



A Preliminary Comparative Analysis of Face-to-Face Bonded Long-Term RAM (LtRAM) and SRAM Chipllets vs. Off-Chip HBM for Agentic Inference

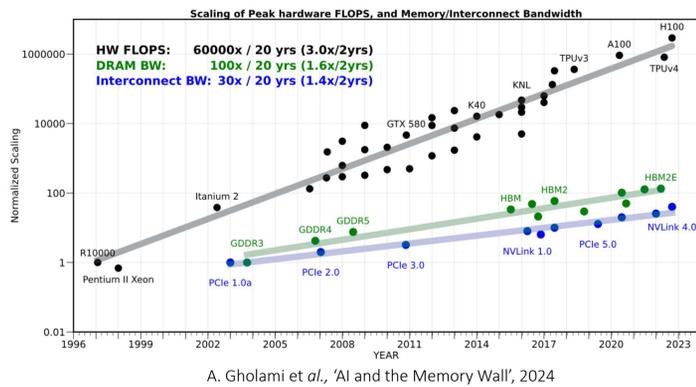
Samuel A. Dayo, Alex Aranburu*, Victor J.B. Jung, Colin H. Drewes, Shridhar Mukund*

Prof. Robert M. Radway, Prof. Christos Kozyrakis, Prof. Luca Benini, Prof. H.S Philip Wong, Prof. Subhasish Mitra

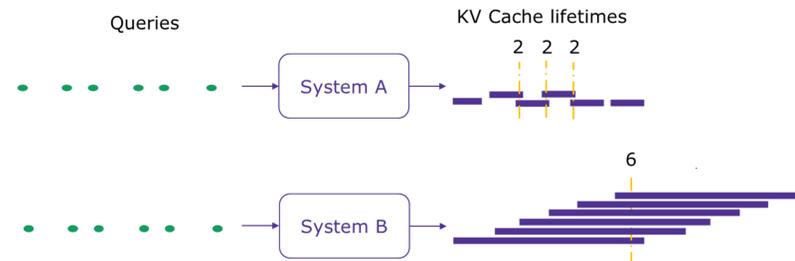


*EMD Electronics

Motivations: Memory-Bandwidth Bottleneck



Motivations: Reduce Active Memory Footprint



Quantitative Analysis Results

Performance Metrics:

Query Response Time

Time to complete one inference query (I + O)

Aggregate Token Throughput

Tokens processed per second (I + O)

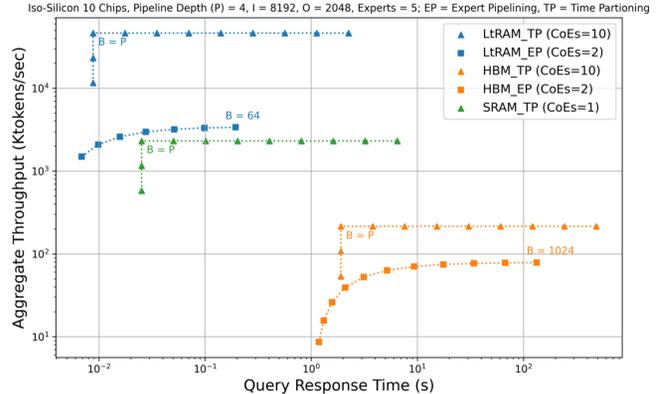
Memory Capacity per Query Throughput

Total memory (KV cache + weights) needed per query/sec

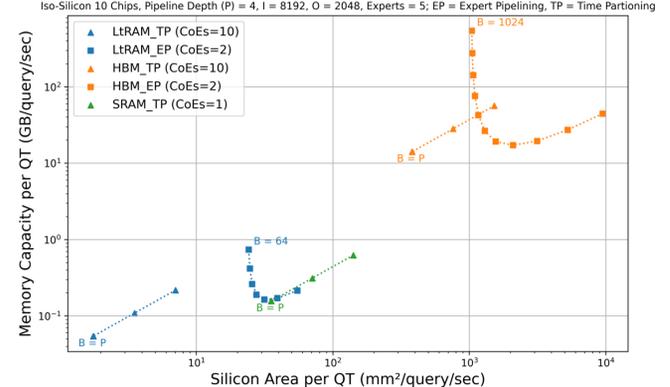
Silicon Area per Query Throughput

Total accelerator area (mm²) needed per query/sec

Aggregate Throughput vs Query Response Time

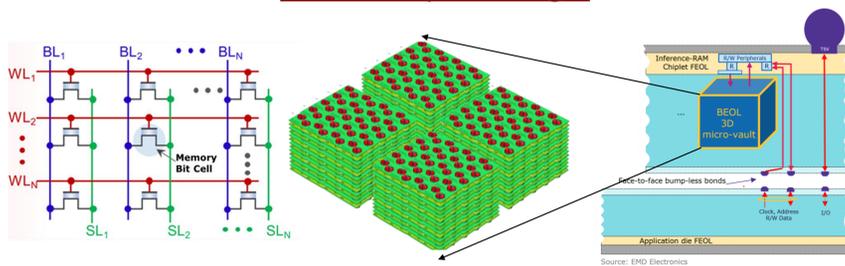


Memory Capacity vs Silicon Area per Query



Long-term RAM (LtRAM)

LtRAM Chipllet Design:



LtRAM Chipllet Summary:

Feature	Specification
Density & Granularity	0.32 GB/mm ² ; 1250 ports/mm ² ; 32 KB Banks
Area-scalable Bandwidth	< 2.32 TB/s/mm ² (scalable from 40 mm ² to wafer-scale)
Read Latency	< 8ns (4-stage pipelined random-access, < 2ns per stage)
Read Energy	~0.01 pJ/b
Read Voltage	< 500 mV
Write Latency	< 16 ns (8-cycle write)
Write Energy	0.04 - 0.06 pJ/b
Write Voltage	1.5 V
Write Endurance	10 ¹² cycles

LtRAM Chipllet Density and Energy Benefits:

$$\frac{\text{Memory Bandwidth}}{\text{mm}^2} \geq \frac{\text{Max MACOps/sec}}{\text{mm}^2} \geq \left(\frac{\text{Memory Capacity (Weights + KV Cache)}}{\text{mm}^2} \right) \downarrow \left(\frac{\text{Number of Tokens/sec}}{\text{mm}^2} \right) \uparrow$$

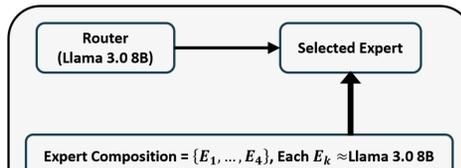
Spatially Scalable Key Performance Indicators (KPIs)	On-Chip F2F Hybrid Bonded Memory Type			Off-Chip HBM
	DRAM Wide I/O ports	SRAM	Inference RAM	Perimeter Limited 4 per 800 mm ² 1.6 TB/sec per HBM
Port Density (ports/mm ²)	Very Low ~0.01 8b per 800 mm ²	Low ~83 32KB Macro	High 1250 32KB Macro	Very Low 0.005
Data Granularity (bytes/transaction)	Coarse 2,048	Fine 4, flexible	Fine 4, flexible	Coarse 2,048
Memory Density (Usable Gb/mm ²)	High ~0.2	Very Low ~0.025	Higher 0.32	Highest ~0.2 x 12
Bandwidth Density (TB/sec per mm ²)	Low ~0.016	Medium ~0.116	High 2.32	Very Low 0.008
Read Energy (pJ/bit)	High ~1	Low ~0.01	Low 0.01	Very High ~20

