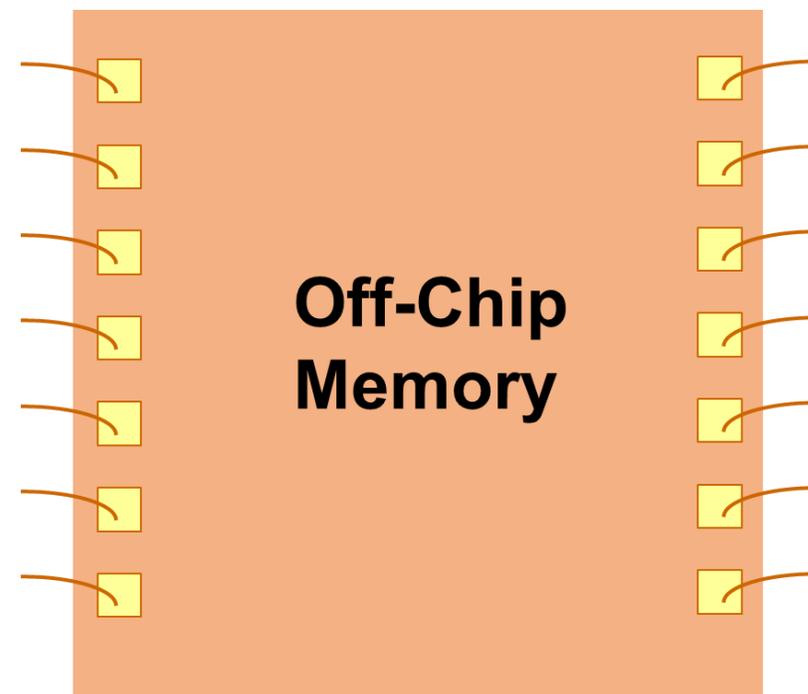
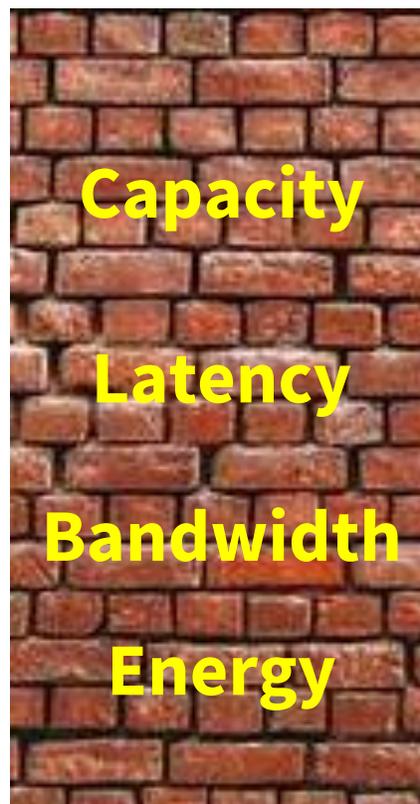
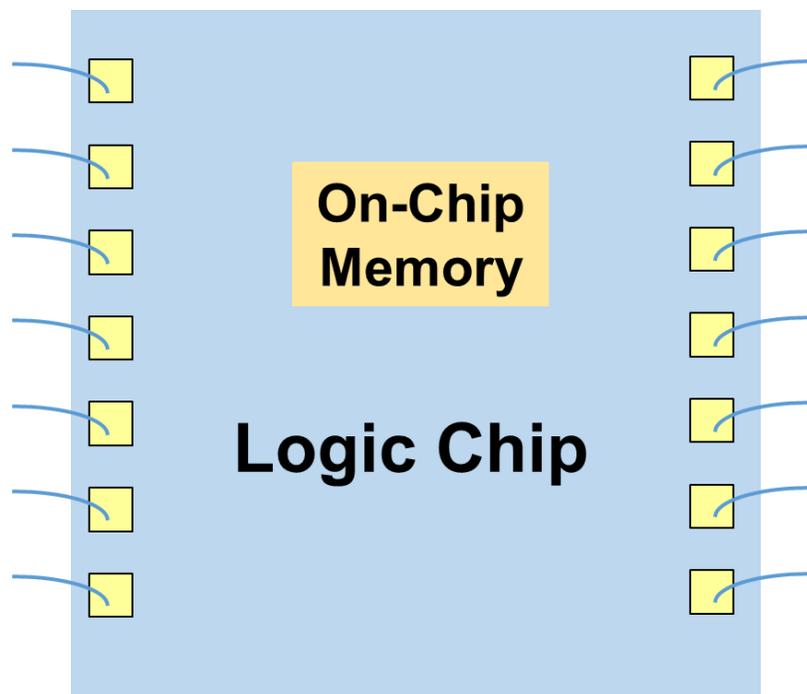


OpenGCRAM: An Open-Source Gain Cell Compiler Enabling Design-Space Exploration for AI Workloads

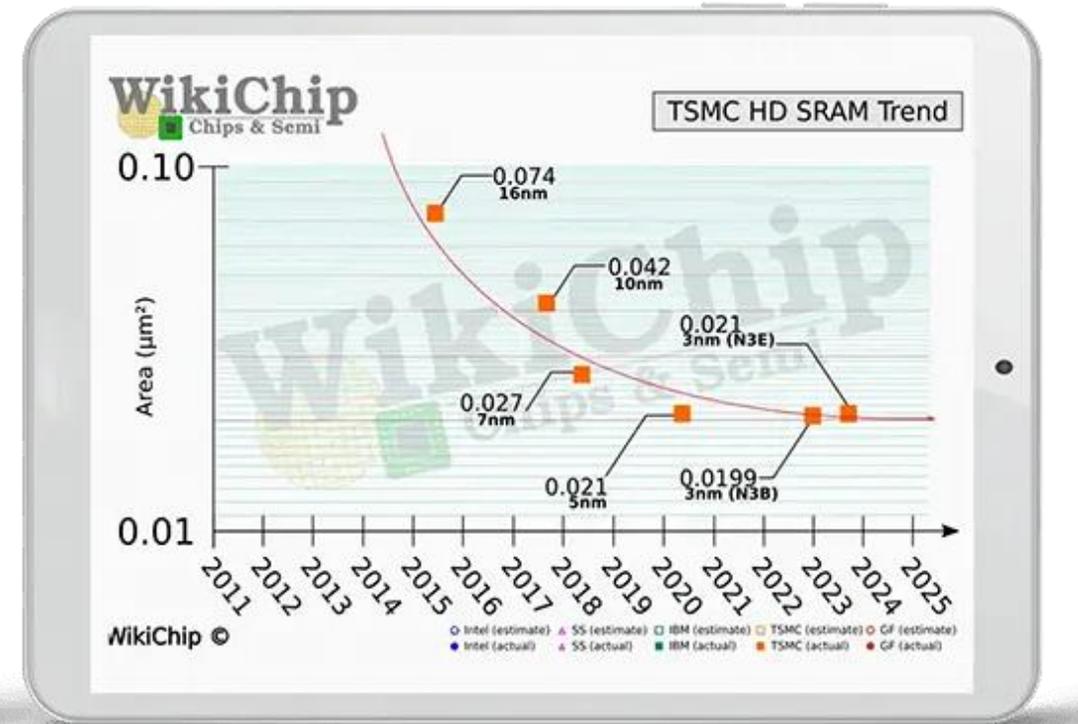
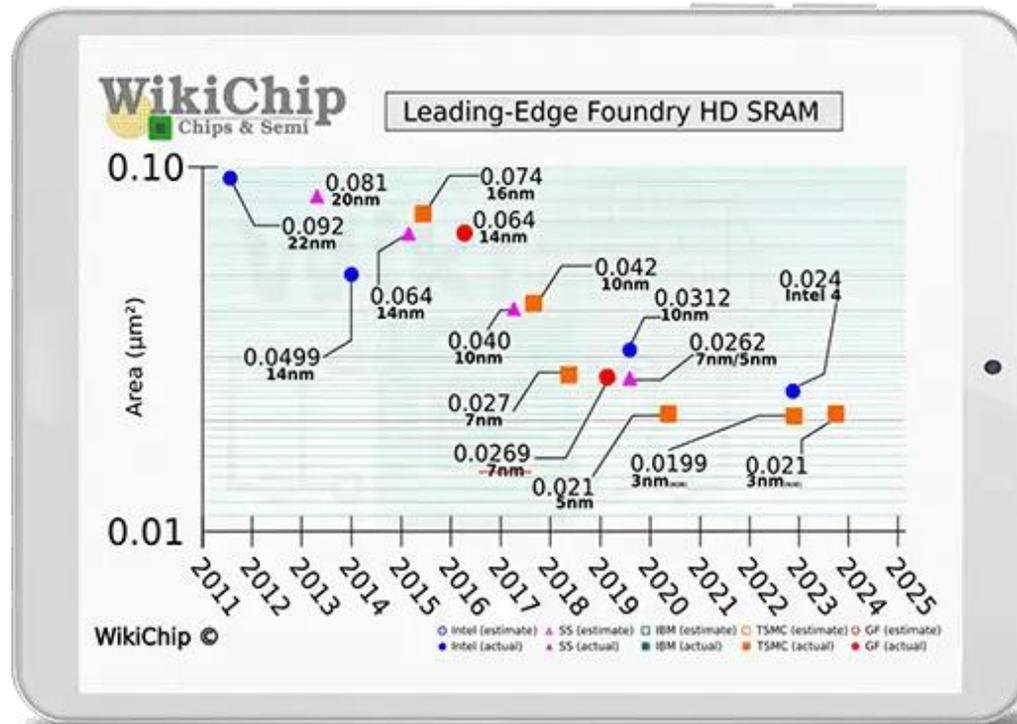
Xinxin Wang

