='WAGE', help='WAGE FP')</pre>
Resolution	parser.add_argument('wl_weight', type=int, default=8) parser.add_argument('wl_grad', type=int, default=8) parser.add_argument('wl_activate', type=int, default=8) parser.add_argument('wl_error', type=int, default=8) V1.4 update: user defines number of rows read # Hardware Properties
Resolution	<pre># if do not consider hardware effects, set inference=0 in parallel in the sub-arrays parser.add_argument('inference', type=int, default=0, help='run hardware inference simulation') parser.add_argument('subArray', type=int, default=128, help='size of subArray (e.g. 128*128)')</pre>
Hardware properties	<pre>parser.add_argument('parallelRead', type=int, default=128, help='number of rows read in parallel (<= subArray e.g. 32)') parser.add_argument('ADCprecision', type=int, default=5, help='ADC precision (e.g. 5-bit)') parser.add_argument('cellBit', type=int, default=1, help='cell precision (e.g. 4-bit/cell)') parser.add_argument('onoffratio', type=float, default=10, help='device on/off ratio (e.g. 6max/Gmin = 3)')</pre>
Nonideal effects	<pre># if do not run the device retention / conductance variation effects, set vari=0, v=0 parser.add_argument('vari', type=float, default=0, help='conductance variation (e.g. 0.1 standard deviation to generate random variation)') parser.add_argument('t', type=float, default=0, help='retention time')</pre>
(variation, retention)	parser.add_argument('v', type=float, default=0, help='drift coefficient') parser.add_argument('detect', type=int, default=0, help='if 1, fixed-direction drift, if 0, random drift') parser.add_argument('target', type=float, default=0, help='drift target for fixed-direction drift, range 0-1') current time = datetime.now().strftime('%Y %m %d %H %M %S')

• Circuit Modules in NeuroSim

3 master - DNN_NeuroSim_V1.3 / Inference_pyton	rch / NeuroSIM /	Go to file Add file - ···
neurosim Update Technology.cpp		168523b on Nov 6, 2022 🕥 History
C Adder.cpp	add hw effect for resnet	2 years ago
🗅 Adder.h	add hw effect for resnet	2 years ago
AdderTree.cpp	add hw effect for resnet	2 years ago
🗅 AdderTree.h	add hw effect for resnet	2 years ago
🗅 BitShifter.cpp	add hw effect for resnet	2 years ago
🗋 BitShifter.h	add hw effect for resnet	2 years ago
🗅 Buffer.cpp	Update Buffer.cpp	5 months ago
🖺 Buffer.h	add hw effect for resnet	2 years ago
🗅 Bus.cpp	Update Bus.cpp	2 years ago
🗅 Bus.h	add hw effect for resnet	2 years ago
Chip.cpp	Update Chip.cpp	3 months ago
Chip.h	add hw effect for resnet	2 years ago
Comparator.cpp	add hw effect for resnet	2 years ago
🗅 Comparator.h	add hw effect for resnet	2 years ago
CurrentSenseAmp.cpp	add hw effect for resnet	2 years ago
CurrentSenseAmp.h	add hw effect for resnet	2 years ago
DFF.cpp	add hw effect for resnet	2 years ago

Ex: 3-bit ripple carry adder



-> Extendible to N-bit ripple carry adder

Circuit Modules in NeuroSim

```
void Adder::Initialize(int _numBit, int _numAdder, double _clkFreq){
    if (initialized)
        cout << "[Adder] Warning: Already initialized!" << endl;</pre>
```

numBit = _numBit; numAdder = _numAdder; clkFreq = _clkFreq;

```
widthNandN = 2 * MIN_NMOS_SIZE * tech.featureSize;
widthNandP = tech.pnSizeRatio * MIN_NMOS_SIZE * tech.featureSize;
```

EnlargeSize(&widthNandN, &widthNandP, tech.featureSize*MAX_TRANSISTOR_HEIGHT, tech);

initialized = true;

Basic setting of the circuit is initialized at the beginning of the simulation based on the input parameters