© Winbond DRAM CUBE
 2 On the Model of Computation, William Dally, 2022
 3 Best guess based on known granularity limited by peripheral overhead and Serdes interface

Quantitative Analysis Framework

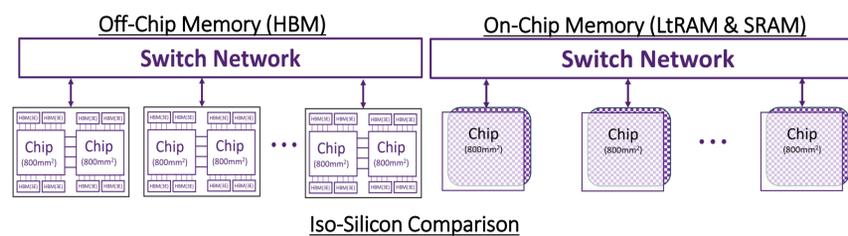
Agentic Workload Modeling:

10 Expert Composition of Experts



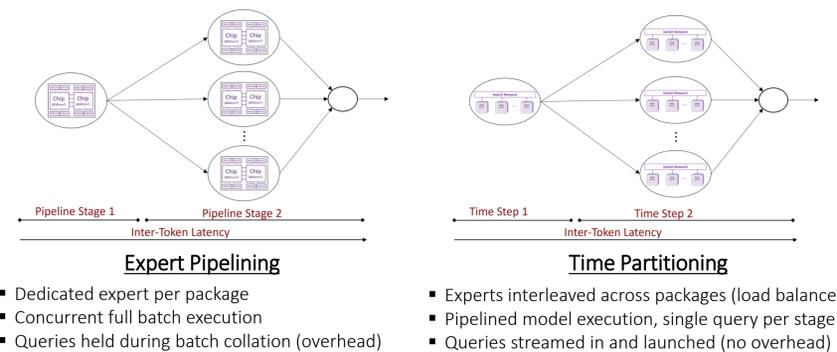
S. Jain et al., "Composition of Experts: A Modular Compound AI System Leveraging Large Language Models", 2024

Memory-Compute Integration:



Iso-Silicon Comparison

Composition of Experts Execution Modes:



Long-Term RAM – Key Benefits

- Compute-Bound Query Processing**
 - Enables sustained throughput scaling at near iso-latency
 - Unlocks peak compute utilization across pipeline depths
- Compact Active Memory Footprint**
 - Reduced per-query memory footprint
 - Multi-query concurrency with minimal memory overhead
- Throughput-Dense Silicon**
 - Near-linear throughput scaling with chip count
 - Iso-throughput targets with significantly fewer accelerators



BEOL-Compatible IGZO FeFET Nanodendrite Enabling Sequence Decoding for Dendrocentric Computing

Hugo J.-Y. Chen¹, Chi-Hsin Huang¹, Kwabena Boahen^{1,2,3}, and H.-S. Philip Wong¹

¹Department of Electrical Engineering, ²Department of Bioengineering, ³Department of Computer Science, Stanford University, Stanford, CA, USA

Background knowledge & Motivation

Goal: New energy-efficient AI hardware for neuromorphic computing

Accelerate data retrieval in language model

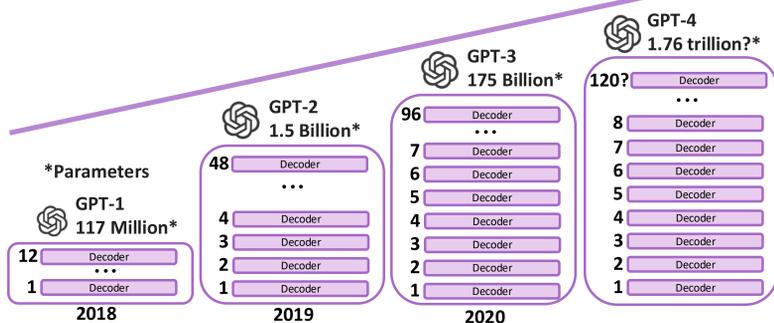


Fig. 1. The scaling of parameters and decoder layers in large language model [1]

Run GPT on the edge

- Needs a **15X faster** processor and a **170X bigger** battery to run GPT-3*
*compared to today's smartphone
- High memory density
- Parallel operation in 3D chip

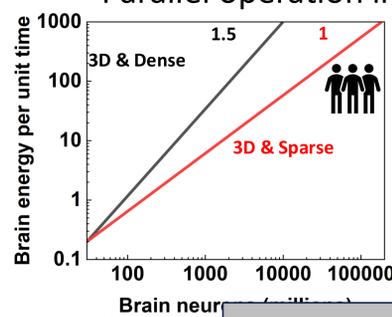
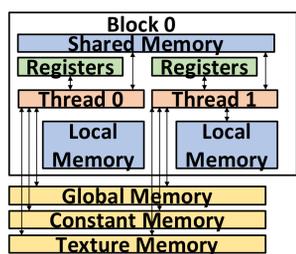


Fig. 2. 3D and sparse

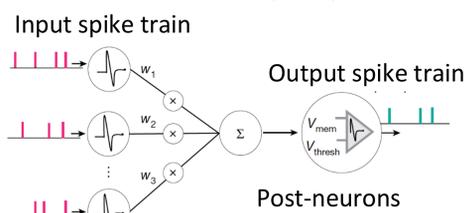
- Heat **generates** as **volume**: n^3
- Heat **dissipates** as **surface**: n^2
→ scales as power 1.5 → overheat → series operation
- Applied pulse to m out of n wires:
 - Keep m constant, increase n → activity scales as $1/n$
- Heat **generates** as **volume**: n^2
- Heat **dissipates** as **surface**: n^2
→ scales as power 1 → **parallel** operation

Synaptocentric (GPU) [3]



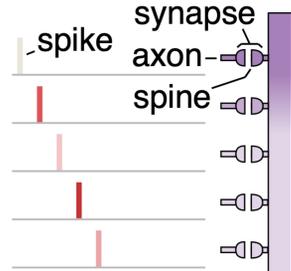
- Matrix multiplication
 - 2D chip
 - Synchronized clock

Axocentric (SNN) [4]



- Pre-neurons
 - Time dependent
 - Event driven
 - Run in parallel but overheat

Dendrocentric [5]



- Sequence dependent
 - 3D chip
 - Run in parallel w/o overheat

Fig. 3. AI hardware. Synaptocentric, Axocentric, and Dendrocentric

(preliminary unpublished results)

Experimental results

(preliminary unpublished results)

Conclusions & References

(preliminary unpublished results)



GEMMA: Gain Cell Embedded DRAM-based MAMBA-like Model Accelerator

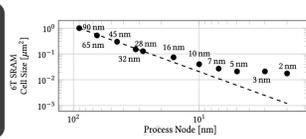
Hery Shin (hshin24@stanford.edu), Thierry Tamba (ttamba@stanford.edu)

Motivation

1. GCRAM: An Emerging Solution for the Growing On-Chip Memory Demands of ML Applications



[On-chip memory requirement for AI workloads [1]]

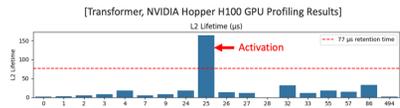


[End of SRAM scaling [2]]

	SRAM	Short-term RAM
Structure	6T	2T or 3T gain cells
Benefits	Fast, easy to integrate, low static power	Dense, low energy
Drawbacks	Sparse	Short retention times, expensive refreshes, active research
Uses	Fast R/W caches	Fast write-and-read operations

[Alternative memory technology]

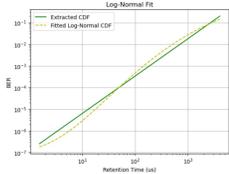
2. Analysis of Kernels Based on Data Lifetime



- Kernels such as nonlinear functions, residual layers are categorized as generating long-lived data [3].