04/28/2025

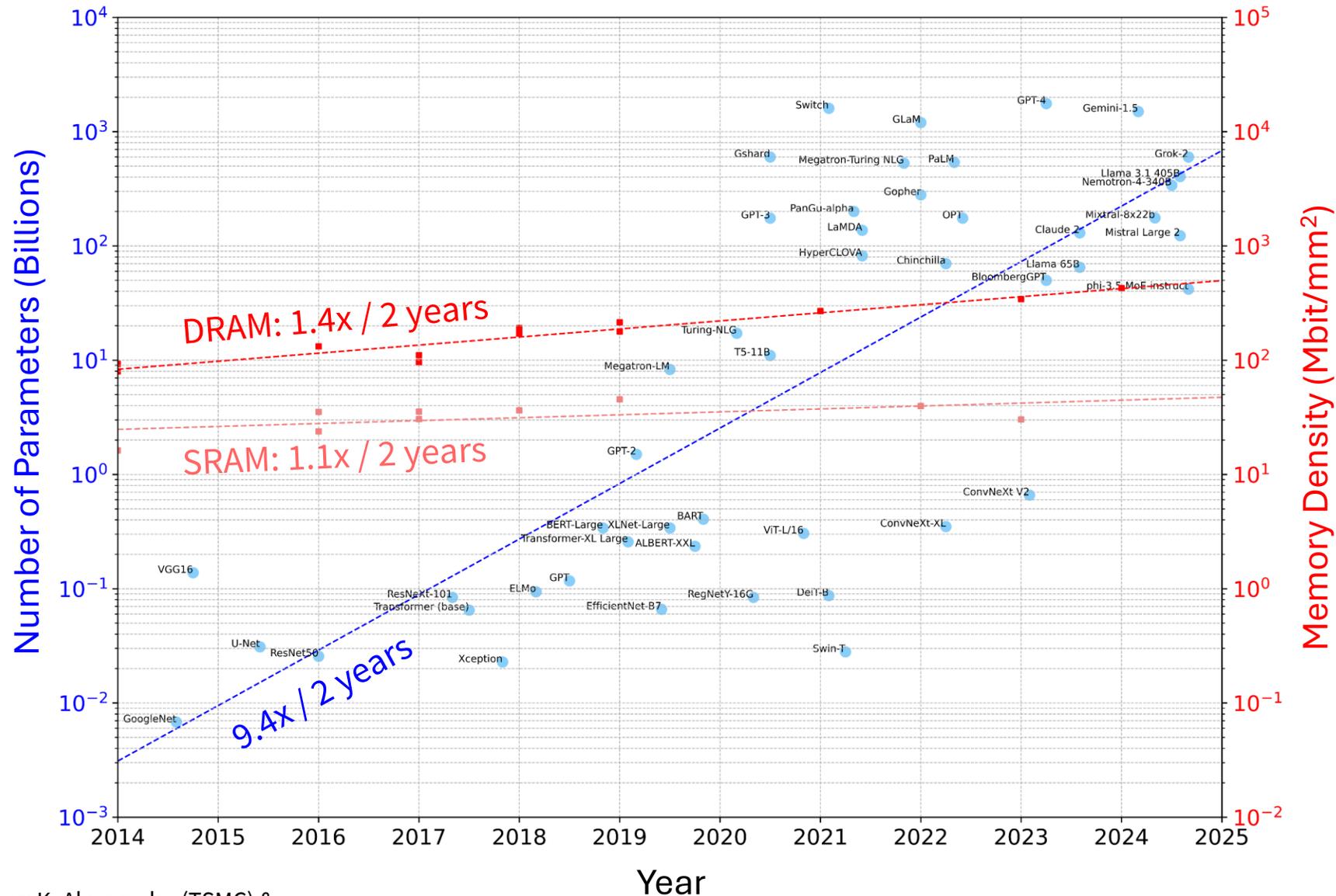
Memory Wall



SRAM scaling miniaturization

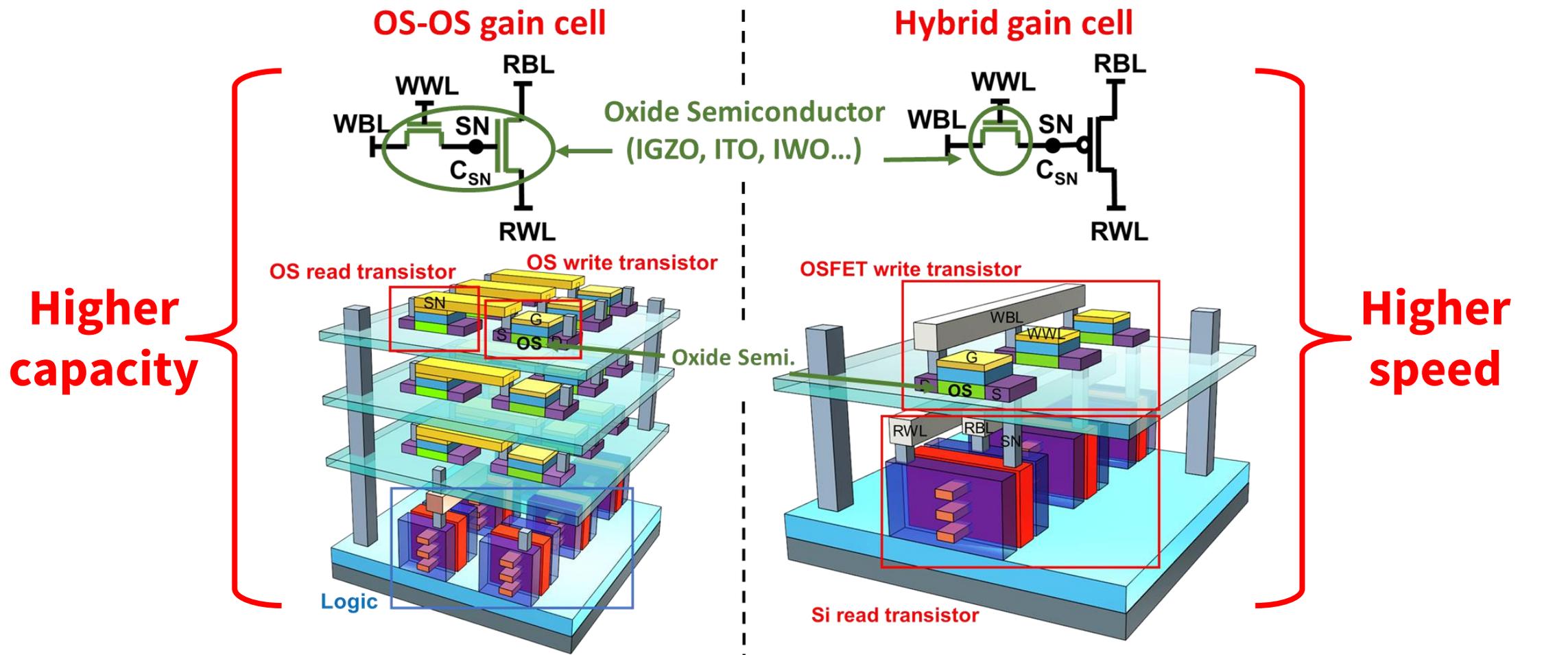


Memory Needs Outpace Memory Advances

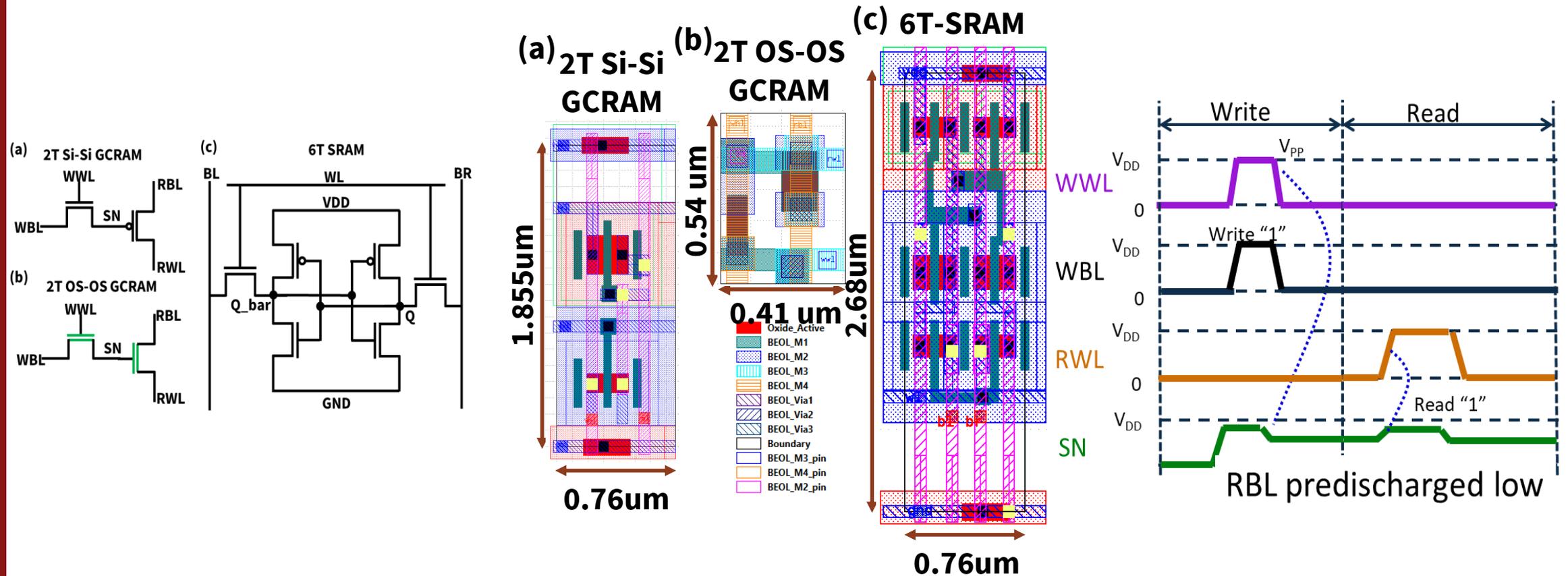


**We need
innovations for
high-density
on-chip memory**

Gain Cell memory: higher density than SRAM



Gain Cell memory: higher density than SRAM



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

We need a Gain Cell compiler

Related work

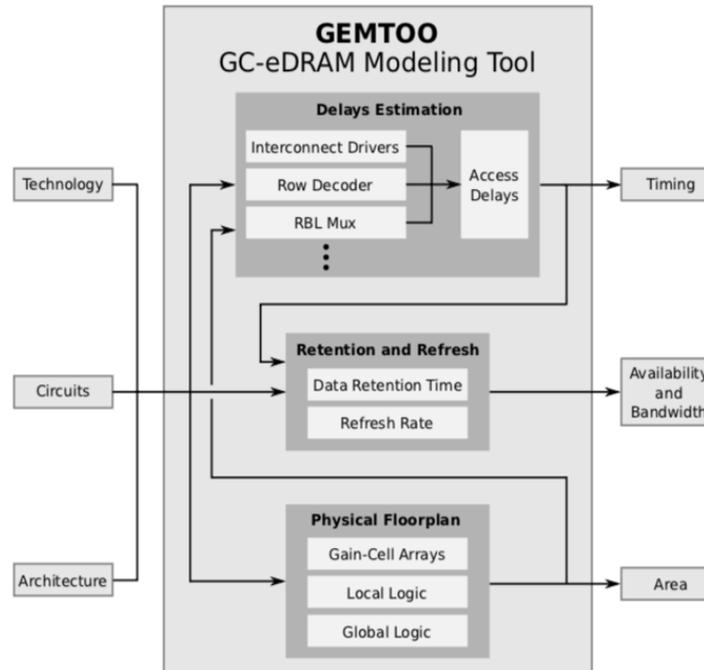
• RAAAM GCRAM

- Up-to **2X Higher density** vs. SRAM
- Up-to **50% area reduction** vs. SRAM
- Standard **SRAM interface**
- **Extended interface options** vs. SRAM
- **Single/two** ported
- Up to **2Mbit** instances
- **Standard CMOS** process
- **Single cycle** operation
- **Customizable**

Limitations:

- 3T GCRAM
- commercialized
- not open-access
- **Compiler in development**

• GEMTOO Simulator



Limitations:

