

Monolithic 3D FPGA Design and Synthesis with Back-End-of-Line Configuration Memories

Faaiz Waqar¹, Jiahao Zhang², Anni Lu¹, Zifan He²,
Jason Cong², Shimeng Yu¹

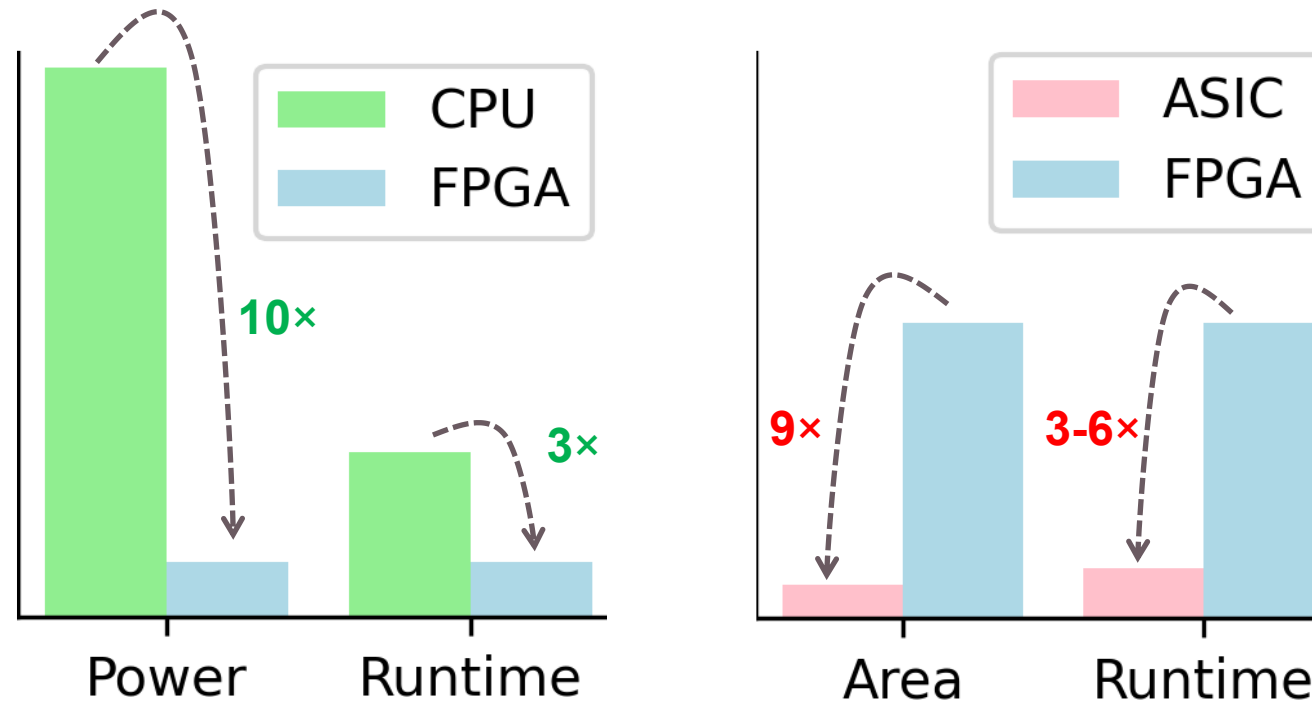
¹Georgia Institute of Technology; ²University of California, Los Angeles



SPONSORED BY

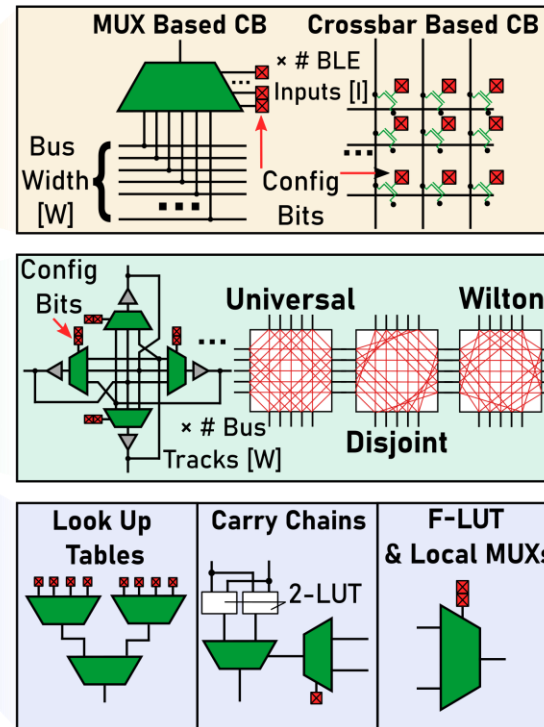
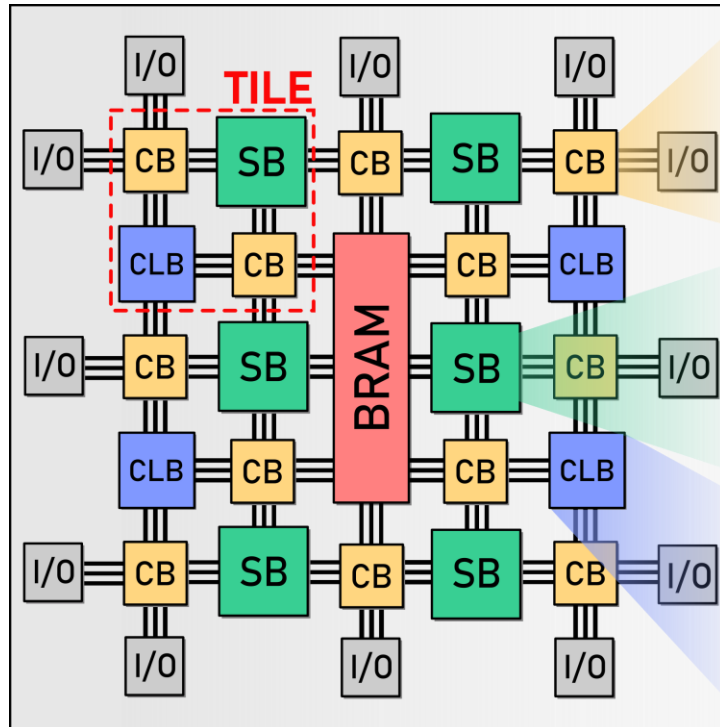
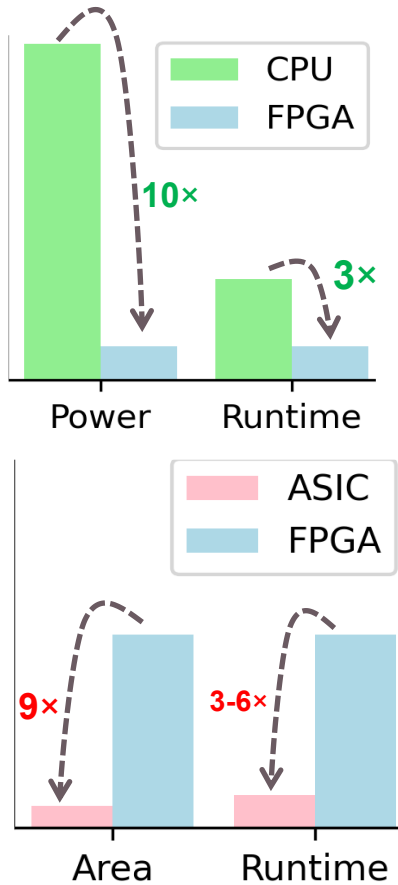


The Cost of Reconfigurability



- **FPGAs** can be desirably **reconfigured** for specialized execution
- Despite innovations, **FPGAs lag ASICs** in Power, Performance, and Area (PPA)

The Cost of Reconfigurability



CBs connect signals on bus to logic (CLB) inputs.

SBs steer signals at intersections to route inputs/outputs.