-> transistor sizing/adder bit etc. are determined based on the input arguments

• Circuit Modules in NeuroSim

```
void Adder::Initialize(int _numBit, int _numAdder, double _clkFreq){
         if (initialized)
                 cout << "[Adder] Warning: Already initialized!" << endl;</pre>
         numBit = _numBit;
         numAdder = _numAdder;
         clkFreq = _clkFreq;
         widthNandN = 2 * MIN_NMOS_SIZE * tech.featureSize;
         widthNandP = tech.pnSizeRatio * MIN_NMOS_SIZE * tech.featureSize;
         EnlargeSize(&widthNandN, &widthNandP, tech.featureSize*MAX_TRANSISTOR_HEIGHT, tech);
         initialized = true;
void Adder::CalculateArea(double _newHeight, double _newWidth, AreaModify _option) { ···
void Adder::CalculateLatency(double _rampInput, double _capLoad, double numRead){...
void Adder::CalculatePower(double numRead, int numAdderPerOperation) { ···
```

Basic setting of the circuit is initialized at the beginning of the simulation based on the input parameters

-> transistor sizing/adder bit etc. are determined based on the input arguments

Each circuit module include area, latency and power calculations.

V1.4 update: calculations support GAA special layout

Circuit Modules in NeuroSim

```
void Adder::CalculateArea(double _newHeight, double _newWidth, AreaModify _option) {
    if (!initialized) {
       cout << "[Adder] Error: Require initialization first!" << endl;</pre>
    } else {
                                                                                     V1.4: Area for GAA special lavout
       double hNand, wNand;
       // 1.4 update: GAA special layout
       if ((tech.featureSize <= 2e-9) && param->speciallayout) {CalculateGateArea(NAND, 2, MIN NMOS SIZE * tech.featureSize,
                                                                MIN_NMOS_SIZE * tech.featureSize, tech.featureSize * MAX_TRANSISTOR_HEIGHT,
                                                                tech, &hNand, &wNand);}
       else {CalculateGateArea(NAND, 2, widthNandN, widthNandP, tech.featureSize * MAX TRANSISTOR HEIGHT, tech, &hNand, &wNand);}
       area = 0;
       height = 0:
       width = 0:
       if (_newHeight && _option==NONE) { // Adder in multiple columns given the total height
           hAdder = hNand;
           wAdder = wNand * 9 * numBit:
           if (hAdder > newHeight) {
               cout << "[Adder] Error: A single adder height is even larger than the assigned height ! " << endl;
           } else {
               height = newHeight;
               width = wAdder * hAdder * numAdder / newHeight;
        } else if (_newWidth && _option==NONE) { // Adder in multiple rows given the total width
           hAdder = hNand * numBit:
           wAdder = wNand * 9;
```

Area estimation code

Uses the typical CMOS layout rule for estimating the area and special rule for GAA.

-> Given the number of NMOS/PMOS & the maximum possible width/fin number of NMOS/PMOS of standard cell, the area is calculated.

-> Folding in the CPP direction is applied to support various circuit topologies considering the area efficiency.

89

• Circuit Modules in NeuroSim

<pre>void Adder::CalculateArea(double _newHeight, double _newWidth, AreaModify _option) {</pre>	
<pre>if (!initialized) {</pre>	
<pre>cout << "[Adder] Error: Require initialization first!" << endl;</pre>	
} else {	
double hNand, wNand;	
// 1.4 update: GAA special layout	
<pre>if ((tech.featureSize <= 2e-9) && param->speciallayout) {CalculateGateArea(NAND, 2, MIN_NMOS_SIZE * tech.featureSize, tech.featureSize, tech.featureSize, tech.featureSize, tech.featureSize, tech.featureSize</pre>	eatureSize, e * MAX_TRANSISTOR_HEIGHT,
else {CalculateGateArea(NAND, 2, widthNandN, widthNandP, tech.featureSize * MAX_TRANSISTOR_HEIGHT, tech, &	hNand, &wNand);}
area = 0;	
height = 0 ;	
width = 0;	
<pre>if (_newHeight && _option==NONE) { // Adder in multiple columns given the total height hAdder = hNand; wAdder = wNand * 9 * numBit;</pre>	
WHATCH = WHATA > Homosty	
<pre>if (hAdder > _newHeight) {</pre>	
<pre>cout << "[Adder] Error: A single adder height is even larger than the assigned height ! " << endl;</pre>	
} else {	
<pre>height = _newHeight;</pre>	
<pre>width = wAdder * hAdder * numAdder / _newHeight; }</pre>	Fig.
} else if (_newWidth && _option==NONE) { // Adder in multiple rows given the total width	Gus
hAdder = hNand * numBit;	Gus
wAdder = wNand * 9;	100
<pre>if (wAdder > _newWidth) {</pre>	
<pre>cout << "[Adder] Error: A single adder width is even larger than the assigned width ! " << endl;</pre>	
} else {	
<pre>width = _newWidth;</pre>	
height = wAdder * hAdder * numAdder / newWidth;	