- No netlist and layout
- No power evaluation
- Very rough delay estimation

• OpenRAM Compiler

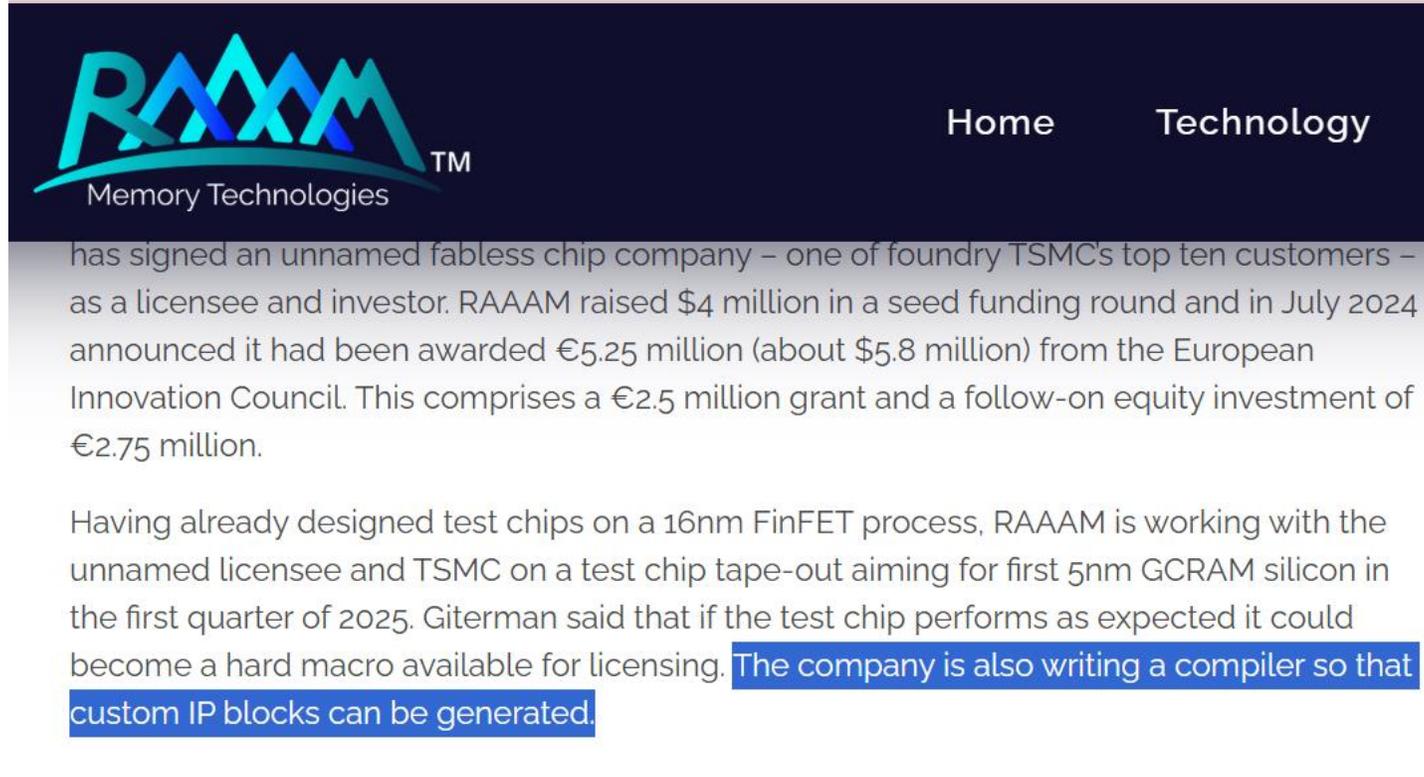
- Open-source SRAM compiler
- Support open-source PDKs
- SRAM bank netlist & layout generation
- Performance simulation

Limitations:

- No Gain Cell support
- No commercial PDK support

Related work

- RAAAM GCRAM

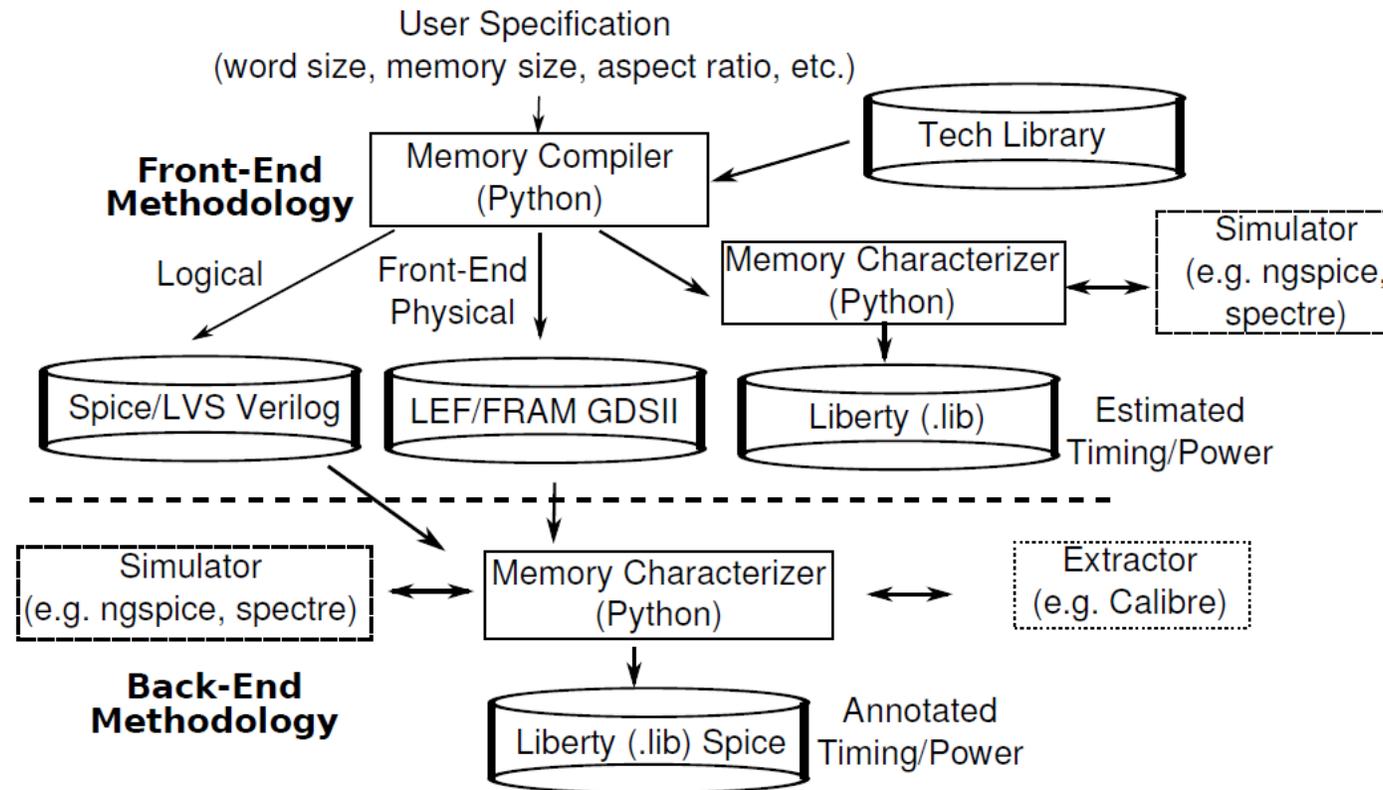


The image shows a screenshot of the RAAAM Memory Technologies website. The header features the RAAAM logo (stylized blue and green letters) with 'Memory Technologies' and a 'TM' symbol below it. To the right of the logo are navigation links for 'Home' and 'Technology'. Below the header, the main text describes RAAAM's partnership with a fabless chip company, its funding, and its work on GCRAM silicon. A blue highlight is present under the sentence: 'The company is also writing a compiler so that custom IP blocks can be generated.'