CLBs contain basic logic elements (BLEs) that emulate K input logical functions

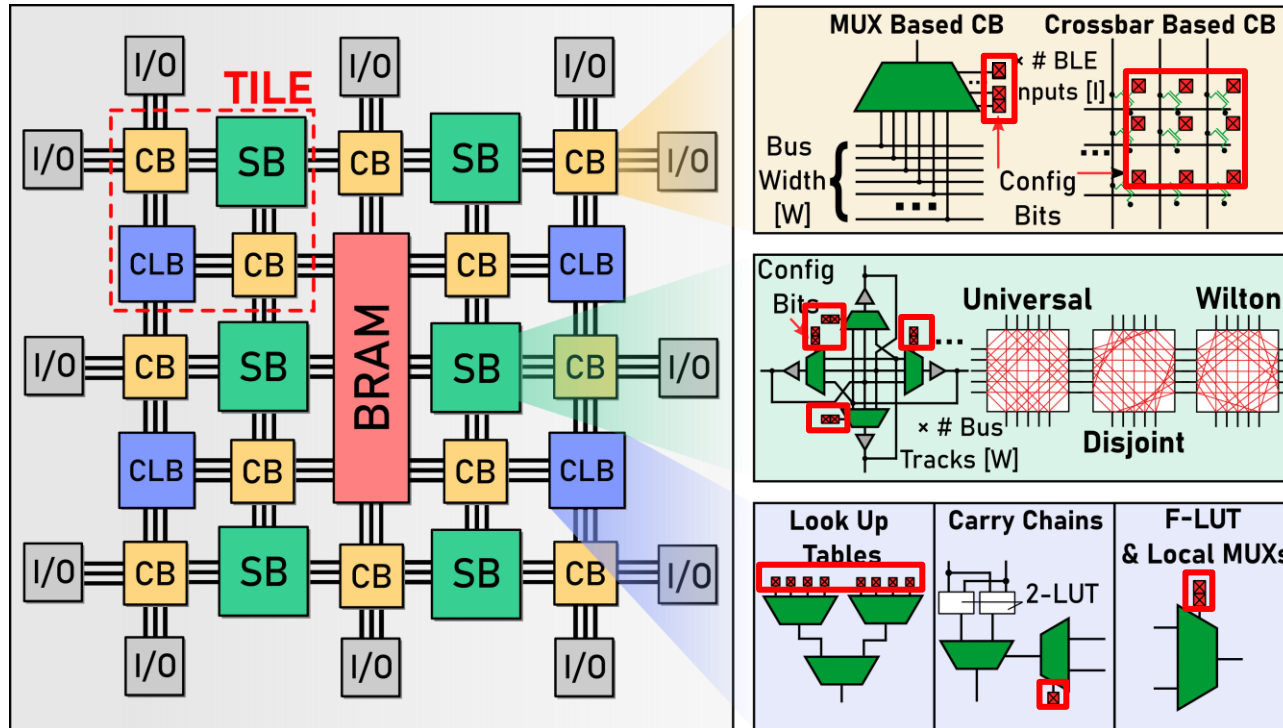
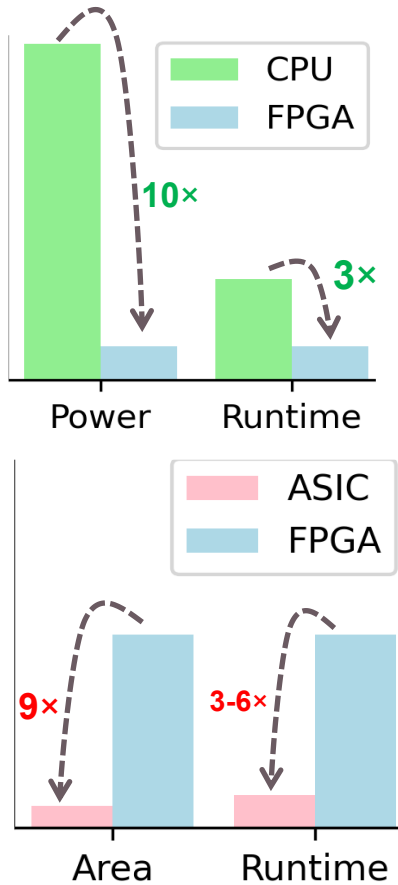
Island-based FPGAs are composed of Configurable **Logic Blocks (CLBs)**, **Switch Blocks (SBs)**, and **Connection Blocks (CBs)**.

A **tile (2 CBs + 1 SB + 1 CLB)** is a modular building block of an FPGA

✓ Reconfigurable

✗ Large PPA Disparity

The Cost of Reconfigurability



CBs connect signals on bus to logic (CLB) inputs.

SBs steer signals at intersections to route inputs/outputs.

CLBs contain basic logic elements (BLEs) that emulate K input logical functions

Island-based FPGAs are composed of Configurable Logic Blocks (CLBs), Switch Blocks (SBs), and Connection Blocks (CBs).

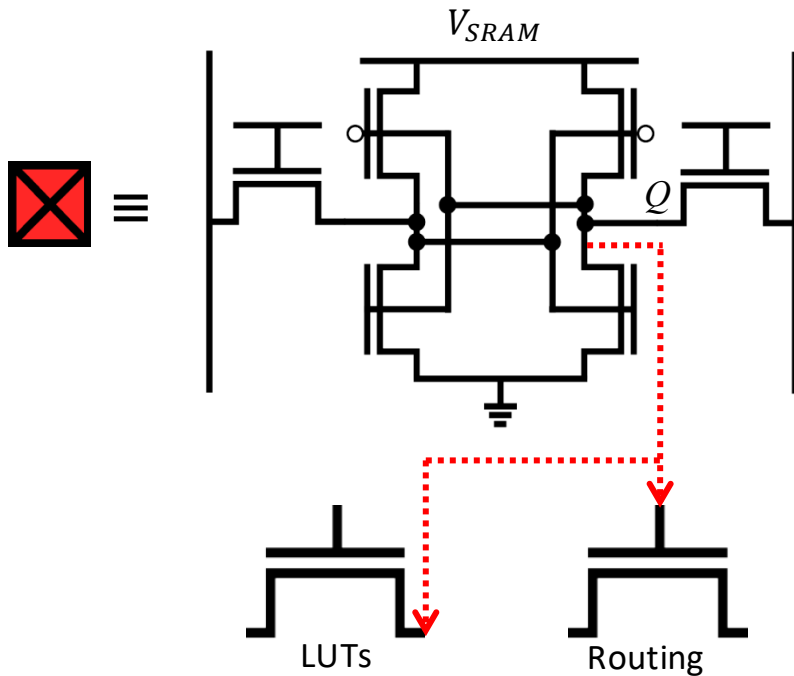
A tile (2 CBs + 1 SB + 1 CLB) is a modular building block of an FPGA

✓ Reconfigurable

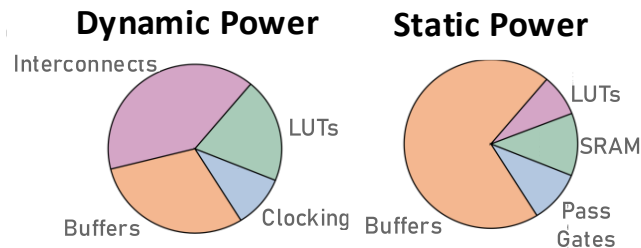
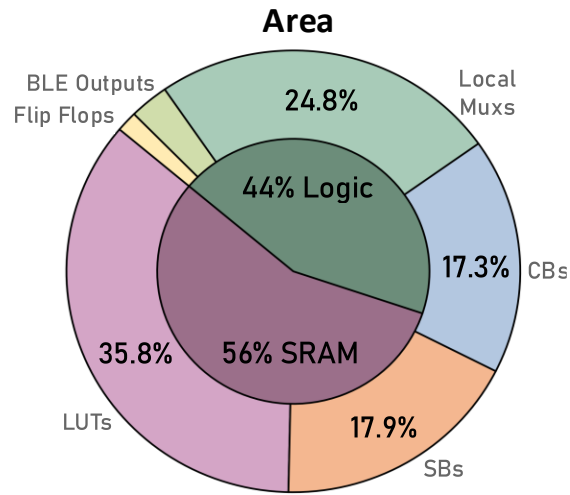
✗ Large PPA Disparity

The Cost of Reconfigurability

Configuration Bit (SRAM)