3. Worst-case Retention-based Static Refresh Policy

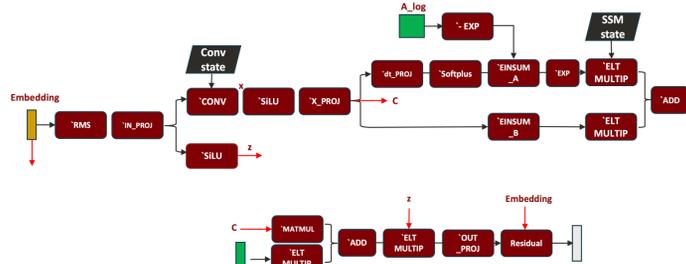
- Over 99% of cells exhibit retention times longer than the worst-case refresh period [4].
- Opportunities to minimize refresh overhead for both long-lived and short-lived data



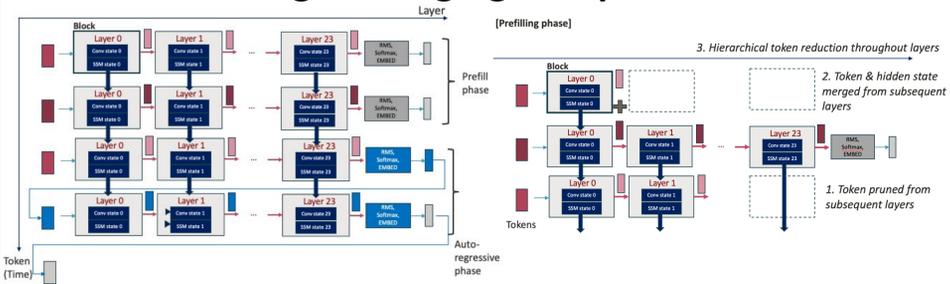
Data Lifetime-Aware HW & SW Co-Design

1. Data Lifetime-Aware Scheduling

- Kernel categorization based on data lifetime

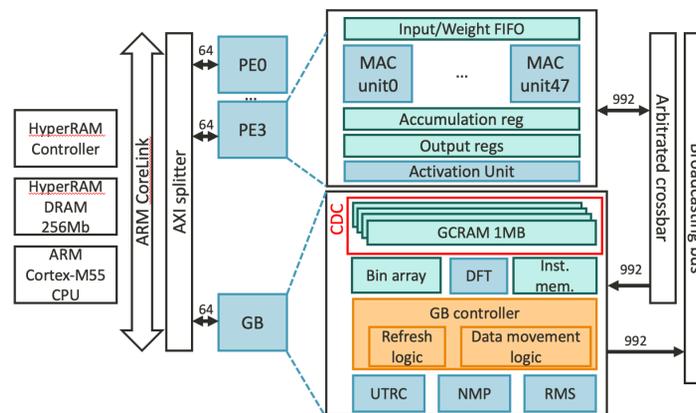


2. Token Pruning & Merging Compression



Overall Architecture

GEMMA Architecture



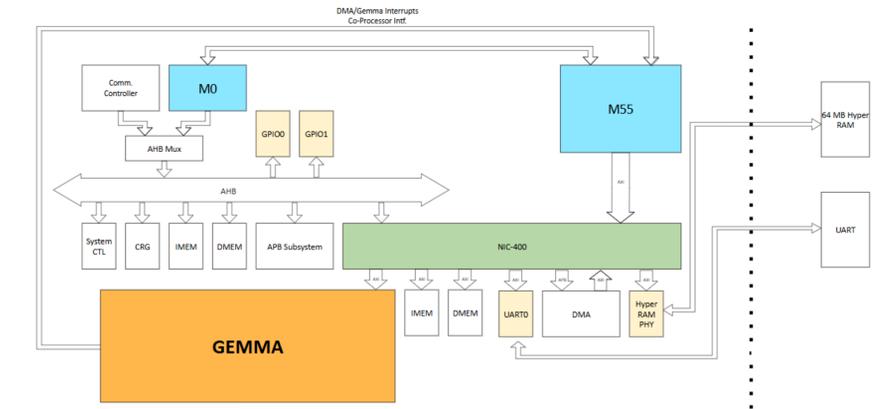
[GEMMA Accelerator Schematic]

1. SIMD Architecture

- Output stationary GEMV PEs

2. Microscaling-Adaptive Floating Point Hybrid Data Format

- Compress weight and activation to ~5 bits per element



[SoC Integration Schematic]

3. Near Memory Processor in GB & Activation Units in PEs

- Reduce the data movement between GB and PEs
- Reduce memory access through kernel fusion

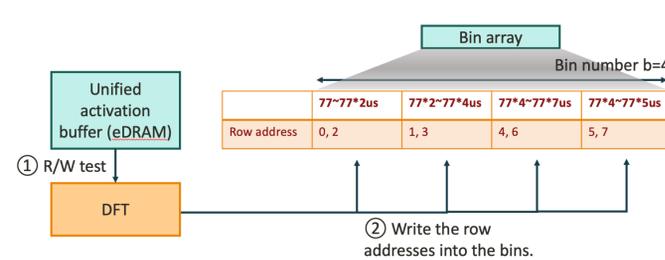
4. Clock Domain Crossing between EDRAM and Accelerator

- Reduce data lifetime of activations

Fine-grained Refresh Policy

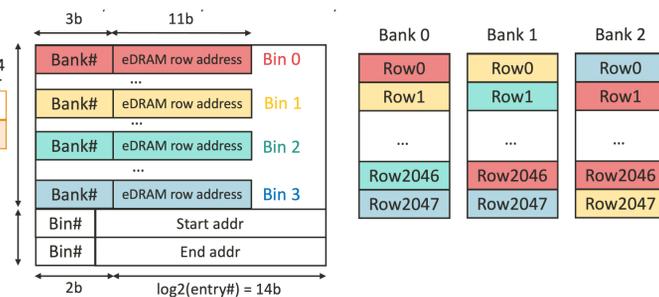
1. Design For Test (DFT)

- Row-wise retention estimation before runtime

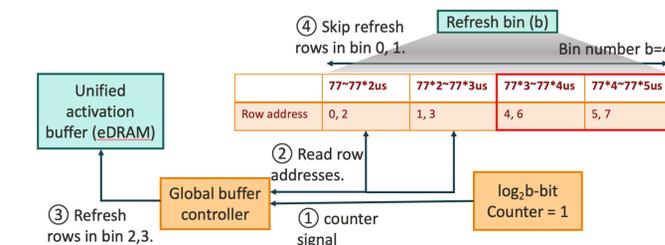


2. Row-wise Refresh With Bin Array

- Virtual-physical address translation book



3. GCRAM Metadata



- Metadata for each GCRAM row: lifetime indicator + error correction code
- Long-lived data mapped onto long retention time rows.
- Short-lived data can skip its refresh or be overwritten after its lifetime.