Peter Clarke is a veteran reporter and analyst covering the global electronics industry.
This article was first published by The Ojo-Yoshida Report on September 17, 2024.
See www.ojoyoshidareport.com for more of such reports.

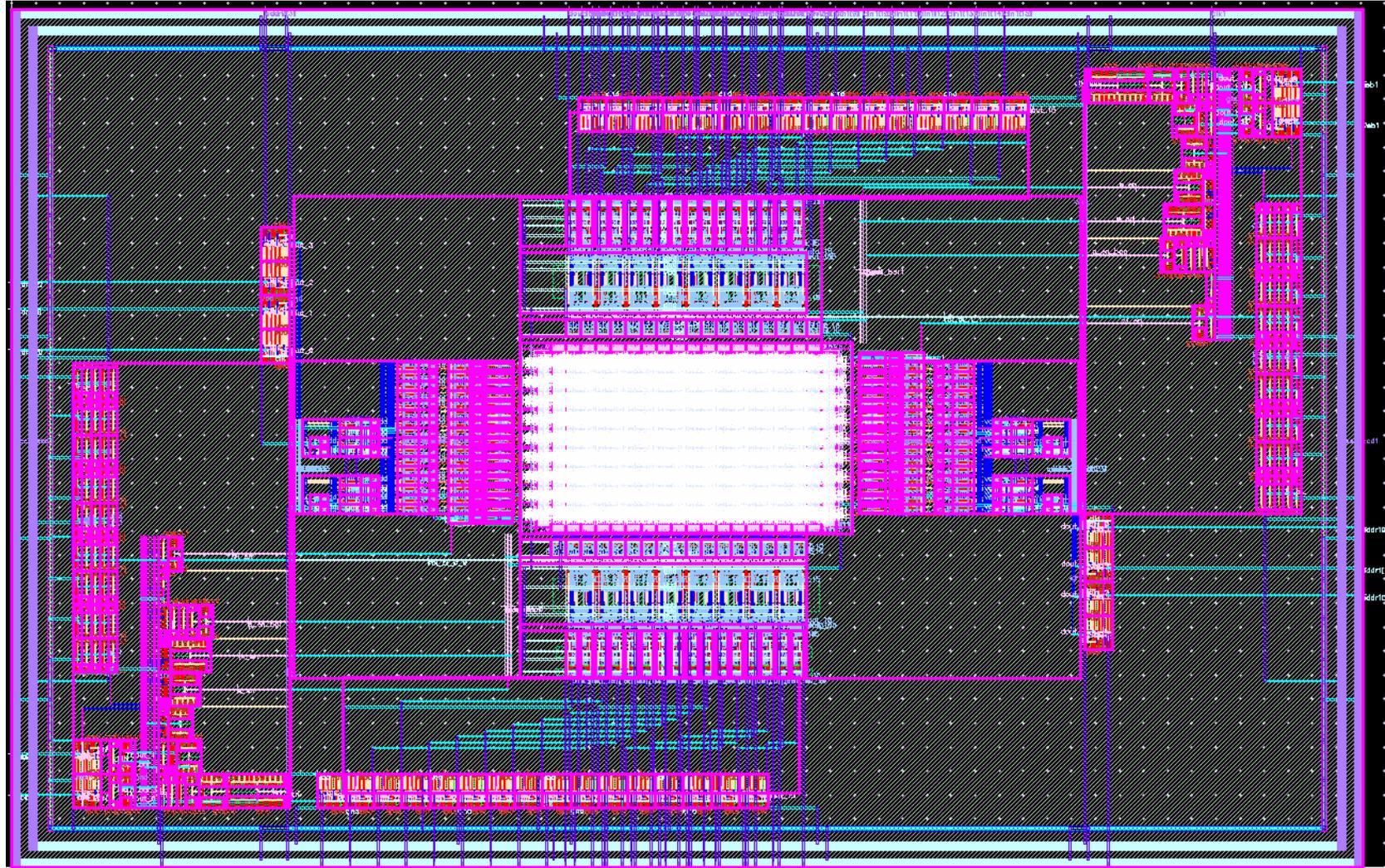
OpenRAM SRAM compiler

- Front-end: Create Netlist and Layout, DRC & LVS check;
- Back-end: Perform simulations on these generated files



OpenRAM SRAM compiler

- Example: 16x16 SRAM macro (file:///E:/OneDrive%20-%20Stanford/CNT-ITO/OpenRAM/sram_16x2.html)

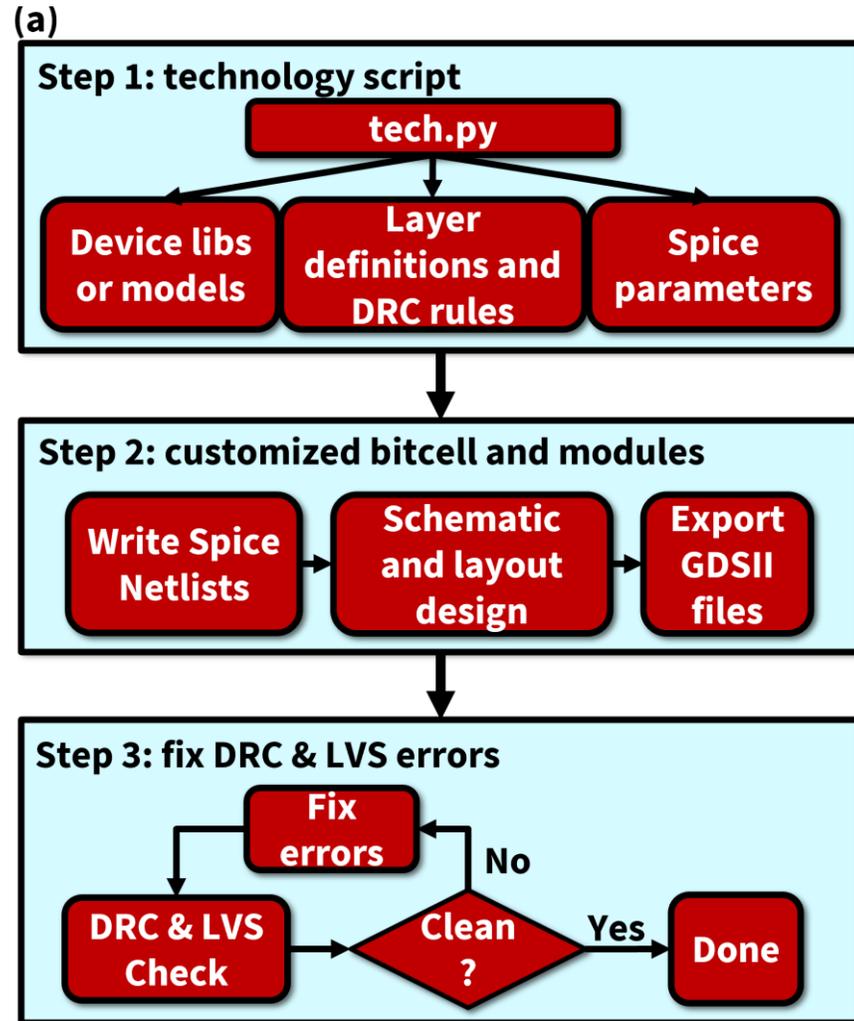


OpenRAM SRAM compiler

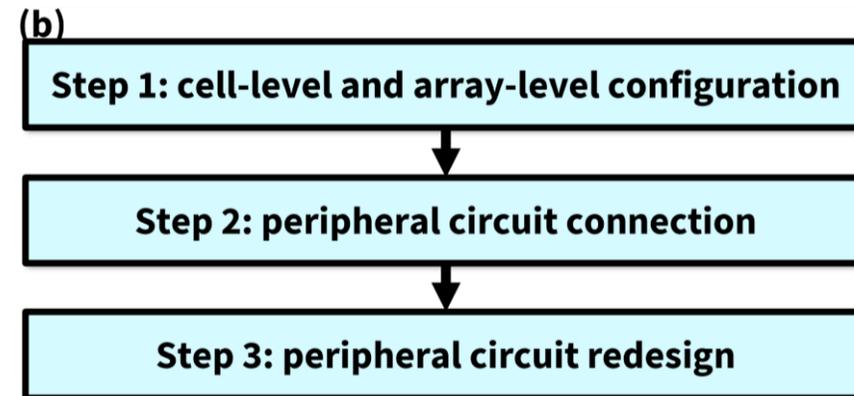
- Limitations:
 - Support NCSU FreePDK 45nm, MOSIS 0.35um (SCN4M_SUBM), Skywater 130nm (sky130), no advanced tech nodes
 - Only SRAM is supported
 - Only one macro architecture is supported