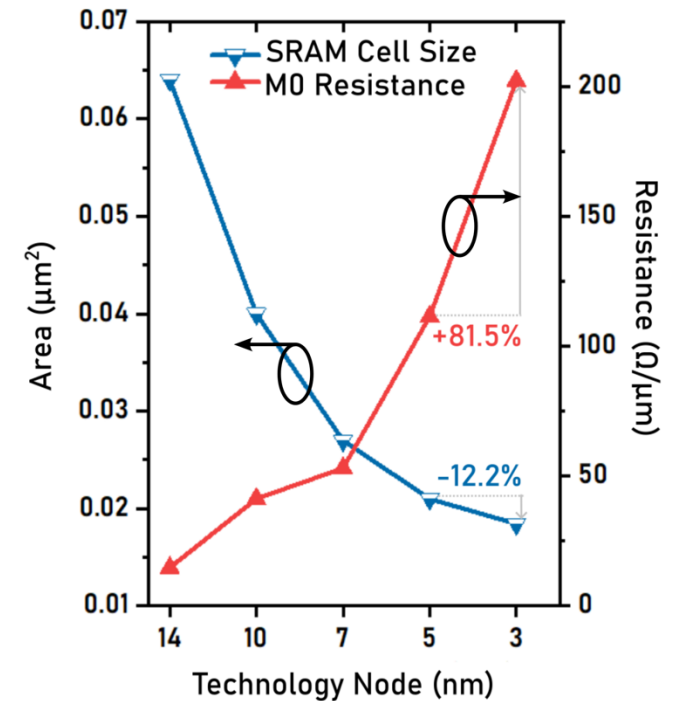


- ✗ Speed
- ✓ Density
- ✓ Stability
- ✓ Low Static Power



SRAM dominates area, and interconnects dominate the dynamic power

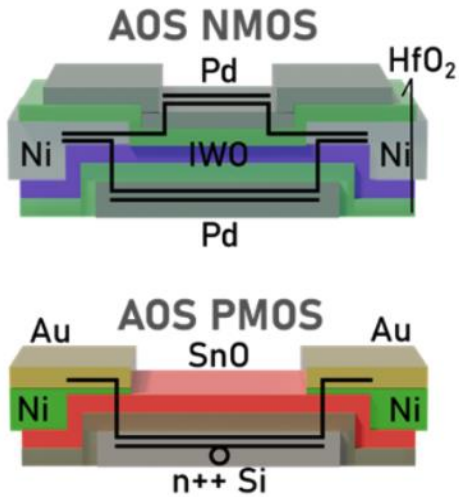
SRAM Area & M0 Unit Resistance vs Technology Node



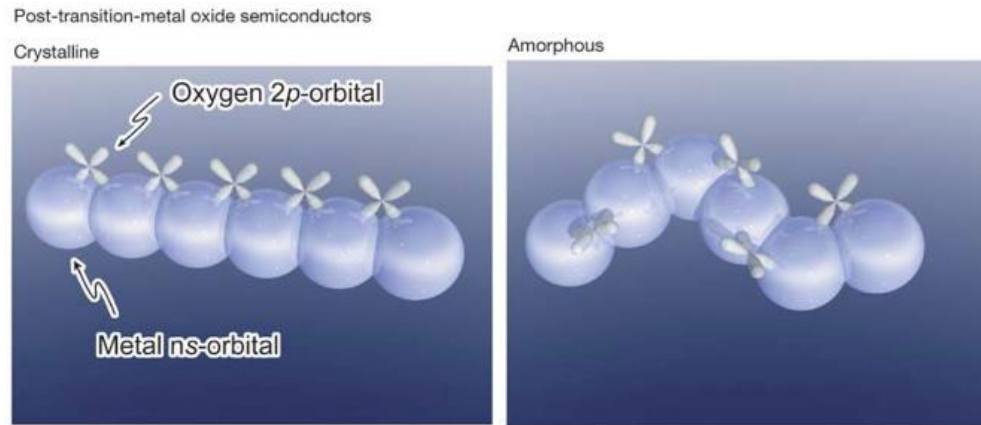
A problem exacerbated by the limits of SRAM/RDL scalability



Opportunity of AOS Channel Devices



K. Nomura et al., *Nature* 2004



What makes it BEOL compatible?

- Equilibrium phase of solids is amorphous at low temperatures
- ➔ No high-temperature annealing
- Isotropic, large S-orbital conduction
- ➔ Eases requirement on long-range order for conduction

✓ Back-End-of-Line (BEOL)
Compatible process (<400 °C)

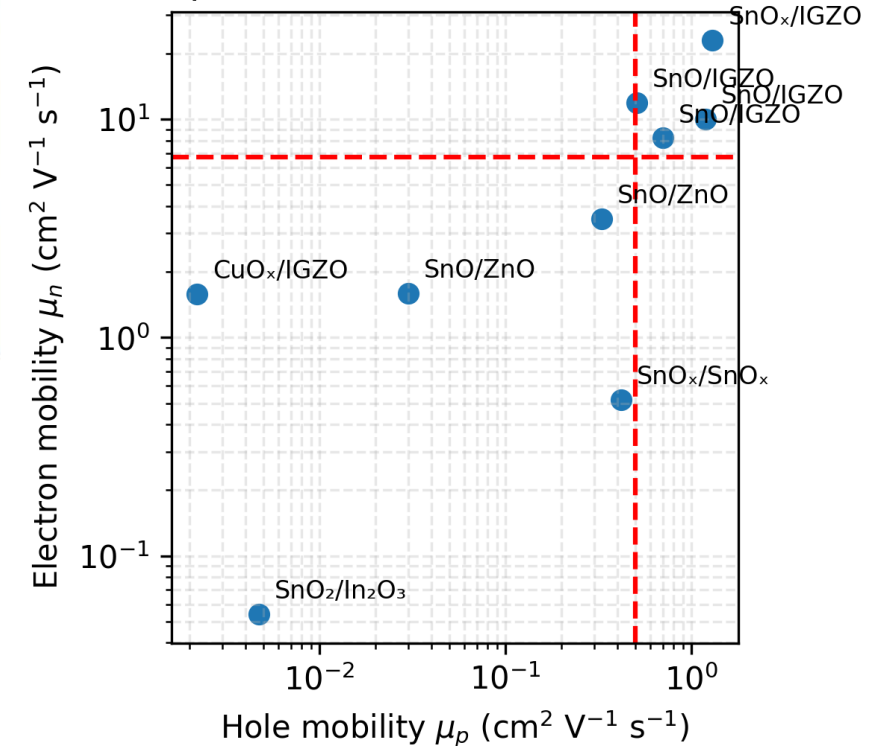
✓ Ultra-low leakage (<10⁻¹⁵ A/μm)

✓ Adequate electron/hole mobility (~20 cm²/V·s, ~2 cm²/V·s)

✓ Strong V_{th} Stability (BTI)