Conclusions & References

- Seamlessly integrate Si-based 3T gain cell RAM (GCRAM) into state space model accelerators to increase on-chip memory density.
- Challenge of Si-GCRAM: Refresh overhead due to short retention time
- **Data Lifetime minimizing HW&SW co-design:** Various techniques to shorten data lifetime of stored weight/activation by clock domain crossing, scheduling, token compression, etc.
- **Retention-aware Fine-grained Refresh Policy:** Row-wise refresh management reduces refresh overhead and enable lifetime-dependent data mapping in the GCRAM.

[1] A. Gholami, et al., "AI and Memory Wall," IEEE MICRO, Mar. 2024.
 [2] K. Zhang, "1.1 Semiconductor Industry: Present & Future," in 2024 IEEE International Solid-State Circuits Conference (ISSCC), Feb. 2024.
 [3] P. Li, T. Tamba, "GainSight: Profile-Guided Design for Heterogeneous On-Chip Memories in Next-Generation AI Accelerators," IEEE MICRO
 [4] RAAAM



GainSight: Application-Guided Profiling for Composing Heterogeneous On-Chip Memories in AI Hardware Accelerators

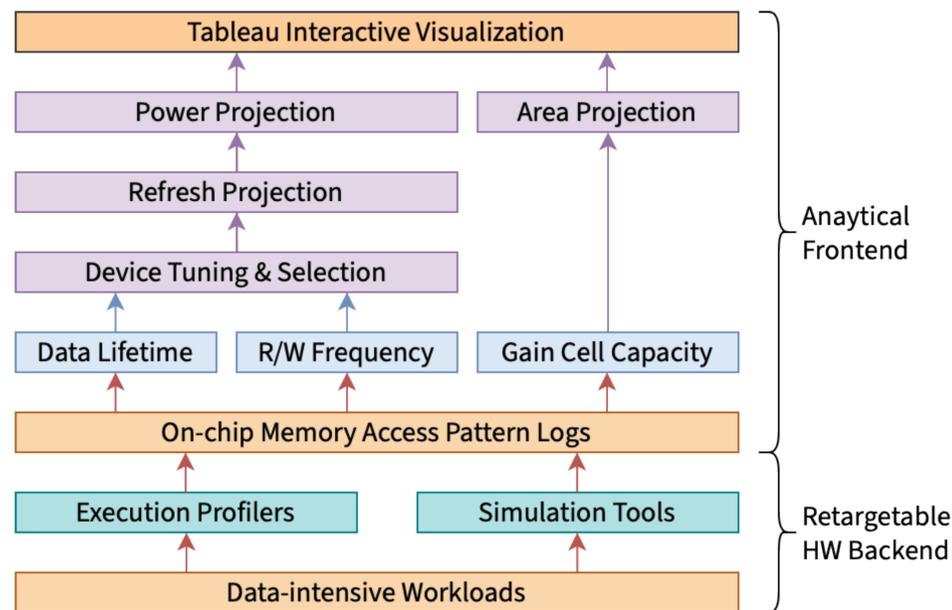
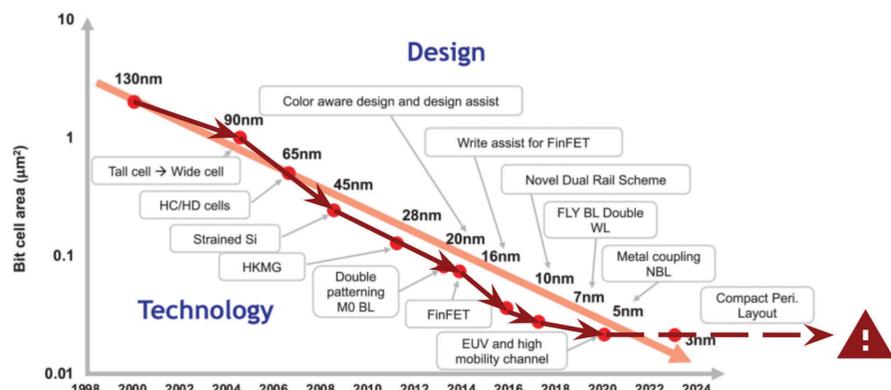
Peijing Li (peli@stanford.edu), Thierry Tamba (ttambe@stanford.edu)
Department of Electrical Engineering, Stanford University, Stanford, CA, United States

Abstract -- Visit our homepage on gainsight.stanford.edu!

As AI workloads drive soaring memory requirements, there is a need for higher-density on-chip memory for domain-specific accelerators that goes beyond what current SRAM technology can provide. We motivate that algorithms and application behavior should guide the composition of heterogeneous on-chip memories. However, there has been little work in factoring dynamic application profiles into such design decisions. We present GainSight, a profiling framework that analyzes fine-grained memory access patterns and computes data lifetimes in domain-specific accelerators. By combining instrumentation and simulation across retargetable hardware backends, GainSight aligns heterogeneous memory designs with workload-specific traffic and lifetime metrics.

Introduction

- Motivation: the “memory capacity wall” for transformer-based AI models; the stopping of SRAM scaling
- Long and short-lived data call for heterogeneous on-chip memory
- Devices such as gain cell random access memory (GCRAM) are promising candidates for short-lived data, but their short retention times and refresh costs pose challenges.
- Need to quantitatively correlate dynamic data lifetime profiles of workloads with characteristics of long and short-term memory
- **We conceived the GainSight Profiler to tackle these challenges, to analyze the memory access patterns and data lifetimes of software running on domain-specific accelerators, in order to inform the design of heterogeneous on-chip memories.**



Retargetable Hardware Backend

- Acquire fine-grained, on-chip memory access patterns for a physical or simulated architecture
- Currently implemented simulators for NVIDIA H100 GPUs (Accel-Sim) and systolic arrays (SCALE-Sim-v2)
- Exploit workload sampling to reduce program sizes to speed up cycle-accurate simulation