Methodology to extend OpenRAM functionality

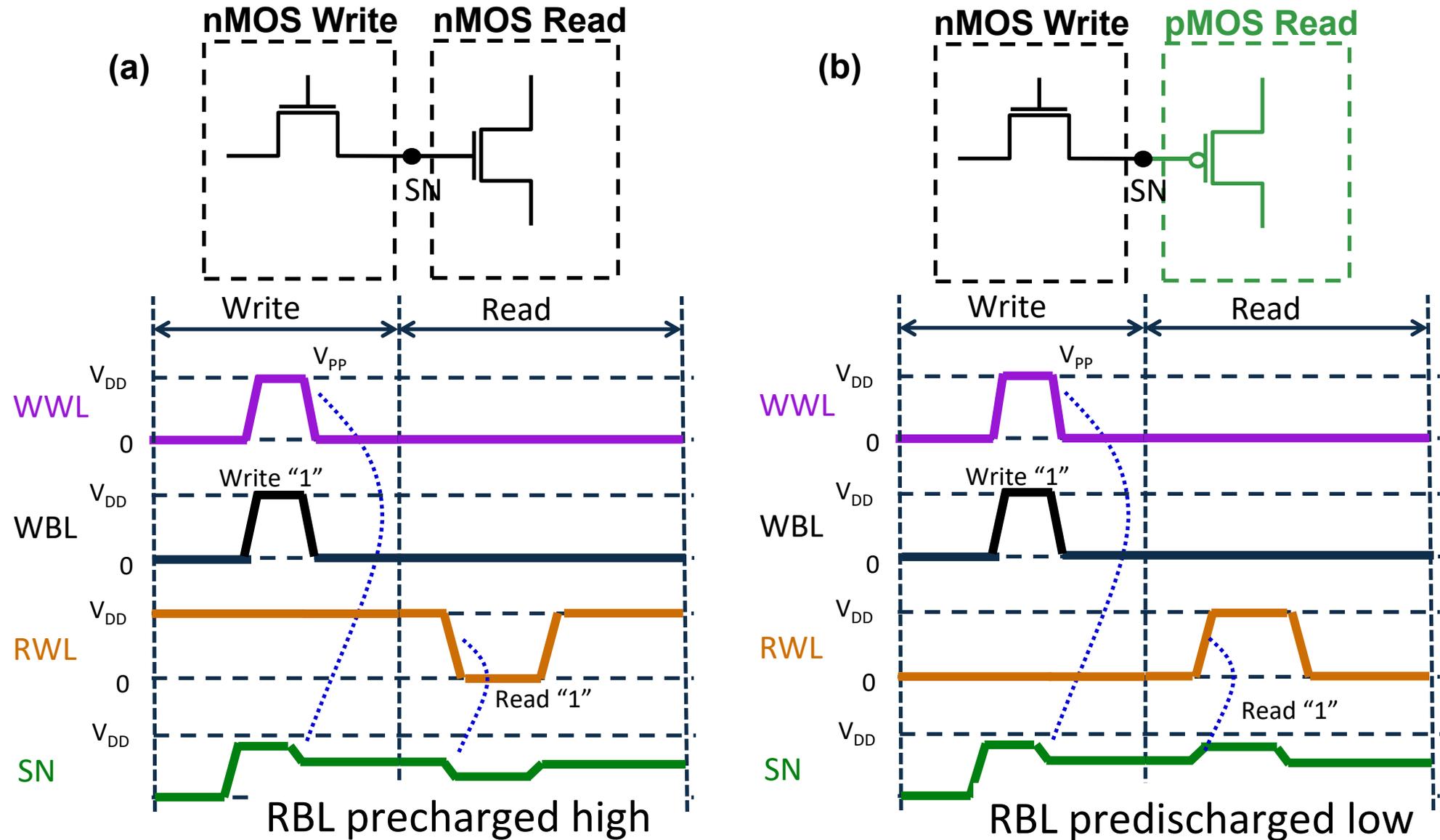
- Step 1: Porting to new PDKs



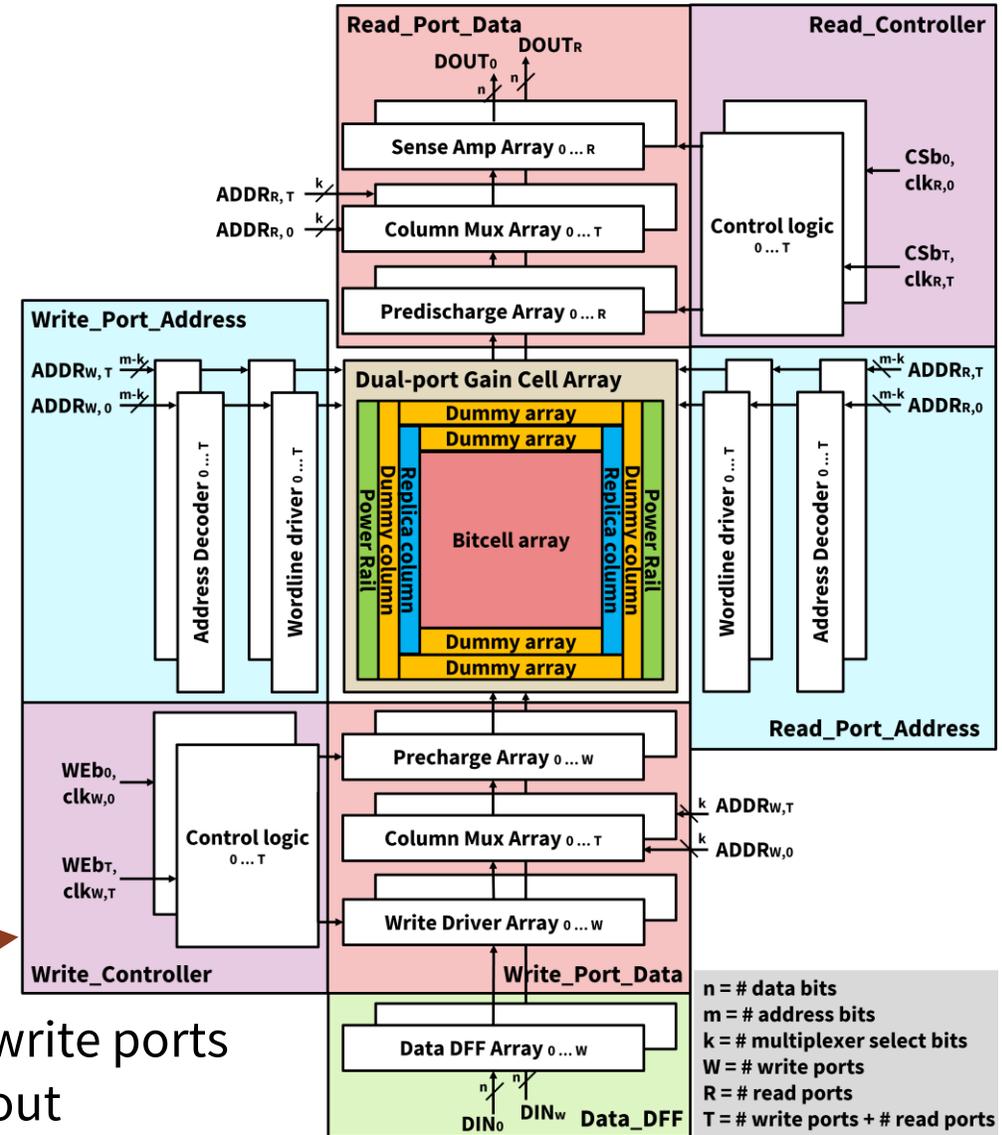
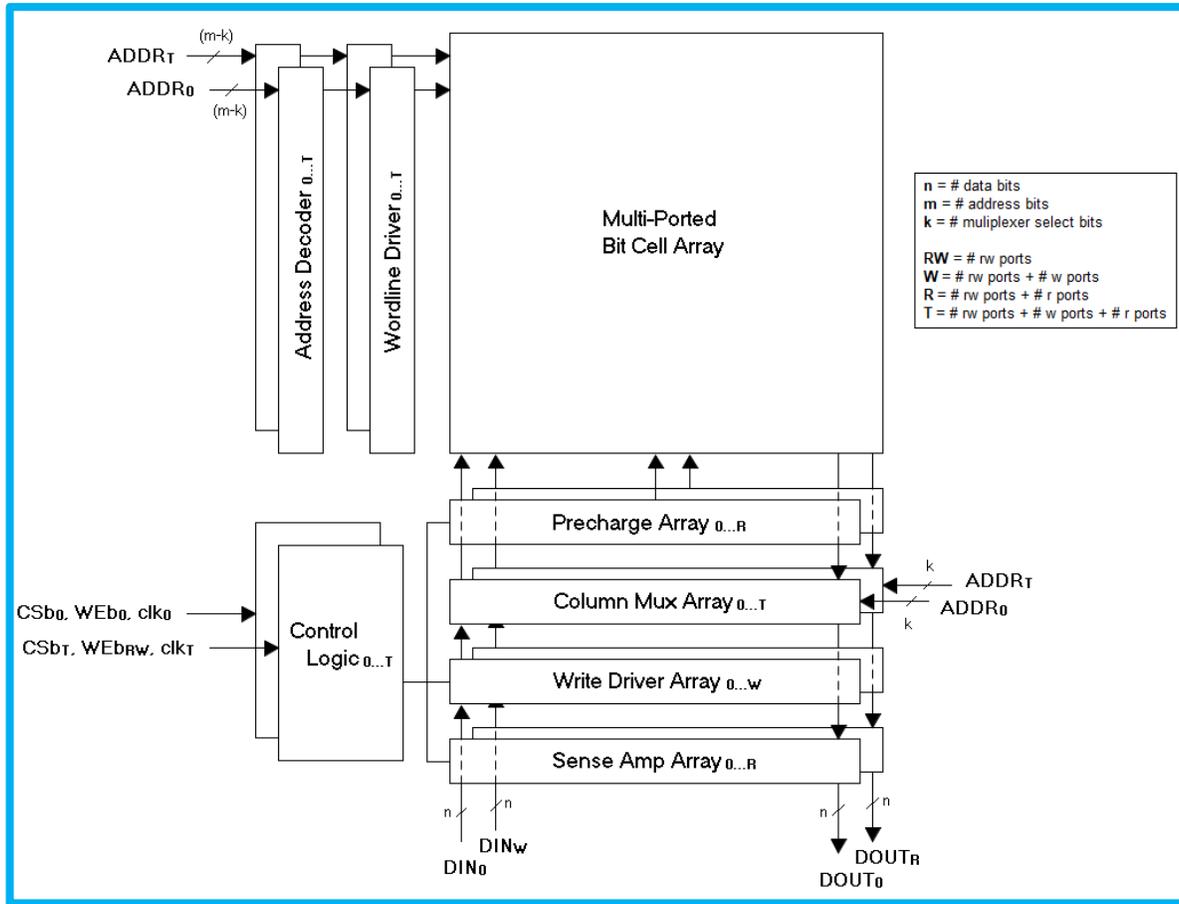
- Step 2: Adding new memory technologies