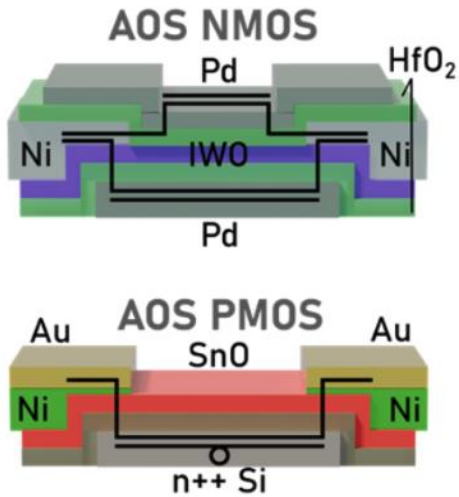


Carrier Mobilities in Reported Oxide-Channel Inverters



Z-W. Shang et al., *Nanotechnology Reviews* 2019

Opportunity of AOS Channel Devices



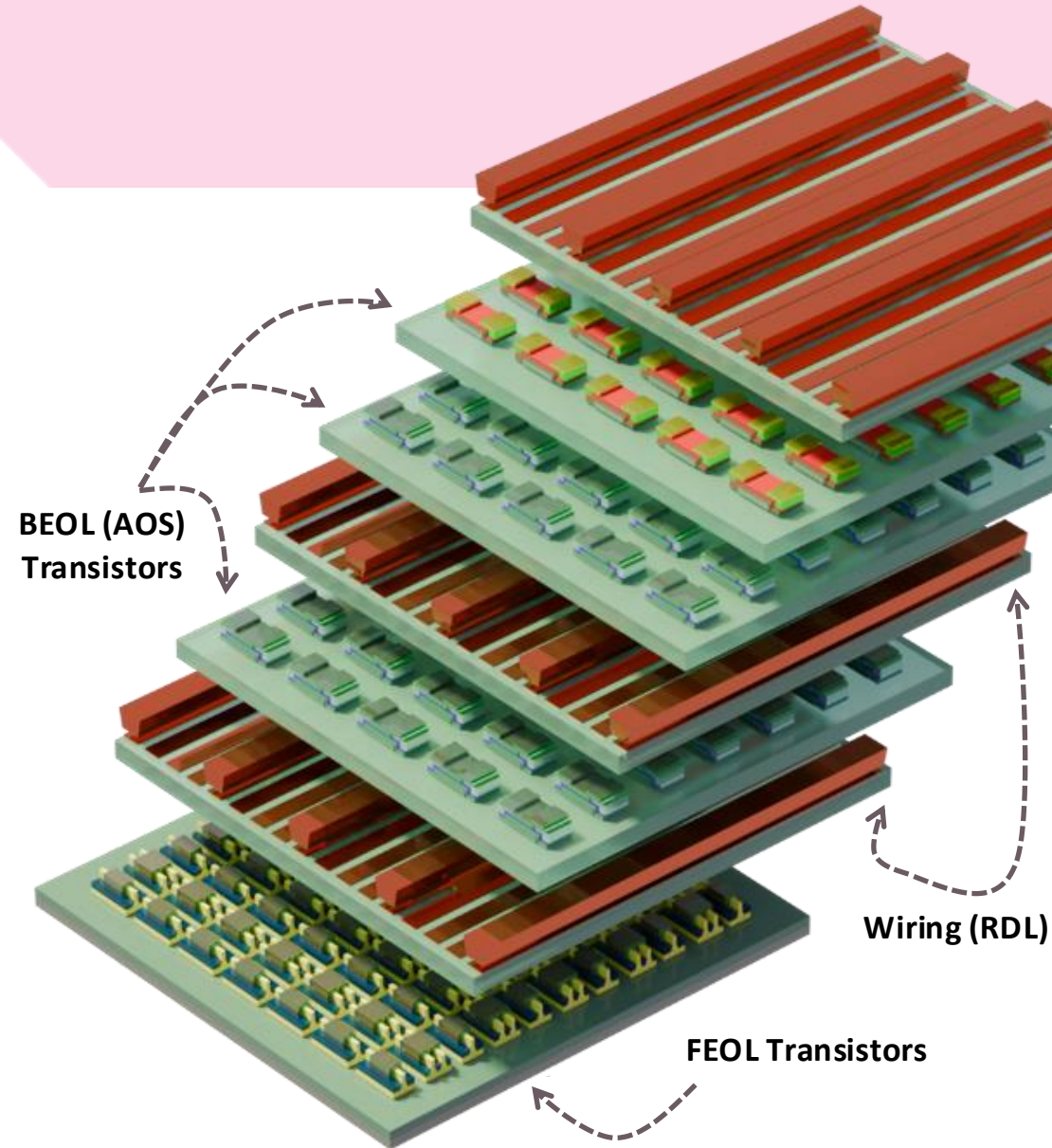
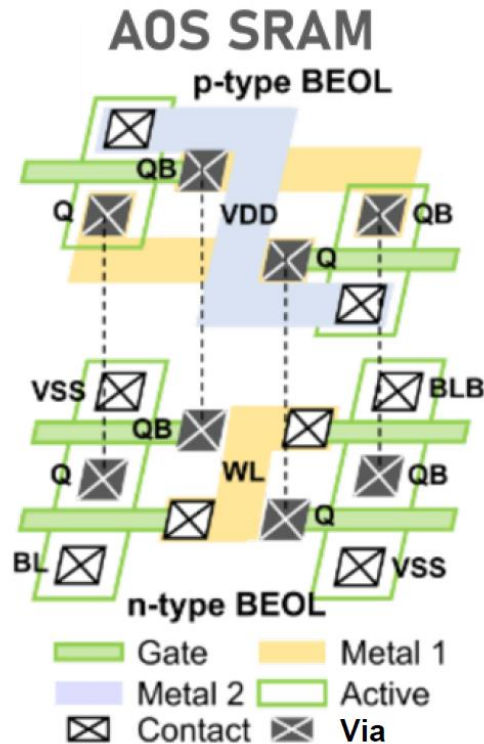
✓ Back-End-of-Line (BEOL)
Compatible process (<math><400\text{ }^\circ\text{C}</math>)

✓ Ultra-low leakage (<math><10^{-15}</math> A/ $\mu\text{m}</math>)$

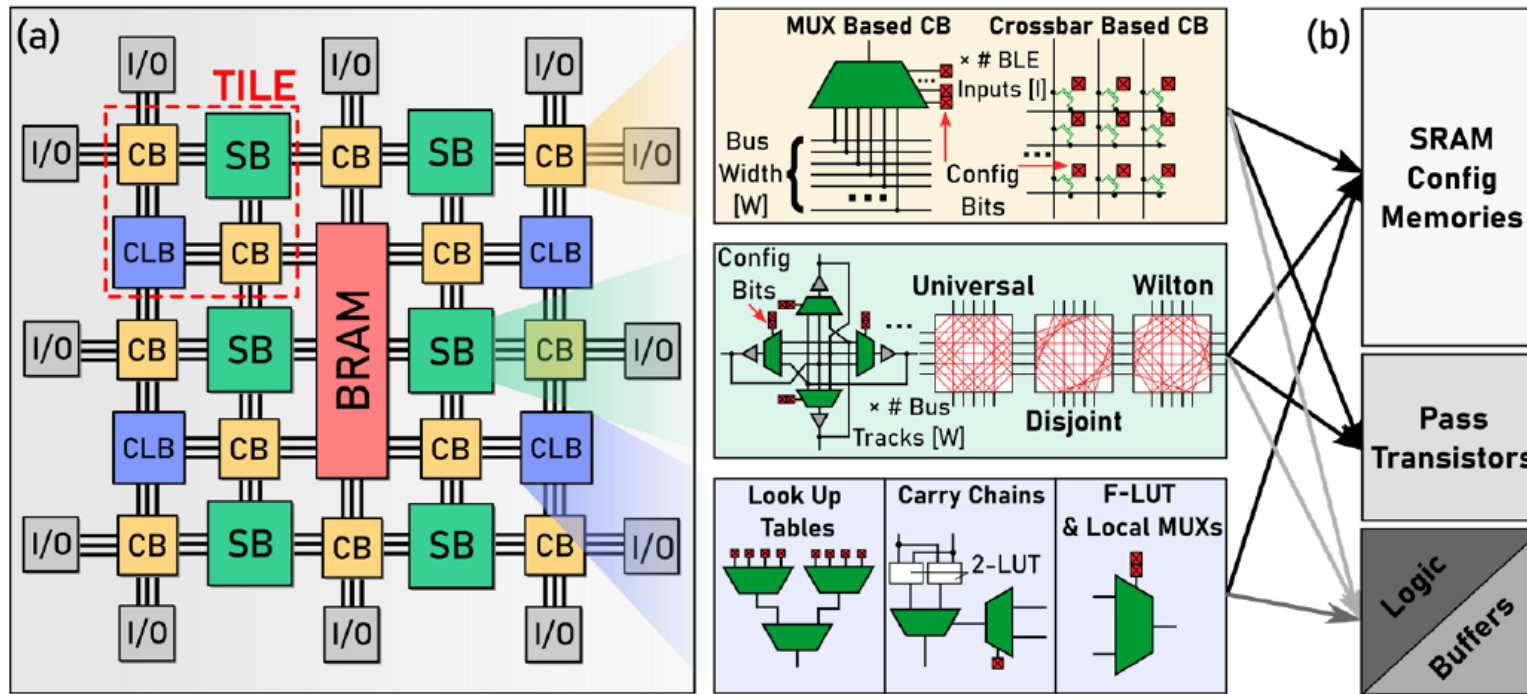
✓ Adequate electron/hole mobility ($\sim 20\text{ cm}^2/\text{V}\cdot\text{s}$, $\sim 2\text{ cm}^2/\text{V}\cdot\text{s}$)