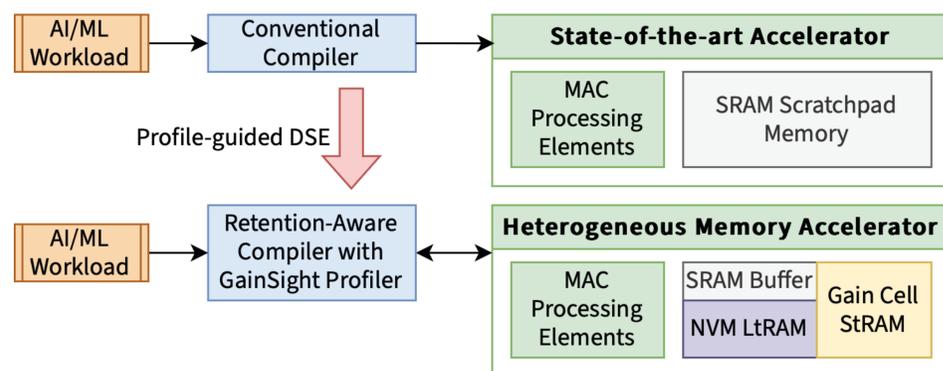
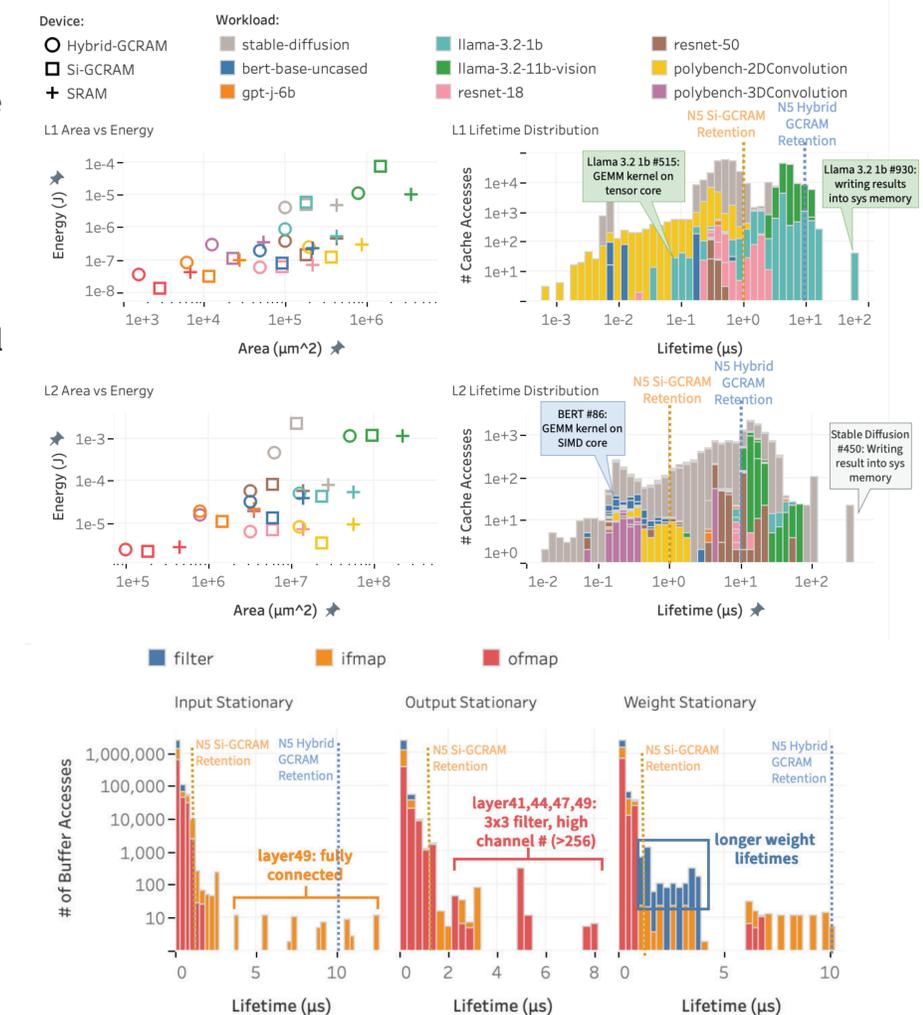
Analytical Frontend

- Consume memory traces generated by backend
- Calculate address-level data lifetimes, read/write frequencies, and on-chip memory capacity utilization
- Visualize lifetime distribution and compare with silicon gain cell retention times in Tableau
- Project number of refreshes, area, and total energy using calculated memory statistics for devices of SRAM, fully Si-Si GCRAM, and hybrid Si-ITO GCRAM.

Case Study

Insights from profiling MLPerf Inference and PolyBench workloads using simulated GPUs and systolic arrays:

1. 40% of L1 and 18% of L2 GPU cache accesses, and 79% of systolic array scratchpad accesses across profiled workloads are short-lived and suitable for Si-GCRAM
2. Si-GCRAM reduces active energy by 11-28% compared to SRAM;
3. Up to 90% of GPU cache fetches are never reused, highlighting inefficiencies in terms of cache pollution.





OpenGCRAM: An Open-Source Gain Cell Compiler Enabling Design-Space Exploration for AI Workloads

Xinxin Wang¹, Lixian Yan¹, Shuhan Liu¹, Luke Upton¹, Zhuoqi Cai¹, Yiming Tan¹, Shengman Li¹, Peijing Li¹, Jesse Cirimelli-Low², Thierry Tambe¹, Matthew Guthaus², and H.-S. Philip Wong¹
¹Stanford University ²University of California, Santa Cruz

I. Motivation:

- Memory Wall Problem [1]
- Gain Cell Memory (GCRAM)
 - vs SRAM: High density, high bandwidth, ultra-low leakage [2]

OS-OS gain cell

Higher capacity

Hybrid gain cell

Higher speed

To enable **fast, accurate, customizable, optimized** Gain Cell bank generation and performance simulation targeting for AI workloads:
We need a Gain Cell compiler

II. Methodology:

- Port memory compiler [3] to new PDK
 - Step 1: technology script
 - tech.py
 - Device libs or models
 - Layer definitions and DRC rules
 - Spice parameters
 - Step 2: customized bitcell and modules
 - Write Spice Netlists
 - Schematic and layout design
 - Export GDSII files
 - Step 3: fix DRC & LVS errors
 - DRC & LVS Check
 - Fix errors
 - Clean? (Yes/No)
 - Done
- Add support for new memory technology
 - Step 1: cell-level and array-level configuration
 - Step 2: peripheral circuit connection
 - Step 3: peripheral circuit redesign

III. OpenGCRAM

- GCRAM vs SRAM: schematic
 - 2T Si-Si GCRAM
 - 2T OS-OS GCRAM
 - 6T SRAM

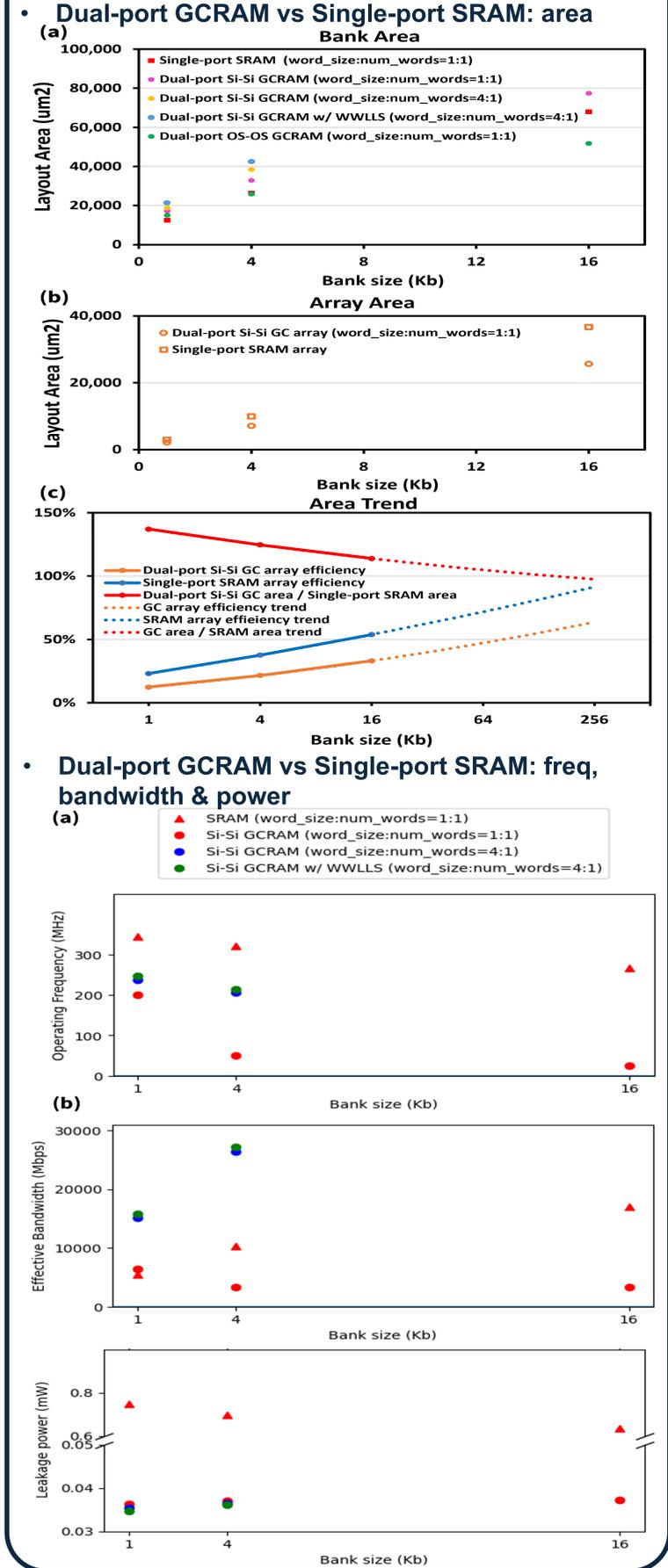
Acknowledgement: This work was supported in part by SRC JUMP 2.0 PRISM and CHIMES centers, Stanford SystemX Alliance, DAM project, and NMTRI.
References: [1] A. Gholami, Micro 2024; [2] S. Liu, IEDM 2023; [3] M. R. Guthaus, et al., ICCAD 2016; [4] P. Li, et al., arXiv:2504.14866