Gain Cell memory operations



OpenGCRAM: bank architecture



SRAM bank

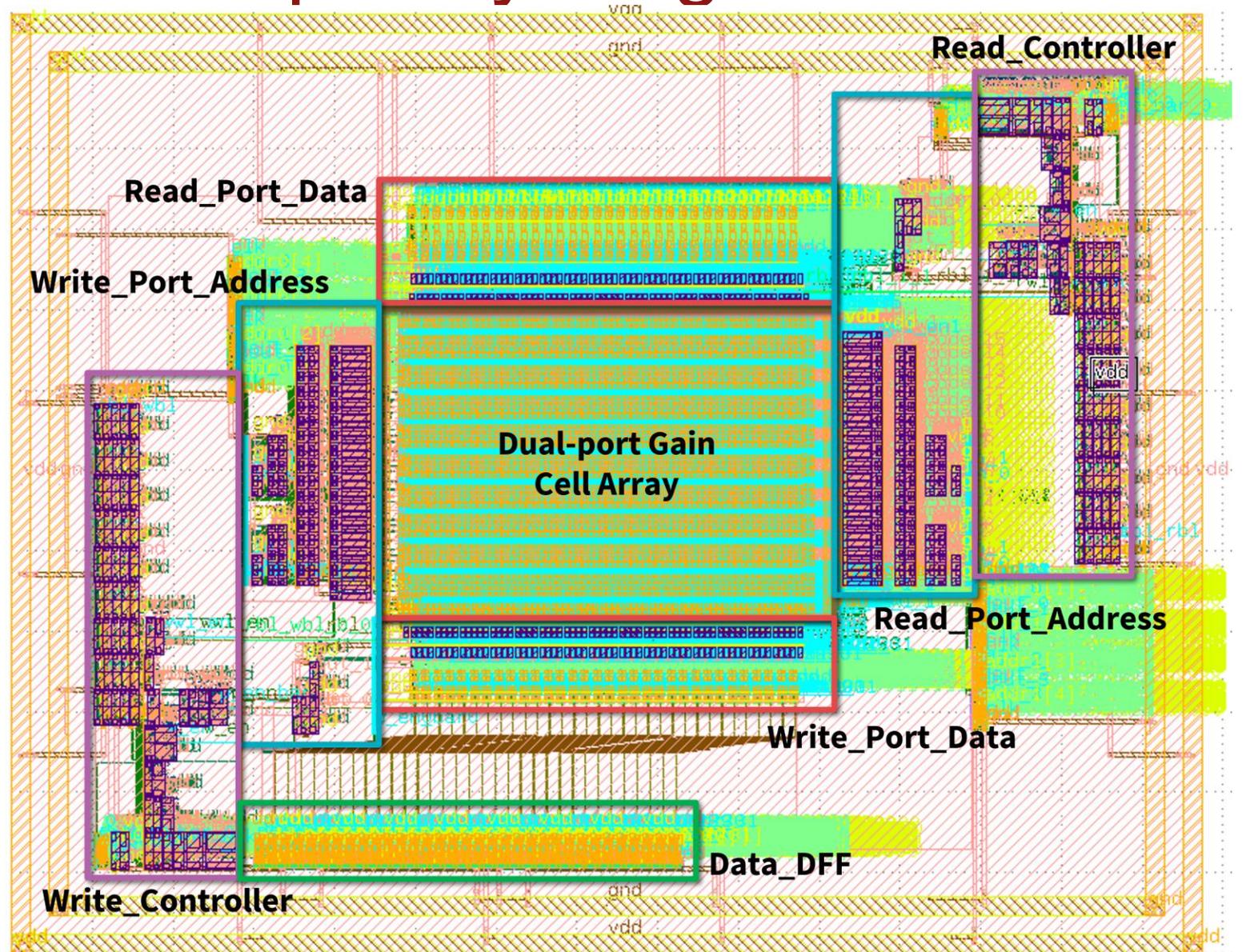


Gain Cell bank

- separate read and write ports
- single-ended read out
- predischarge array for read port

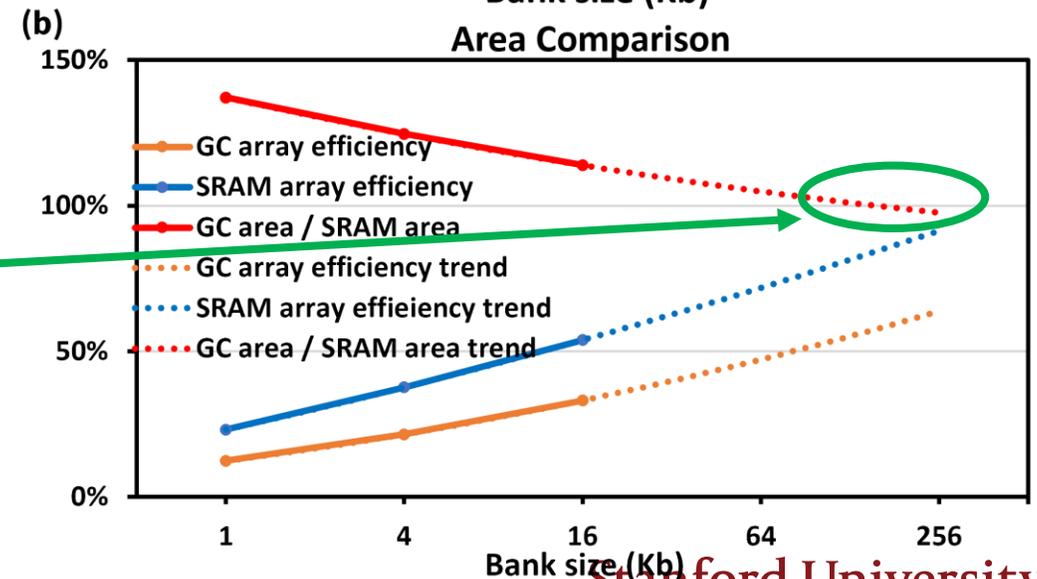
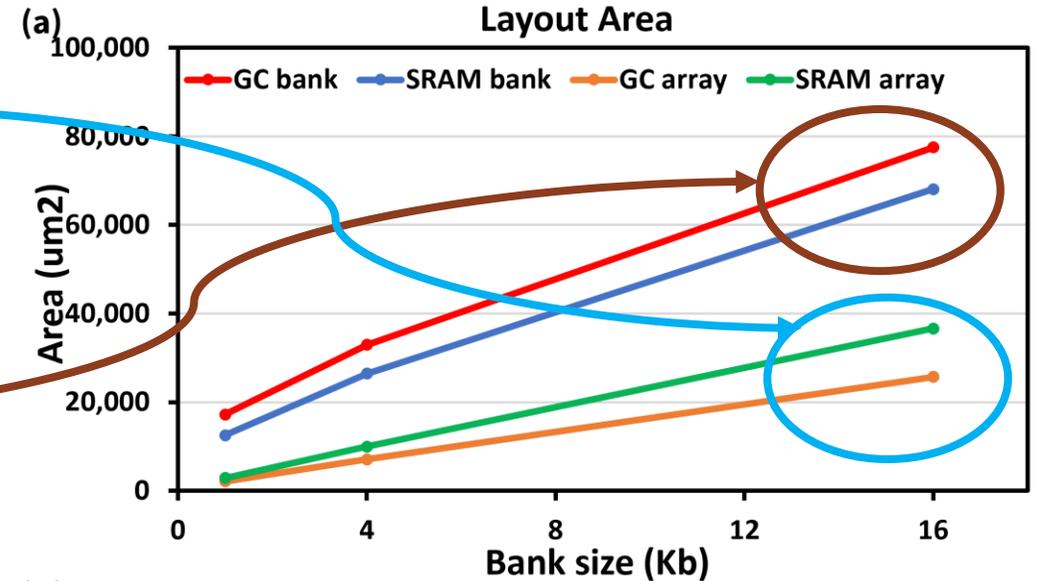
OpenGCRAM: Example layout generation