✓ Strong V_{th} Stability (BTI)

Leverage BEOL compatibility to build a **monolithic 3D (M3D) FPGA** with **tiered active devices**

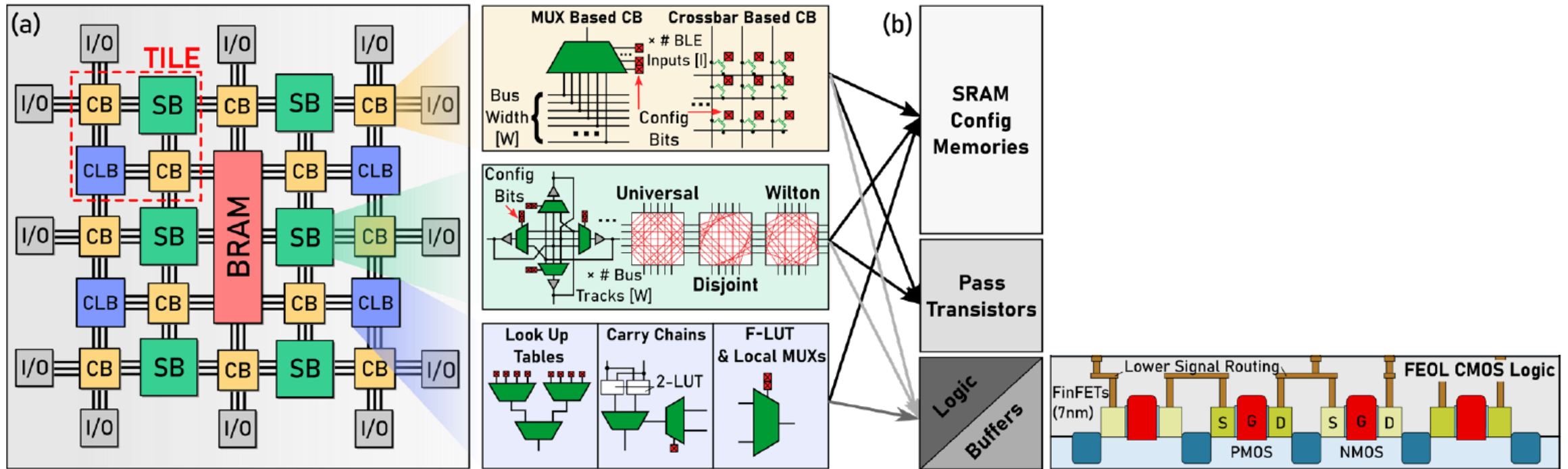


Our Proposed Design: M3D FPGA



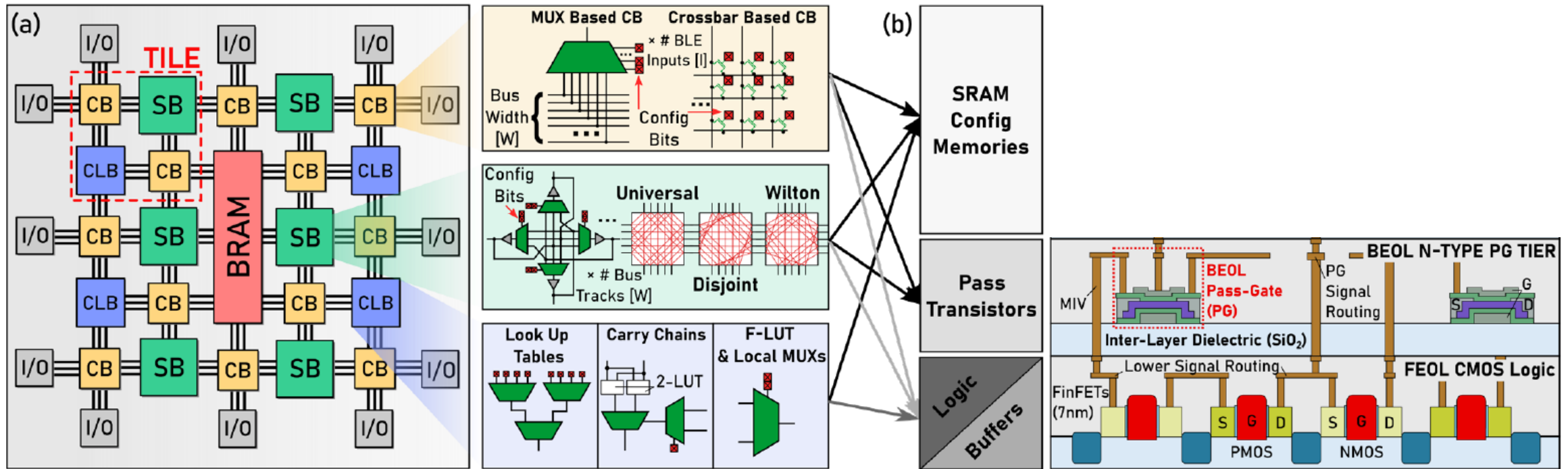
Observe the composition of each block: CBs and SBs provide connectivity to CLBs. Made mostly of configuration bits and pass transistors.

Our Proposed Design: M3D FPGA



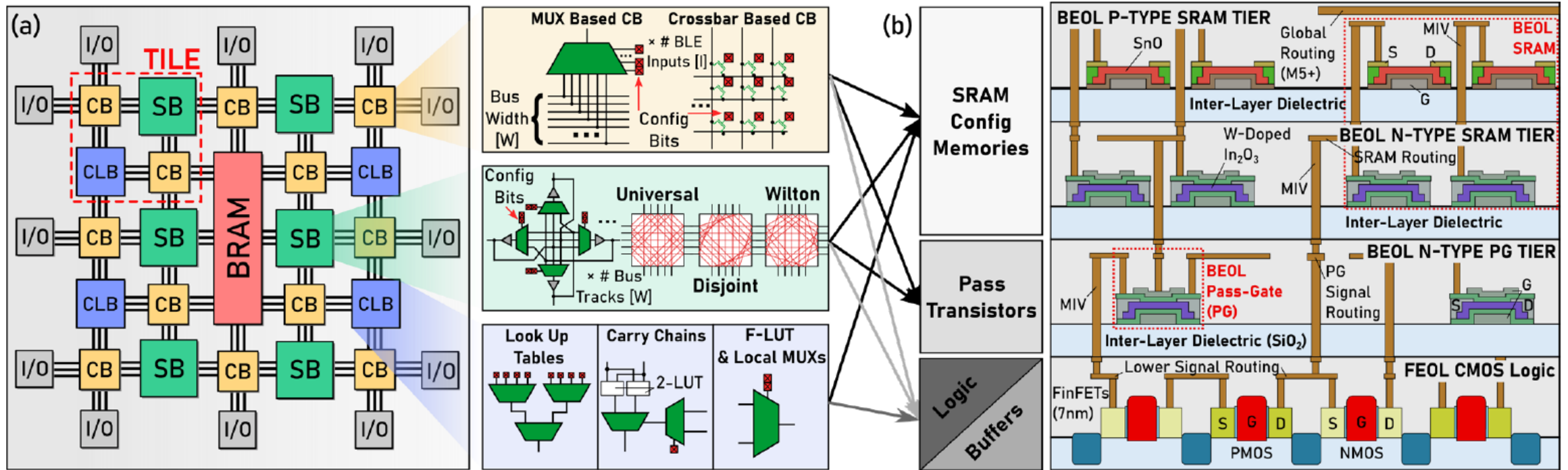
First, implement critical path logic using silicon CMOS in FEOL, such as CLBs and Buffers