III. OpenGCRAM:

- GCRAM vs SRAM: layout
 - 2T Si-Si GCRAM
 - 2T OS-OS GCRAM
 - 6T-SRAM
- GCRAM bank architecture [2]

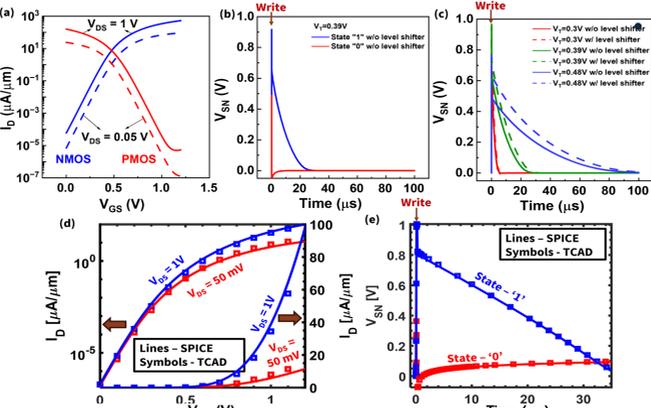
- Example GCRAM bank layout (32x32)

III. Results:

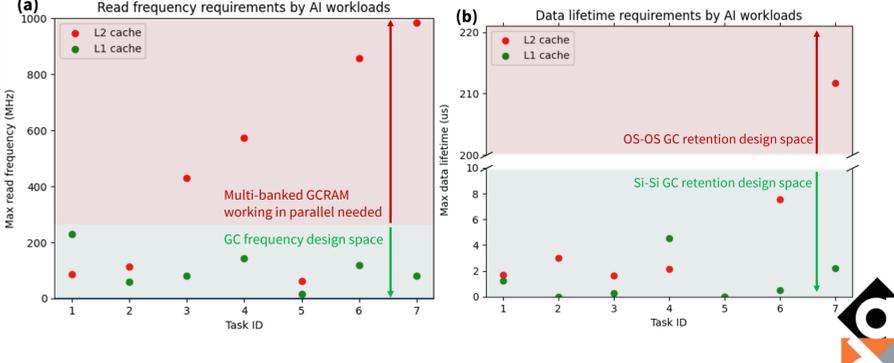


IV. Design Space Exploration:

GCRAM retention can be boosted with higher V_{th} or with WWL level shifter



OpenGCRAM enables design space exploration for AI tasks [4]



N-P Oxide Semiconductor Complementary Gain Cell (CGC) Memory

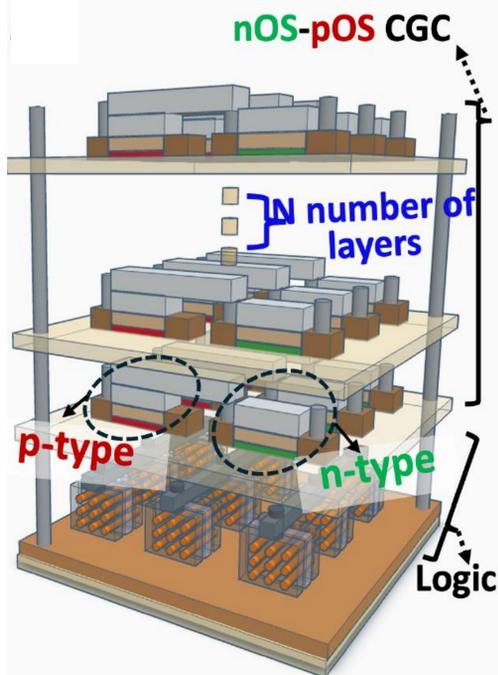
F. F. Athena, E. Ambrosi, K. Jana, Ch. H. Wu, J. Hartono, Y. M. Lee, C. C. Kuo, S. Liu, B. Saini, C. C. Wang, C. F. Hsu, G. Zeevi, X. Wang, J. Kang, E. Pop, T. Y. Lee, P. C. McIntyre, H.-S. P. Wong, X. Y. Bao

(contact: fathena@stanford.edu)

PI: H.-S. Philip Wong

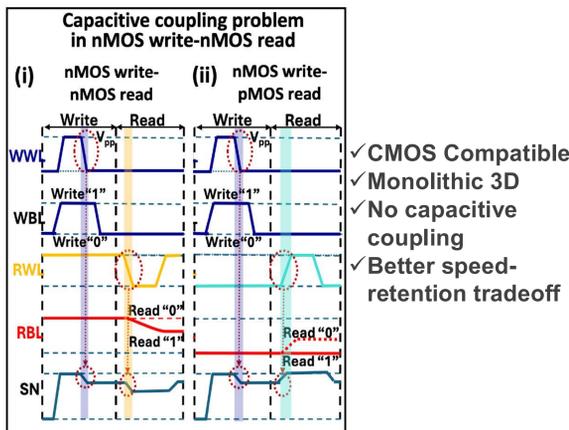
Objectives

- Complementary GC (CGC), with its n- and p-type OS configuration, results in high density, good retention, CMOS compatibility, low capacitive coupling, and adequate read speeds
- CGC with complementary back-end-of-line technology enables future monolithic multiple layers 3D integration