A 32x32 Gain Cell bank generated by OpenGC



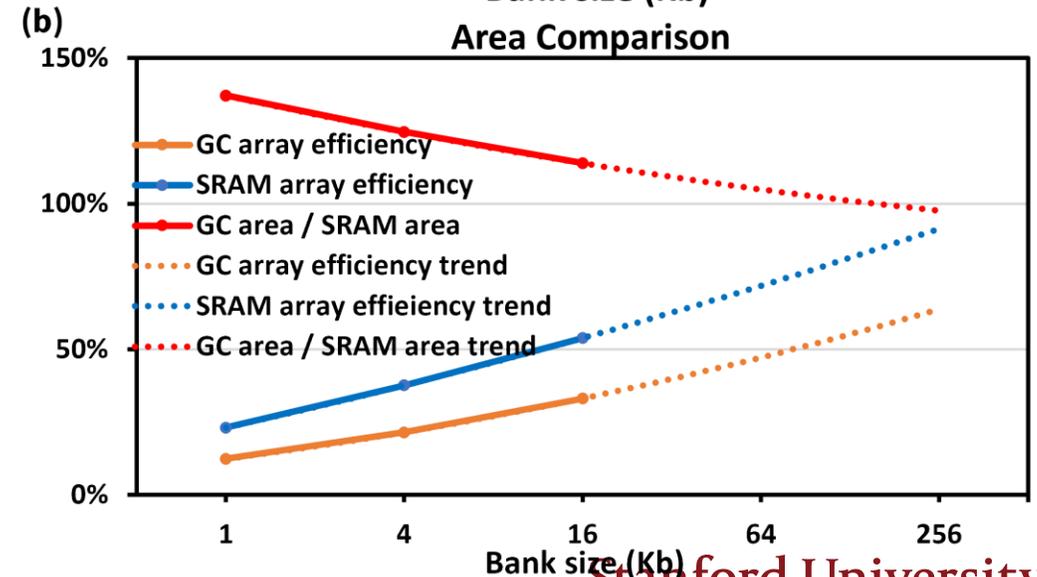
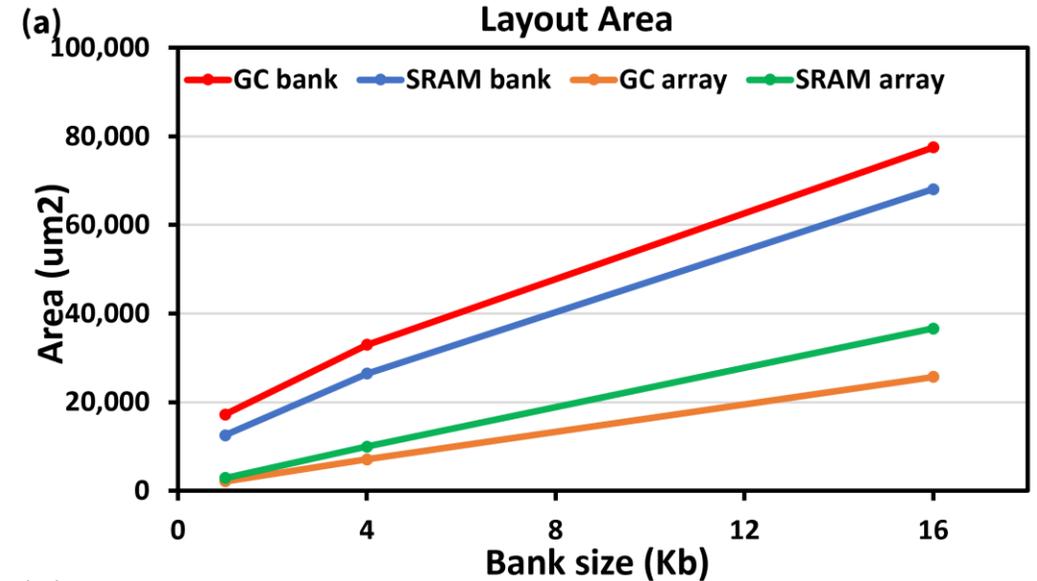
Gain Cell vs. SRAM: area

- GC array area < SRAM array area
→ GC cell size < SRAM cell size
- GC bank area > SRAM bank area for small bank size
→ GC has separate peripheral circuits for different ports
- GC bank area < SRAM bank area for large bank size
→ Peripheral circuits are amortized



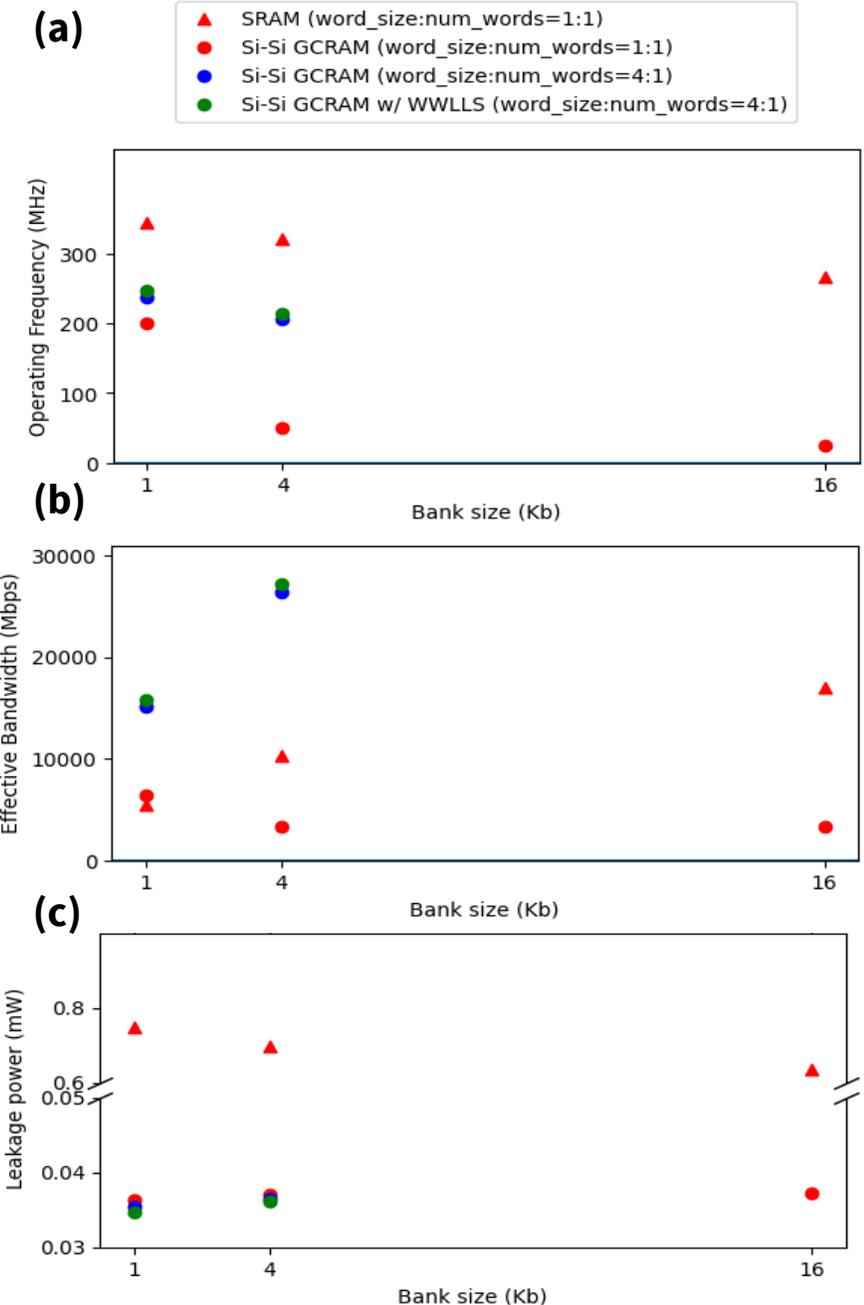
Gain Cell vs. SRAM: area

- Dual-port peripheral of GC bank allows simultaneous read and write operations, resulting in high bandwidth
- Dual-port SRAM bank area is $\sim 2x$ of single-port SRAM bank area, which is larger than GC bank.