Our Proposed Design: M3D FPGA



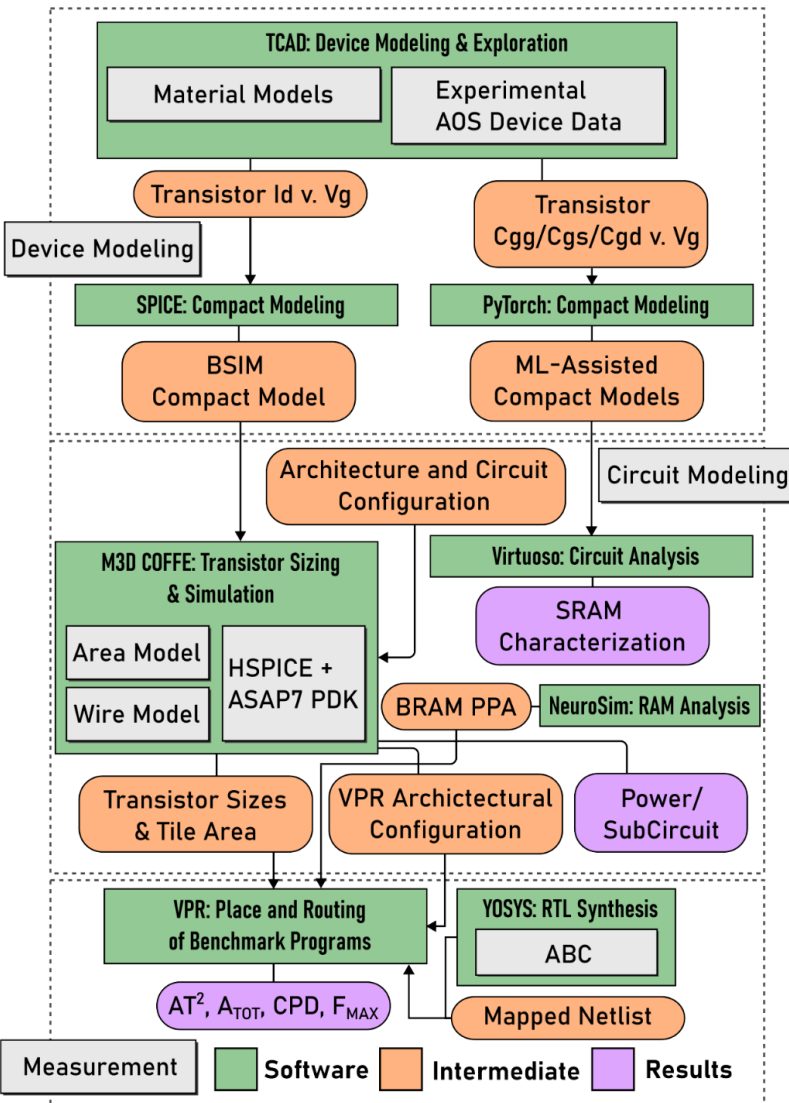
Second, build **BEOL** pass transistors (N-Type, W-doped In_2O_3) above logic and buffers. Connect using MIVs

Our Proposed Design: M3D FPGA



Finally, construct **BEOL SRAM** above Pass Transistors (P-type SnO, N-type IWO)

Methods: A TCAD to VTR Framework

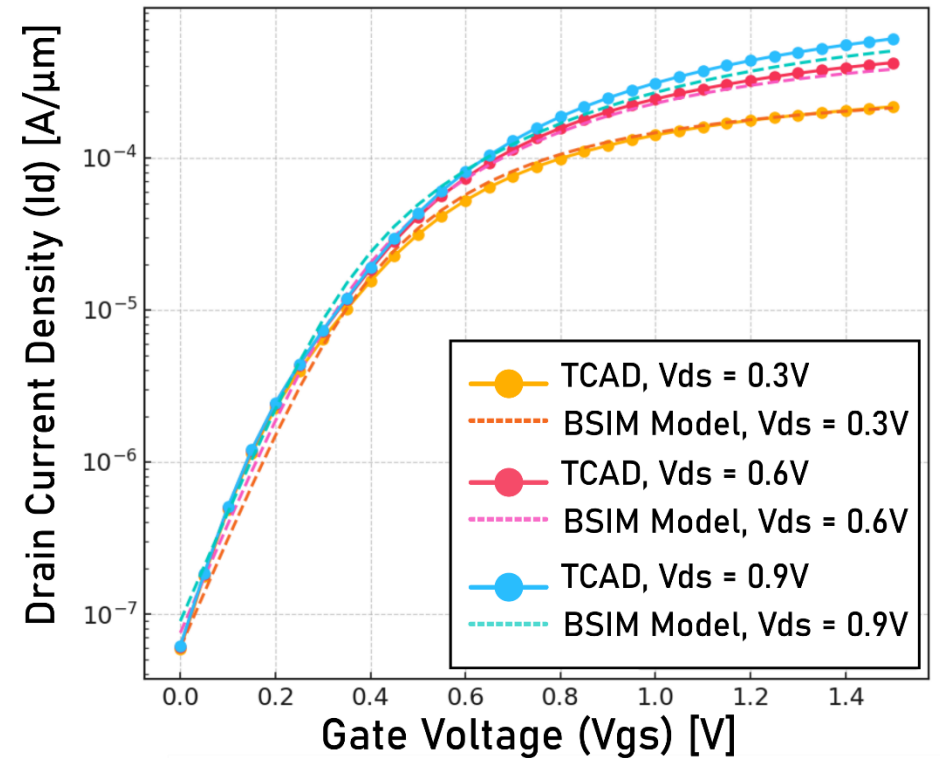


Device Modeling: Sentaurus TCAD, transferred to SPICE using BSIM & ML-Assisted Compact Models

Circuit Modeling: Use modified COFFE to optimize transistors and tile sizing using ASAP7

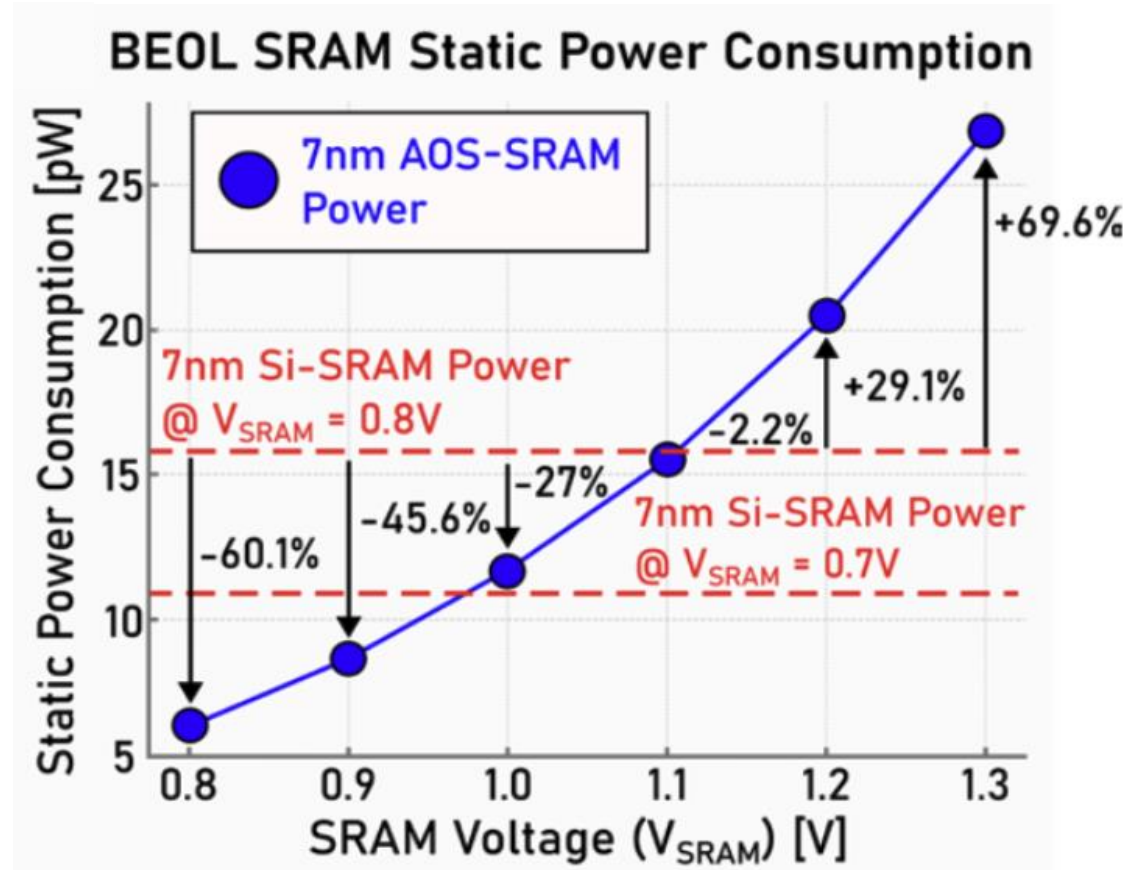
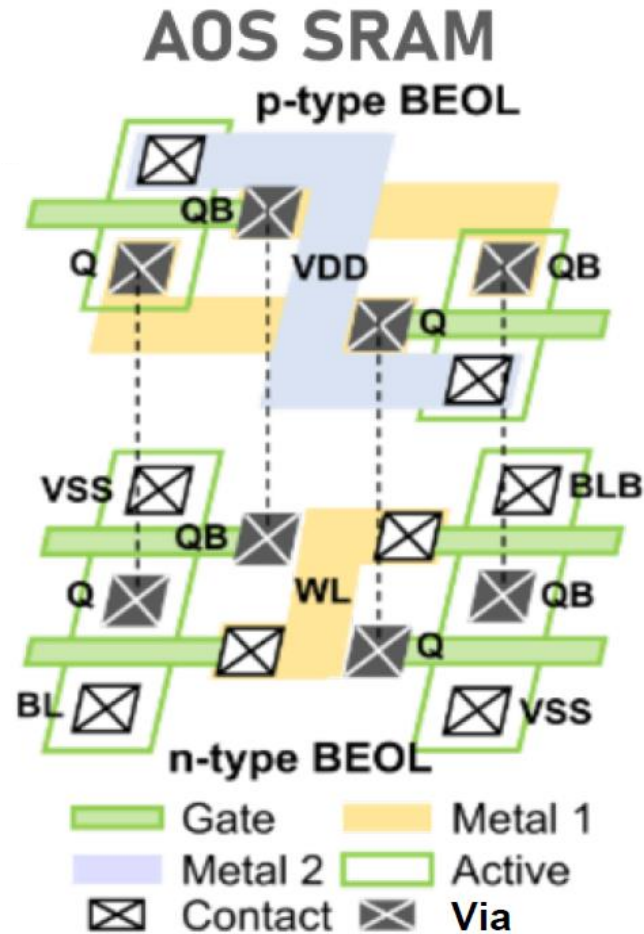
System Modeling: Use VTR tools to synthesize benchmarks and evaluate system performance