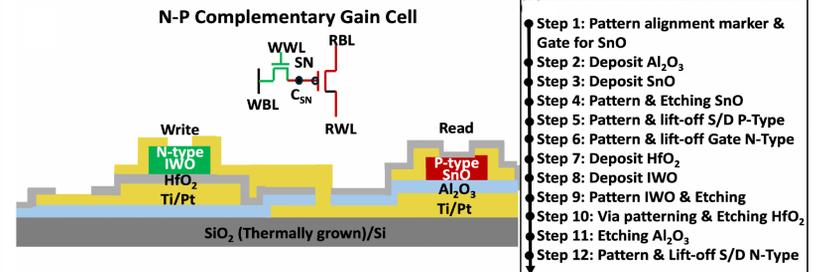


GC Comparison

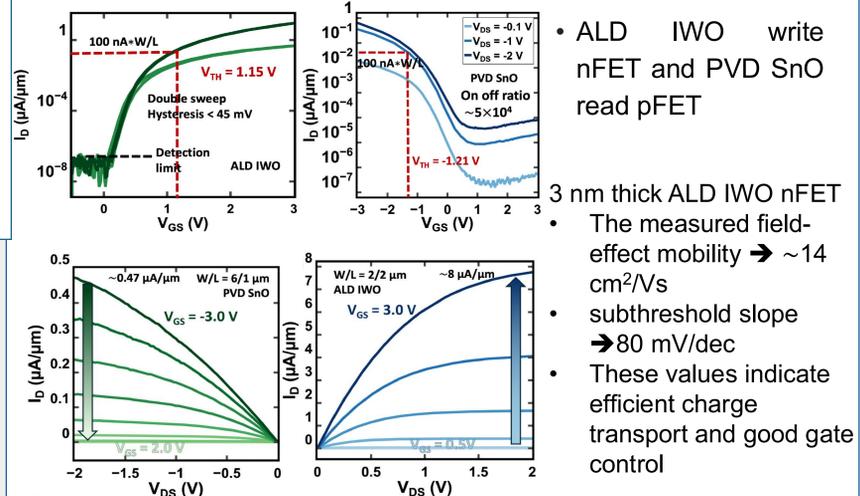
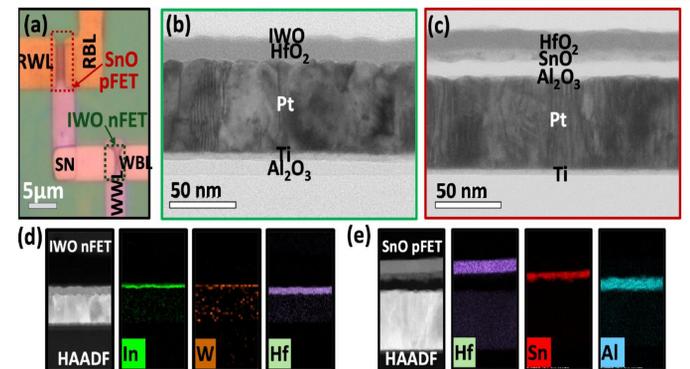
nOS-nOS	nOS-pSi (HGC)	nOS-pOS (CGC)
<ul style="list-style-type: none"> ✓ CMOS compatible ✓ Monolithic (N-layer) integration ✗ Issue with capacitive coupling ✗ Poor speed-retention tradeoff 	<ul style="list-style-type: none"> ✓ CMOS compatible ✗ FEOL Si FET for read → multiple memory layers not possible ✓ No capacitive coupling problem ✓ Better speed-retention tradeoff 	<ul style="list-style-type: none"> ✓ CMOS compatible ✓ Monolithic multiple-layer integration possible ✓ Good p-OSFET performance is needed ✓ No capacitive coupling problem ✓ Better speed-retention tradeoff



Technical Approach

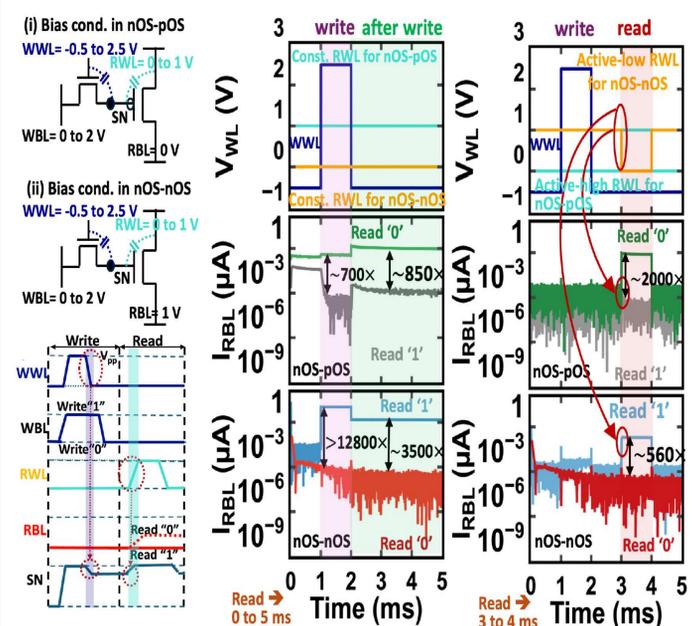


- Step 1: Pattern alignment marker & Gate for SnO
- Step 2: Deposit Al₂O₃
- Step 3: Deposit SnO
- Step 4: Pattern & Etching SnO
- Step 5: Pattern & lift-off S/D P-Type
- Step 6: Pattern & lift-off Gate N-Type
- Step 7: Deposit HfO₂
- Step 8: Deposit IWO
- Step 9: Pattern IWO & Etching
- Step 10: Via patterning & Etching HfO₂
- Step 11: Etching Al₂O₃
- Step 12: Pattern & Lift-off S/D N-Type



- ALD IWO write nFET and PVD SnO read pFET
- 3 nm thick ALD IWO nFET
- The measured field-effect mobility → ~14 cm²/Vs
- subthreshold slope → 80 mV/dec
- These values indicate efficient charge transport and good gate control

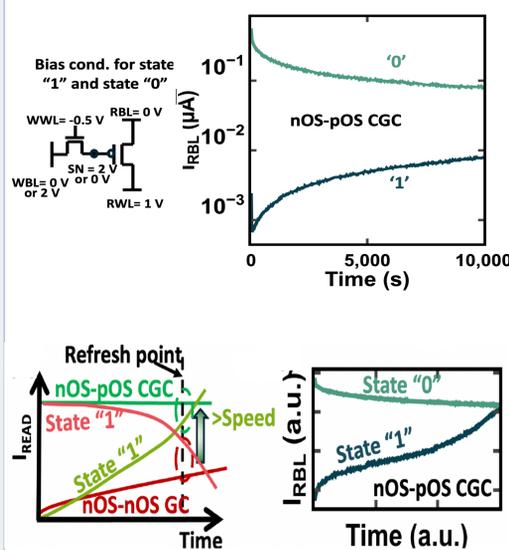
Counterbalanced Capacitive Coupling in CGC



RWL-SN capacitive coupling

- causes further degradation in nOS-nOS GC
- improves/retains the margin between state "0" and state "1" for nOS-pOS CGC (performing read only from 3 ms to 4 ms)

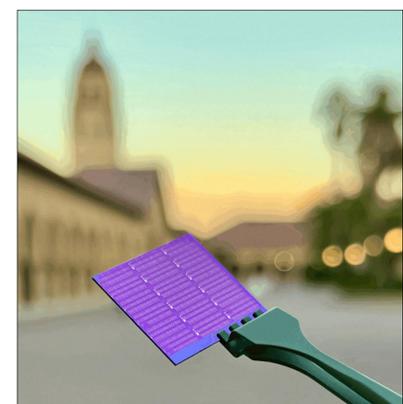
Optimized State Allocation



- p-type read transistor → reversed polarity
- state "1" is allocated to low read current while state "0" to high read current
- read current at low V_{GS} refresh point is high, providing high speed read
- data states "0" and "1" remain clearly distinguishable for over 10,000 seconds with a read current (~0.1 μA)

Key Accomplishments, Showstoppers & Next Steps

- Key results and metrics:**
 - Counteracting capacitive coupling
 - P-type read enhances voltage sensing speed and sense margin
 - Better speed-retention tradeoff
 - 10,000 s retention, with state "0" as a high-current state
- Grand challenge applications:**
 - M3D integrated high-density on-chip memory
- Next Steps:**
 - Improve SS for p-type OS
 - Variation control
 - Multi-layer integration and scalability to advanced nodes



This work was supported by

Stanford ENERGY Precourt Institute for Energy
Stanford ENERGY Postdoctoral Fellowship

DAM
Stanford Differentiated Access Memories Project



Stanford Non-Volatile Memory Technology Research Initiative (NMTRI)

Stanford SystemX Alliance

CHIMES
Center for Heterogeneous Integration of Micro Electronic Systems

Publications: F. F. Athena et al., "First Demonstration of an N-P Oxide Semiconductor Complementary Gain Cell Memory," IEEE International Electron Devices Meeting, 2024 (SRC ID: 397536)

DAM
Stanford Differentiated Access Memories Project

Retreat | July 1-2, 2025

Conserving Memory Bandwidth in SmartNICs with Flow-Based Memory Addressing

Agur Adams, Colin Drewes, David Shim, and Ben O’Keefe
Advised by Phil Levis and Zakir Durumeric

Problem

SmartNICs cannot easily perform complex network analysis at line rate due to limited memory bandwidth.