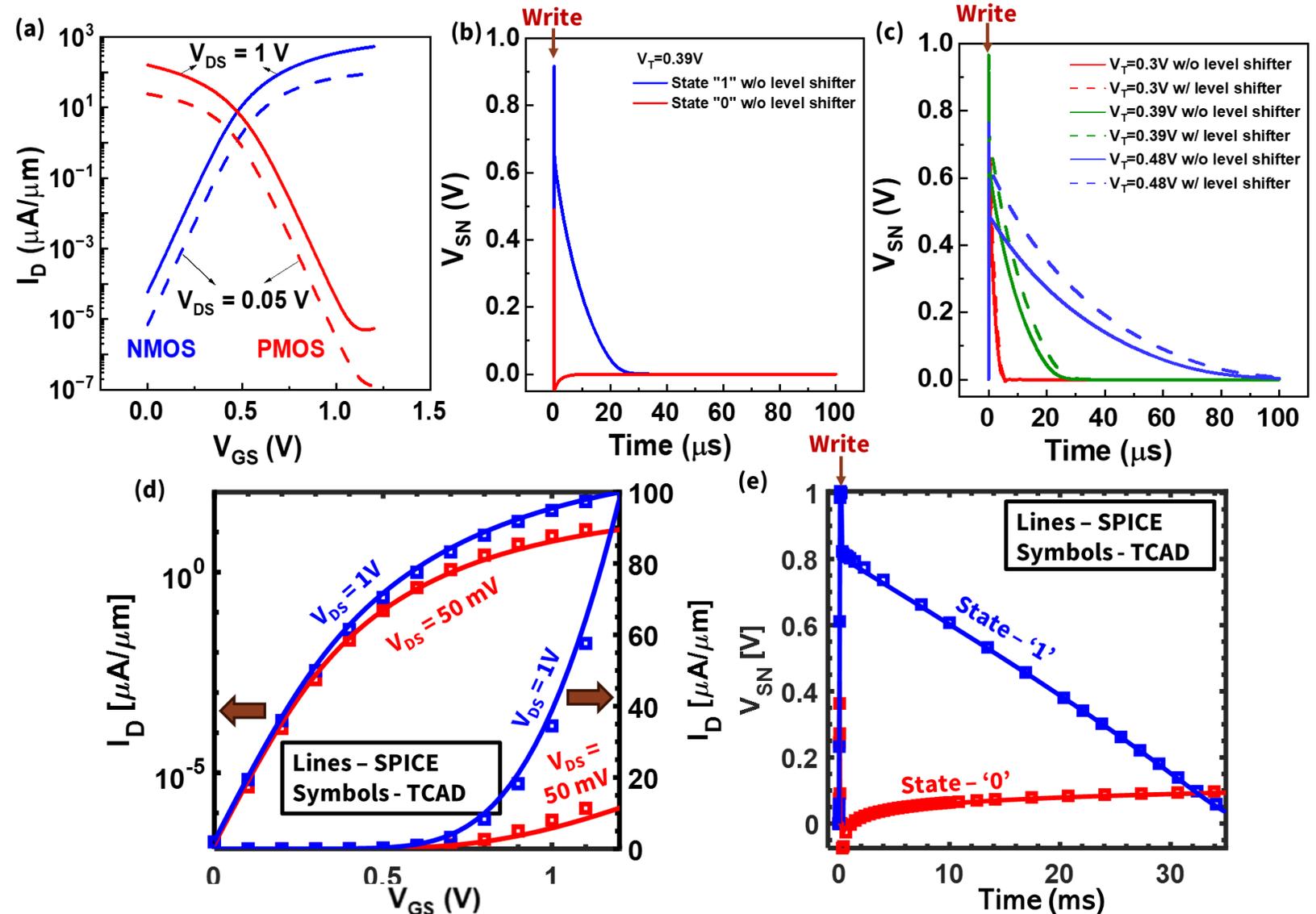
Gain Cell vs. SRAM: delay and power

- GC frequency < SRAM frequency
 - GC uses single-ended readout
 - GC storage node voltage is $V_{dd} - V_{th}$ when writing "1"
 - GC frequency can be boosted with design space exploration
- GC leakage power < SRAM leakage power
 - There is no direct path from Vdd to GND in GC
- An abrupt change of GC freq and power from 1 Kb to 4 Kb
 - Increase in delay chain stages



Gain Cell retention modulation

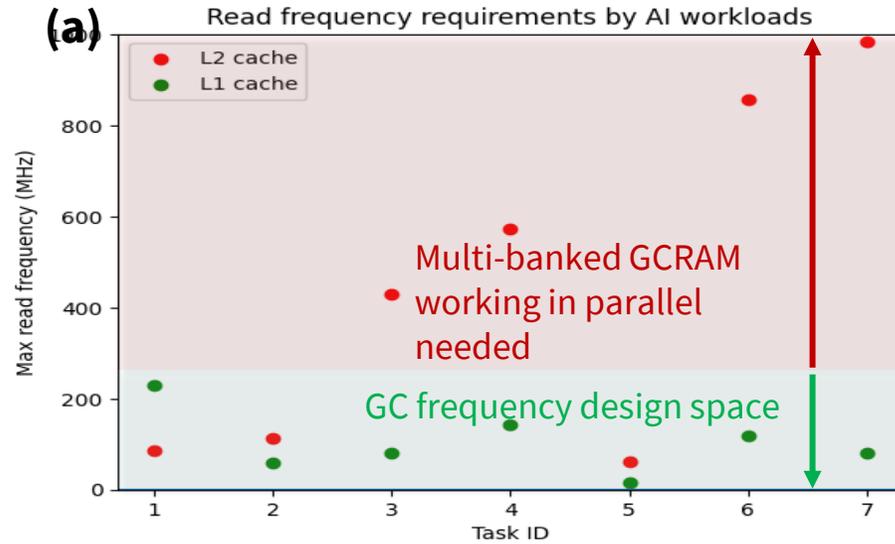
By adjusting the transistor design (like V_{th} and channel material), the retention can be tuned to accommodate for different applications (such as activation caches and weight memory in AI inference).



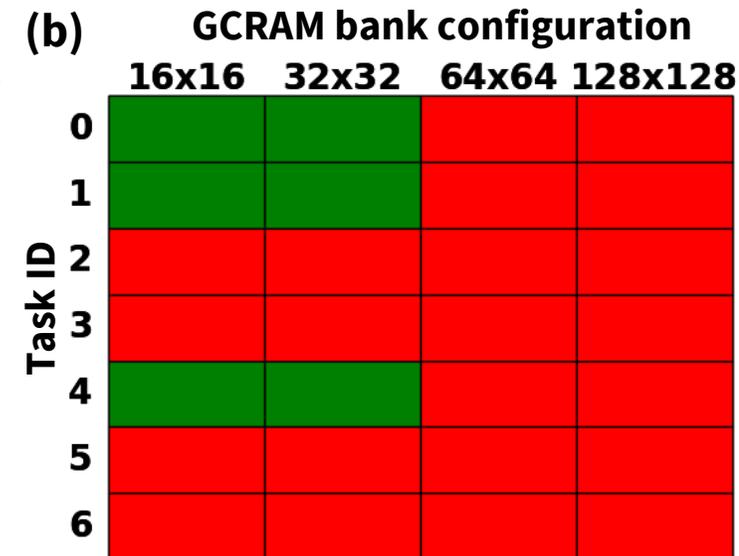
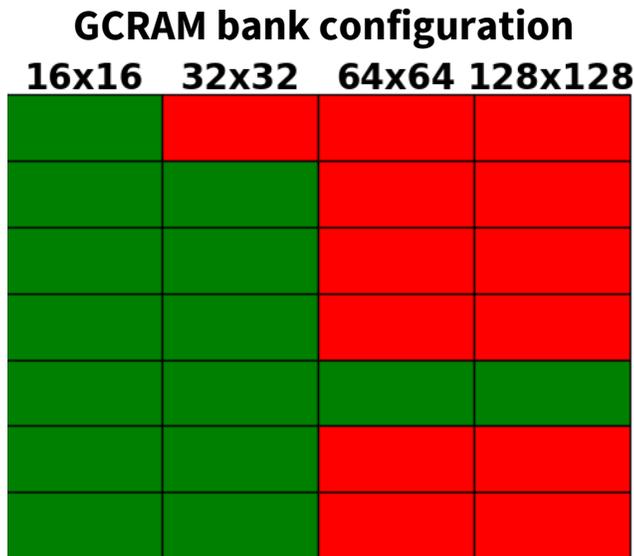
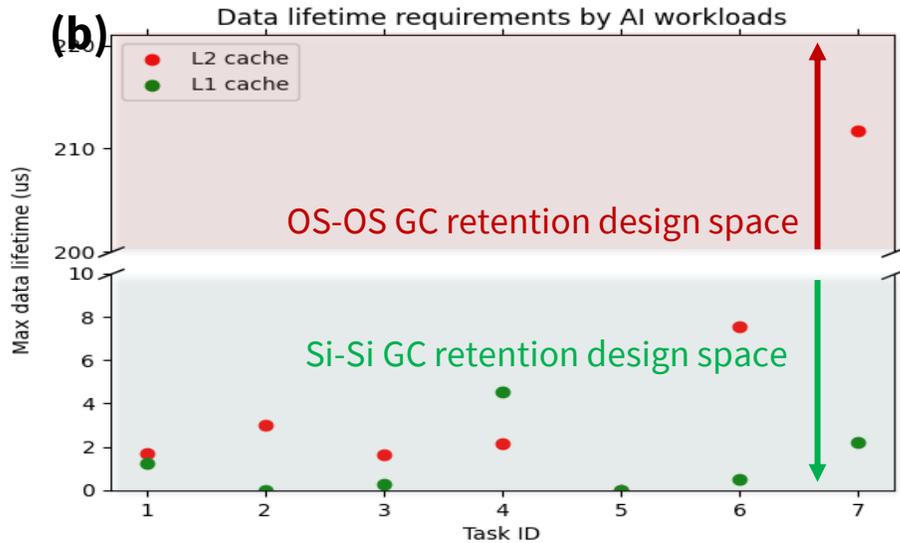
Design space exploration for AI workloads

Task	L1 lifetime Max (us)	L2 lifetime Max (us)	L1 read freq (MHz)	L2 read freq (MHz)
polybench- 2DConvolution	1.25	1.68	230.43	84.93
llama-3.2-11b-vision	4.50	2.14	142.82	574.17
bert-base-uncased	0.48	7.53	117.11	856.99
llama-3.2-1b	0.28	1.63	79.78	428.50
stable-diffusion	2.17	211.76	79.33	984.24
polybench- 3DConvolution	0	3.01	57.66	113.18
resnet-18	0	0	16.31	62.47

Design space exploration for AI workloads



Task ID	1	2	3	4	5	6	7	
Task name	2DConvolution	Polybench-vision	Llama-3.2-11b-uncased	Bert-base-uncased	Llama-3.2-1b	Stable-diffusion	Polybench-2DConvolution	Resnet-18



Summary

- We developed OpenGC to enable fast, accurate, customizable, and optimized Gain Cell bank design as high-density on-chip memory
- We introduced a standard methodology for porting OpenRAM compiler to new PDKs and memory technologies
- OpenGC supports Gain Cell bank design generation with TSMC N40 PDK and precise Spice simulations for performance evaluation
- By following the proposed methodology, OpenGC can be extended to support additional types of Gain Cell memory and other PDKs
- Design space exploration is enabled to accommodate for different AI workloads