BSIM-CMG Compact Modeling Fit I_d - V_g



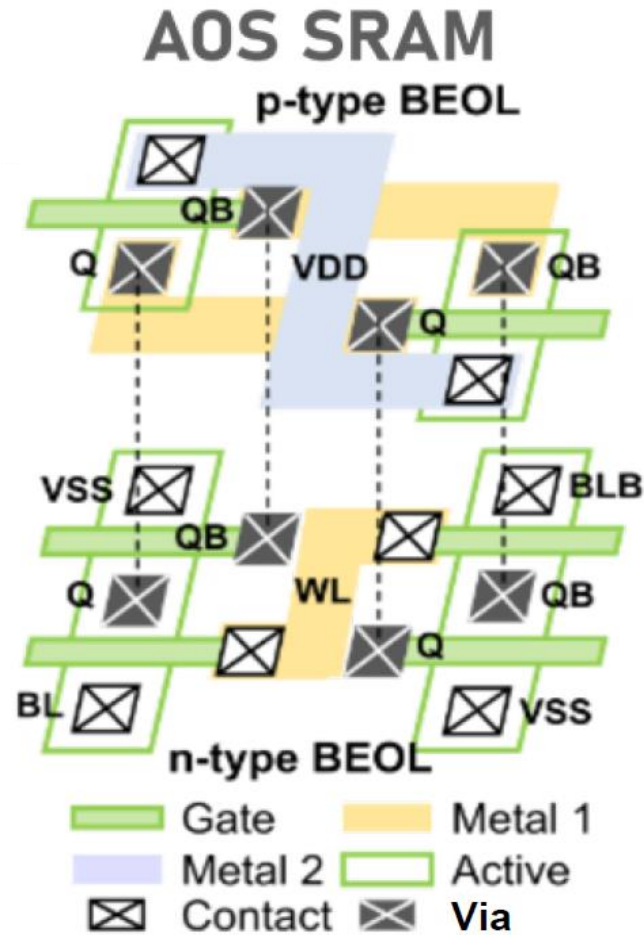
Capture C_{gs}/C_{gd} with **95.4%** accuracy. Capture I_d - V_g/I_d - V_d with **91.6%** accuracy

Characterization of AOS SRAM

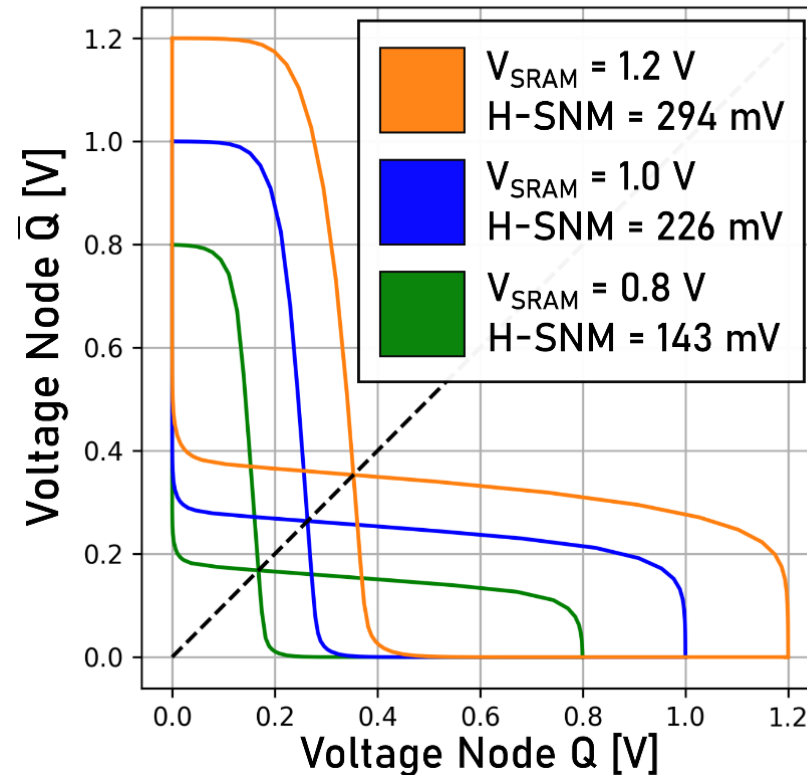


At a typical PG overdrive, static power is reduced by 60.1%
 V_{SRAM} can be increased to 1.1 V w/o a static power penalty

Characterization of AOS SRAM

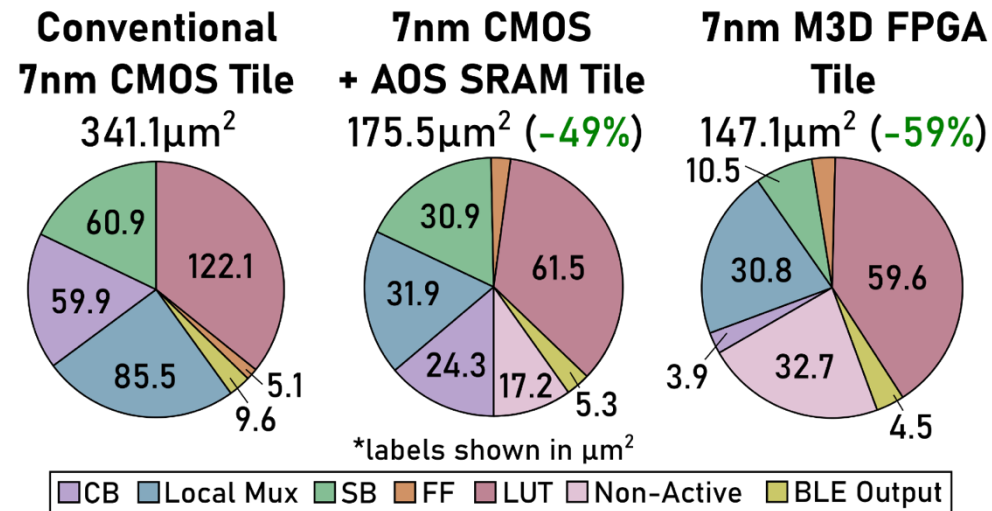


AOS SRAM Butterfly Curves & H-SNM



H-SNM is 226 mV at 1.0 V V_{SRAM} , Lower than CMOS (336 mV @ 0.8 V)
 Sufficient (> 200 mV), can be remedied with better-matched NMOS (e.g. IZO)