Network line rates have significantly increased:

- Increased speeds has led to **reduced visibility**
- High traffic volumes **overwhelms** most analysis tools

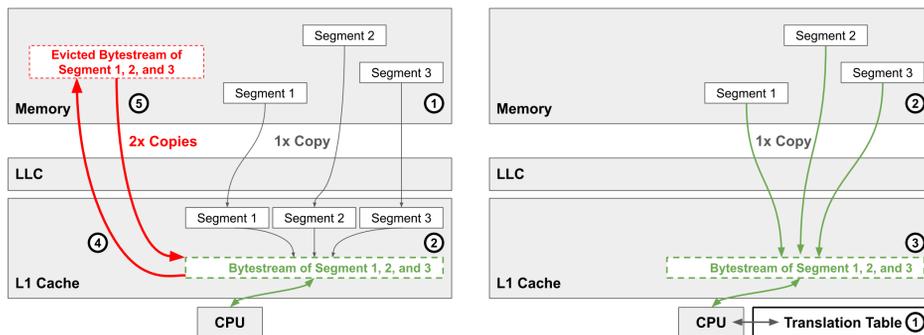
Complex analysis at line rates ≥ 100 Gbps requires on-NIC processing:

- Switches and routers support only simple analysis
- Specialized hardware can be costly or hard to program

SmartNICs have limited memory bandwidth relative to their line rates:

- Copying data to reassemble byte streams is expensive
- Careful cache management required

	Cores	NIC	DRAM	DRAM BW per core	NIC BW per core	Ratio
Google Cloud C3 Sapphire Rapids	176	200 Gbps	2x 8-ch DDR5	3.49 GB/s	0.14 GB/s	24.93
BlueField-3 SmartNIC DDR5	16	400 Gbps	2-ch DDR5	5.60 GB/s	3.13 GB/s	1.79



Goals

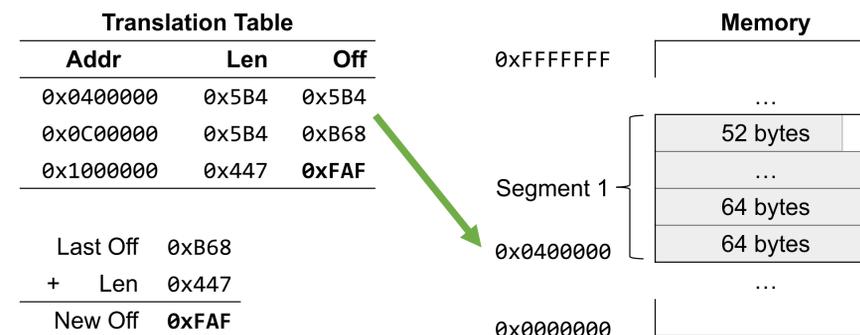
We aim to create a SmartNIC packet processing pipeline capable of arbitrary network analysis on reassembled TCP bytestreams at 200 Gbps using flow-based addressing.

Solution

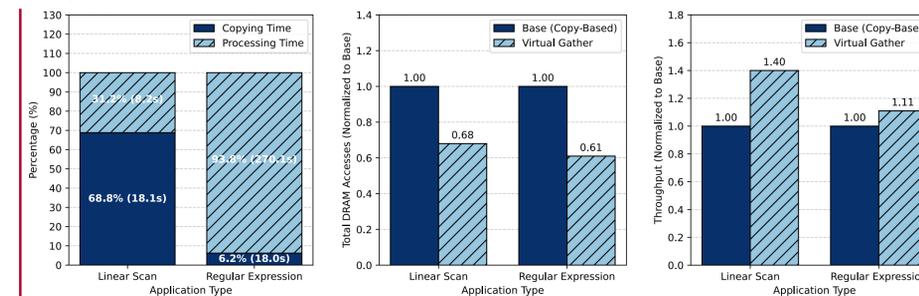
Using **flow-based addressing**, a SmartNIC can make the payloads of a series of packets addressable as a contiguous region of memory without a copy.

The SmartNIC populates and manages a translation table of all transport-layer segments of a flow.

Example: 4015 bytes received in three TCP segments (assuming a maximum segment size of 1460 bytes)



Impact



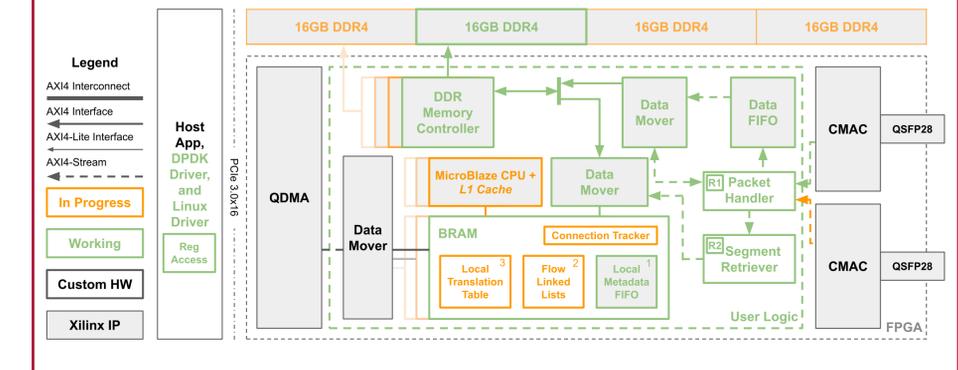
Preliminary Results: To quantify the cost of copying overhead on packet processing and the benefits of our optimization, we emulated the translation table in software and measured the time spent on copying versus application processing in two network applications: *linear scan* and *regular expression* matching.

- DRAM access reduction: 32.2% (linear) / 39.1% (regex)
- Throughput increase: 40% (linear) / 11% (regex)

Design

We are building out the proposed flow-based addressing scheme on an AMD Alveo U250 FPGA-based SmartNIC. Key aspects of the design include:

- Large number of RISC-V CPUs in parallel for flexible network analysis
- TCP connection tracking for flow-level awareness and timeouts
- Linked-list assembly for on-demand memory allocation



Challenges

How do you insert and manage the table at line rate (e.g., fragmentation, garbage collection, etc.)?

- We observed approximately 22,000 new flows a second (~145 Gbps) on the Stanford Campus network.
- We estimate ~32 ns between packets at 200 Gbps.

How large do you make the table?

- ~80% of all flows on the Stanford Campus network are less than 1100 bytes long, which means a translation table with just 18 entries could support most flows.

Should you use virtual or physical addresses in the table?

- We use physical addresses in our FPGA implementation but are considering both options.

How does the translation table interact with the cache?

- Could use the translation table to prefetch for the SmartNIC CPU’s cache