Fig. 2. Example of transistor folding. [9]

Gustavo, et al. Toward Better Layout Design in ASTRAN CAD Tool by Using an Efficient Transistor Folding

Circuit Modules in NeuroSim

void Adder::CalculateLatency(double _rampInput, double _capLoad, double numRead){

```
if (!initialized) {
    cout << "[Adder] Error: Require initialization first!" << endl;
} else {
    readLatency = 0;
    rampInput = _rampInput;
    capLoad = _capLoad;
    double tr;    /* time constant */
    double gm;    /* transconductance */
    double beta;    /* for horowitz calculation */
    double resPullUp, resPullDown;
    double readLatencyIntermediate = 0;
    double ramp[10];
</pre>
```

ramp[0] = rampInput;

V1.4: Latency for GAA special layout

// 1.4 update: GAA special layout
// Calibration data pattern is A=1111111..., B=1000000 .. and Cin=1
// 1st
resPullDown = CalculateOnResistance(widthNandW, NMOS, inputParameter.temperature, tech) * 2;
if ((tech.featureSize <= 2e-9) && param->speciallayout)
{resPullDown = resPullDown *3.0/2.0;}
tr = resPullDown * (capNandOutput + capNandInput * 3);
gm = CalculateTransconductance(widthNandN, NMOS, tech);
beta = 1 / (resPullDown * gm);
readLatency += horowitz(tr, beta, ramp[0], &ramp[1]);

// 2nd

resPullUp = CalculateOnResistance(widthNandP, PMOS, inputParameter.temperature, tech); if ((tech.featureSize <= 2e-9) && param->speciallayout) {resPullUp= resPullUp *3.0/2.0;} tr = resPullUp * (capNandOutput + capNandInput * 2); gm = CalculateTransconductance(widthNandP, PMOS, tech); beta = 1 / (resPullUp * gm); readLatency += horowitz(tr, beta, ramp[1], &ramp[2]);

Latency estimation code

Horowitz delay estimation model is adopted for latency estimation. -> Given the input ramp rate/Elmore delay, the propagation delay/output ramp rate is calculated.



Input/output pulse has a certain ramp rate

Fig. 4. Delay of a logic gate driving a load through an RC line.

$$D = \underbrace{R_o \cdot (C_L + C_w) + p}_{\text{gate delay}} + \underbrace{R_w \cdot \left(\frac{C_w}{2} + C_L\right)}_{\text{wire delay}} \quad (1)$$

Bharadwaj S, et al. Speed and Power Scaling of SRAM's

91

Circuit Modules in NeuroSim

void Adder::CalculateLatency(double _rampInput, double _capLoad, double numRead){

```
if (!initialized) {
```

cout << "[Adder] Error: Require initialization first!" << endl;</pre>

} else {

```
readLatency = 0;
rampInput = _rampInput;
capLoad = _capLoad;
double tr; /* time constant */
double gm; /* transconductance */
double beta; /* for horowitz calculation */
double resPullUp, resPullDown;
double readLatencyIntermediate = 0;
double ramp[10];
```

1. Circuit RC parasitic modeling

Fig. 4. Delay of a logic gate driving a load through an RC line.



$$D = \underbrace{R_o \cdot (C_L + C_w) + p}_{\text{gate delay}} + \underbrace{R_w \cdot \left(\frac{C_w}{2} + C_L\right)}_{\text{wire delay}} \tag{1}$$

Bharadwaj S, et al. Speed and Power Scaling of SRAM's

Latency =
$$\tau_f \sqrt{\ln(v_s)^2 + \frac{2}{rampInput \times \tau_f}\beta(1-v_s)}$$
 (3)

rampOutput = $(1-v_s)/Latency$

Pai-Yu Chen, et al. NeuroSim: A Circuit-Level Macro Model for Benchmarking Neuro-Inspired Architectures in Online Learning

2. Elmore delay (D) calculation

3. Propagationdelay/output ramp rateestimation throughHorowitz delay model

ramp[0] = rampInput;

// 1.4 update: GAA special layout

// Calibration data pattern is A=1111111..., B=1000000... and Cin=1

resPullDown = CalculateOnResistance(widthNandN, NMOS, inputParameter.temperature, tech) * 2; if ((tech.featureSize <= 2e-9) && param->speciallayout) {resPullDown = resPullDown *3.0/2.0;} tr = resPullDown * (capNandOutput + capNandInput * 3); gm = CalculateTransconductance(widthNandN, NMOS, tech); beta = 1 / (resPullDown * gm); readLatency += horowitz(tr, beta, ramp[0], &ramp[1]);

// 2nd

resPullUp = CalculateOnResistance(widthNandP, PMOS, inputParameter.temperature, tech); if ((tech.featureSize <= 2e-9) && param->speciallayout) {resPullUp= resPullUp *3.0/2.0;} tr = resPullUp * (capNandOutput + capNandInput * 2); gm = CalculateTransconductance(widthNandP, PMOS, tech); beta = 1 / (resPullUp * gm); readLatency += horowitz(tr, beta, ramp[1], &ramp[2]);

• Circuit Modules in NeuroSim

void

A	dder::CalculatePower(double numKead, int numAdderPerOperation) {
if	(!initialized) {
	<pre>cout << "[Adder] Error: Require initialization first!" << endl;</pre>
}	else {
	leakage = 0;
	readDynamicEnergy = 0;
	V1 4 [.] Power for GAA special layout
	// 1.4 update - updated
	/* Leakage power */
	<pre>if ((tech.featureSize == 2e-9) && param->speciallayout) {</pre>
	/* Leakage power */
	<pre>leakage += CalculateGateLeakage(NAND, 2, MIN_NMOS_SIZE * tech.featureSize, MIN_NMOS_SIZE * tech.featureSize,</pre>
	inputParameter.temperature, tech) * tech.vdd * 9 * numBit * numAdder * 2/3; }
	else if ((tech.featureSize == 1e-9) && param->speciallayout) {
	/* Leakage power */
	leakage += CalculateGateLeakage(NAND, 2, MIN_MMOS_SIZE * tech.teatureSize, MIN_MMOS_SIZE * tech.teatureSize,
	inputParameter.temperature, tech.vdd * 9 * numBit * numAdder * 2/3;}
	/* Leakage power */
	Teakage +- CalculateGateLeakage(WAHD) 2, WitchWahdW, WitchWahdW, MitchWahdW, * aka waka * aka * ak
	1
	r r
	/* Read Dynamic energy */
	// Calibration data pattern of critical path is A=1111111 B=1000000 and Cin=1
	// Only count 0 to 1 transition for energy
	// First stage
	readDynamicEnergy += (capNandInput * 6) * tech.vdd * tech.vdd: // Input of 1 and 2 and Cin
	readDynamicEnergy += (capNandOutput * 2) * tech.vdd * tech.vdd; // Output of S[0] and 5
	// Second and later stages
	readDynamicEnergy += (capNandInput * 7) * tech.vdd * tech.vdd * (numBit-1);
	readDynamicEnergy += (capNandOutput * 3) * tech.vdd * tech.vdd * (numBit-1);
	// Hidden transition
	// First stage
	readDvnamicEnergy += (capNandOutput + capNandInput) * tech.vdd * tech.vdd * 2: // #2 and #3

Power estimation code

- 1. Dynamic Energy
- -> $C_{node}V_{dd}^2$ * (switching factor) per each cycle is summed up <u>for all nodes</u>.
- -> switching factor considers the average switching event during a cycle.

2. Leakage Energy

-> I_{off} * V_{dd} * (averaging factor) * (Idle time) is summed up <u>for all gates.</u>

-> the averaging factor considers the average leakage energy of the transistors in a gate. It is dependent on the type of gate/number of inputs.

readDynamicEnergy += (capNandOutput + capNandInput * 2) * tech.vdd * tech.vdd; // #4

- Optional: self-defined network models
 - Create your own file in "models" folder: Replace nn.Conv2d / nn.Linear as QConv2d / QLinear self.conv1 = make layers([('C', 3, num planes, 3, 'same', 1)], args, logger)
 - make_layers function (just call it, no need to change)

```
def make_layers(cfg, args, logger):
   global name
    layers = []
    for i, v in enumerate(cfg):
       if v[0] == 'M':
           layers += [nn.MaxPool2d(kernel size=v[1], stride=v[2])]
        if v[0] == 'C':
           in channels = v[1]
           out channels = v[2]
           if v[4] == 'same':
               padding = v[3]//2
            else:
           if args.mode == "WAGE":
                conv2d = QConv2d(in channels, out channels, kernel size=v[3], stride=v[5], padding=padding,
                                 logger=logger,wl input = args.wl activate,wl activate=args.wl activate,
                                 wl error=args.wl error,wl weight= args.wl weight, inference=args.inference, onoffratio=args.onoffratio, cellBit=args.cellBit,
                                 subArray=args.subArray,ADCprecision=args.ADCprecision,vari=args.vari,t=args.t,v=args.v,detect=args.detect,target=args.target,
                                 name = 'Conv'+str(name)+' ')
           elif args.mode == "FP":
               conv2d = nn.Conv2d(in channels, out channels, kernel size=v[3], stride=v[5], padding=padding, bias=False)
            name += 1
           batchnorm = nn.BatchNorm2d(out channels)
           non linearity activation = nn.ReLU()
           layers += [conv2d, batchnorm, non linearity activation]
            in channels = out channels
        if v[0] == 'L':
            if args.mode == "WAGE":
               linear = QLinear(in_features=v[1], out_features=v[2],
                                logger=logger, wl input = args.wl activate, wl activate=args.wl activate, wl error=args.wl error,
                                wl weight=args.wl weight, inference=args.inference, onoffratio=args.onoffratio, cellBit=args.cellBit,
                                subArray=args.subArray,ADCprecision=args.ADCprecision,vari=args.vari,t=args.t,v=args.v,detect=args.detect,target=args.target,
                                name='FC'+str(i)+' ' )
            elif args.mode == "FP":
               linear = nn.Linear(in_features=v[1], out_features=v[2], bias=False)
           layers += [linear]
    return nn.Sequential(*layers)
```

- Optional: self-defined network models
 - Create your own NetWork.csv

Example: VGG-8

	IFM Length	IFM Width	IFM Channel Depth	Kernel Length	Kernel Width	Kernel Depth	Followed by pooling or not?
Layer 1	32	32	3	3	3	128	0
Layer 2	32	32	128	3	3	128	1
Layer 3	16	16	128	3	3	256	0
Layer 4	16	16	256	3	3	256	1
Layer 5	8	8	256	3	3	512	0
Layer 6	8	8	512	3	3	512	1
Layer 7	1	1	8192	1	1	1024	0
Layer 8	1	1	1024	1	1	10	0

• Workflow:



• Output example:

 Accuracy 	Test set: Average loss: 1.5865, Accuracy: 9007/10000 (90%) FloorPlan
Floorplan	Tile and PE size are optimized to maximize memory utilization (= memory mapped by synapse / total memory on chip)
	Desired Conventional Mapped Tile Storage Size: 1024x1024
	Desired Conventional PE Storage Size: 512x512
	Desired Novel Mapped Tile Storage Size: 9x512x512
	User-defined SubArray Size: 128x128
	# of tile used for each layer
	layer1: 1
	layer2: 1
	Speed-up of each layer
	layer1: 16
	layer2: 4
	Utilization of each layer
	layer1: 0.210938
	layer2: 1
	Memory Utilization of Whole Chip: 96.8584 %
	FloorPlan Done

• Output example:

• each layer performance breakdown

Hardware Performance
Estimation of Layer 1
layer1's readLatency is: 377982ns
layer1's readDynamicEnergy is: 1.89268e+06pJ
layer1's leakagePower is: 8.47407uW
layer1's leakageEnergy is: 147340pJ
layer1's buffer latency is: 337150ns
layer1's buffer readDynamicEnergy is: 22681.5pJ
layer1's ic latency is: 24831ns
layer1's ic readDynamicEnergy is: 487920pJ

Latency & energy of:
➢ Whole layer
➢ Buffer
➢ Interconnection
➢ ADC
➢ Accumulation
➢ Other
Leakage

----- ADC (or S/As and precharger for SRAM) readLatency is : 7385.06ns

------ Accumulation Circuits (subarray level: adders, shiftAdds; PE/Tile/Global level: accumulation units) readLatency is : 7385.06ns

----- Other Peripheries (e.g. decoders, mux, switchmatrix, buffers, IC, pooling and activation units) readLatency is : 363212ns

----- ADC (or S/As and precharger for SRAM) readDynamicEnergy is : 1.03834e+06pJ

------ Accumulation Circuits (subarray level: adders, shiftAdds; PE/Tile/Global level: accumulation units) readDynamicEnergy is : 155498pJ

------ Other Peripheries (e.g. decoders, mux, switchmatrix, buffers, IC, pooling and activation units) readDynamicEnergy is : 698849pJ

- Output examples:
 - System-level summary

ChipArea : 4.5554e+07um^2 Chip total CIM array : 1.5408e+06um^2 Total IC Area on chip (Global and Tile/PE local): 6.17954e+06um^2 Total ADC (or S/As and precharger for SRAM) Area on chip : 1.46062e+07um^2 Total Accumulation Circuits (subarray level: adders, shiftAdds; PE/Tile/Global level: accumulation units) on chip : 3.13322e+06um^2 Other Peripheries (e.g. decoders, mux, switchmatrix, buffers, pooling and activation units) : 2.00943e+07um^2

Chip clock period is: 2.05141ns Chip layer-by-layer readLatency (per image) is: 1.45613e+06ns Chip total readDynamicEnergy is: 4.03921e+07pJ Chip total leakage Energy is: 1.19011e+06pJ Chip total leakage Power is: 628.781uW Chip buffer readLatency is: 1.00082e+06ns Chip buffer readDynamicEnergy is: 412969pJ Chip ic readLatency is: 141884ns Chip ic readDynamicEnergy is: 8.94262e+06pJ

----- ADC (or S/As and precharger for SRAM) readLatency is : 63462.3ns

- ------ Accumulation Circuits (subarray level: adders, shiftAdds; PE/Tile/Global level: accumulation units) readLatency is : 246001ns
- ------Other Peripheries (e.g. decoders, mux, switchmatrix, buffers, IC, pooling and activation units) readLatency is : 1.14667e+06ns
- ------ ADC (or S/As and precharger for SRAM) readDynamicEnergy is : 2.08642e+07pJ

------ Accumulation Circuits (subarray level: adders, shiftAdds; PE/Tile/Global level: accumulation units) readDynamicEnergy is : 5.16853e+06pJ

------ Other Peripheries (e.g. decoders, mux, switchmatrix, buffers, IC, pooling and activation units) readDynamicEnergy is : 1.43594e+07pJ

- Energy efficiency
- > Throughput
- Compute efficiency

Georgia School of Electrical and Tech Computer Engineering Area

Latency & energy of:

- > Whole chip
- > Buffer
- Interconnection
- > ADC
- > Accumulation
- > Other peripheries

Leakage

References

- J. Lee, A. Lu, W. Li, **S. Yu**, "NeuroSim V1.4: Extending technology support for digital compute-in-memory towards 1nm node", **IEEE TCAS-I**, accepted.
- A. Lu, J. Lee, T.-H. Kim, M. Karim, R. Park, H. Simka, **S. Yu**, "High-speed emerging memories for Al hardware accelerators," **Nature Reviews Electrical Engineering**, accepted.
- Y. Luo, **S. Yu**, "H3D-Transformer: A heterogeneous 3D (H3D) computing platform for transformer model acceleration on edge devices," **ACM TODAES**, accepted.
- P.-K. Hsu, V. Garg, A. Lu, S. Yu, "A heterogeneous platform for 3D NAND-based in-memory hyperdimensional computing engine for genome sequencing applications," IEEE TCAS-I, accepted.
- A. Lu, X. Peng, W. Li, H. Jiang, **S. Yu**, "NeuroSim simulator for compute-in-memory hardware accelerator: validation and benchmark," **Frontiers in Artificial Intelligence**, vol. 4, 659060, 2021.
- X. Peng, S. Huang, H. Jiang, A. Lu, S. Yu, "DNN+NeuroSim V2.0: An end-to-end benchmarking framework for compute-in-memory accelerators for on-chip training," IEEE Trans. CAD, vol. 40, no. 11, pp. 2306-2319, 2021.
- X. Peng, W. Chakraborty, A. Kaul, W. Shim, M. S. Bakir, S. Datta, S. Yu, "Benchmarking monolithic 3D integration for compute-in-memory accelerators: overcoming ADC bottlenecks and maintaining scalability to 7nm or beyond," IEEE International Electron Devices Meeting (IEDM) 2020.
- X. Peng, S. Huang, Y. Luo, X. Sun, S. Yu, "DNN+NeuroSim: An end-to-end benchmarking framework for compute-in-memory accelerators with versatile device technologies," IEEE International Electron Devices Meeting (IEDM) 2019.

Acknowledgement



100