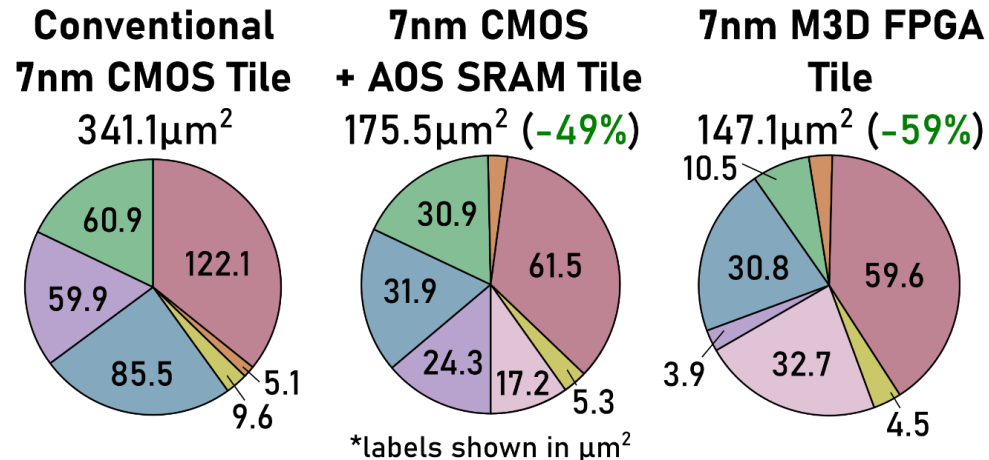
On Area Reduction



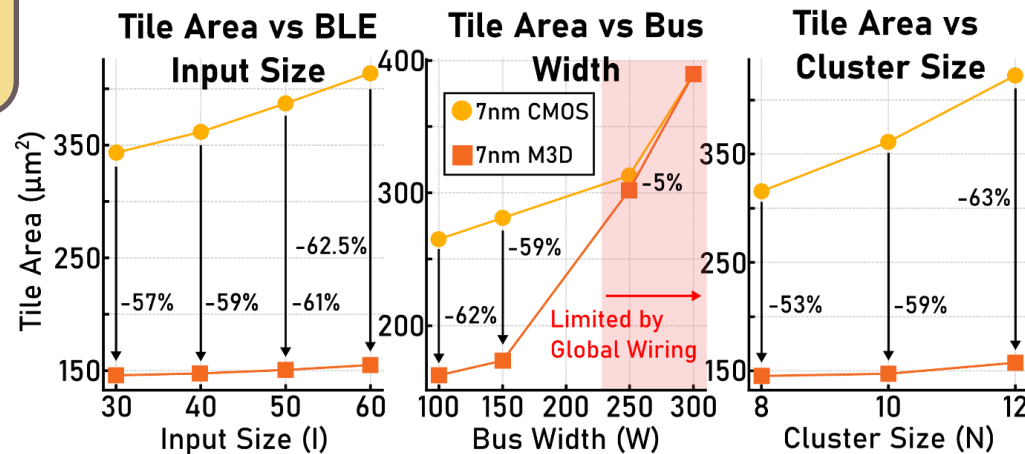
SRAM in BEOL can reduce footprint by ~49%. An additional 10% (Total ~59% reduction) with SB/CB BEOL PGs



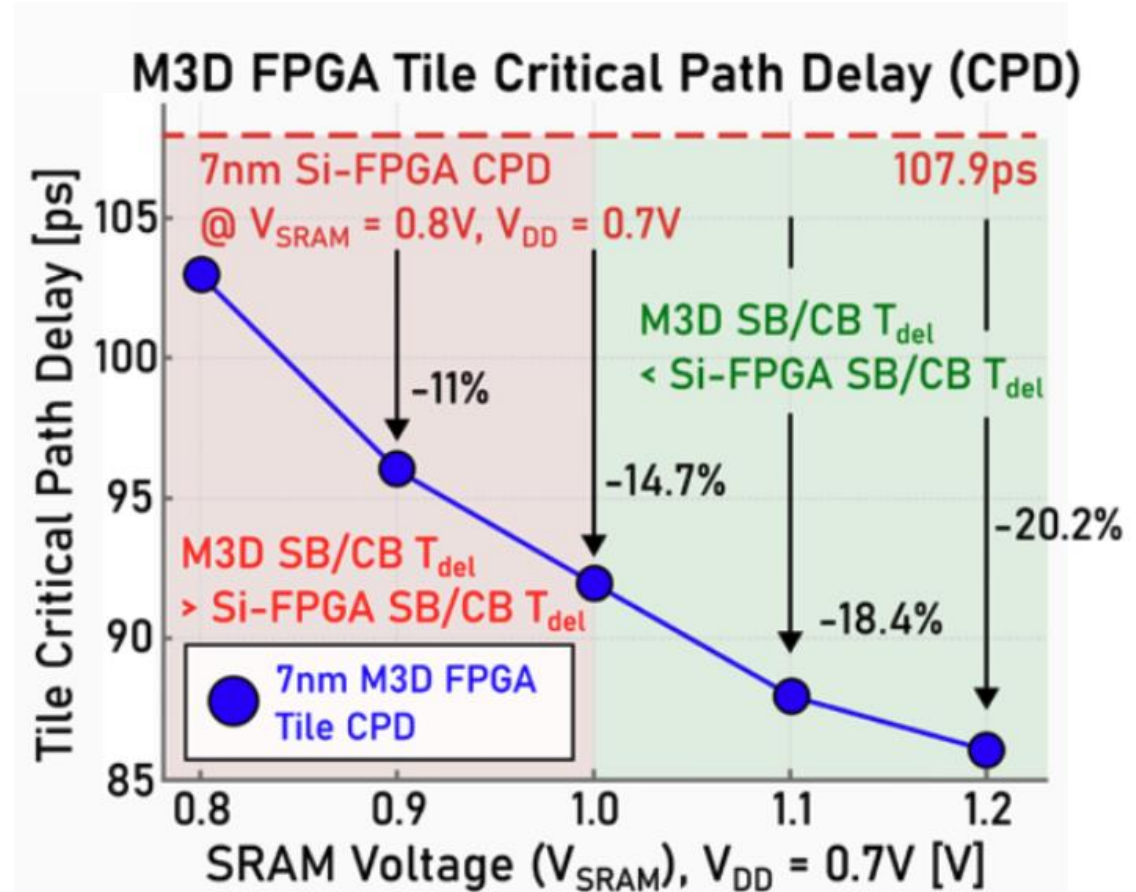
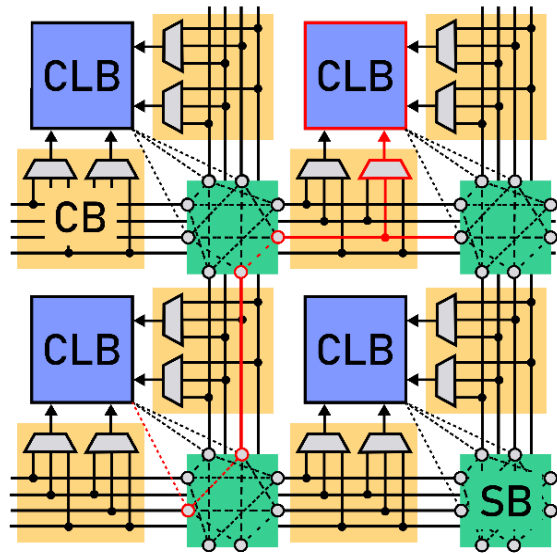
On Area Reduction



M3D Logic is Highly Scalable



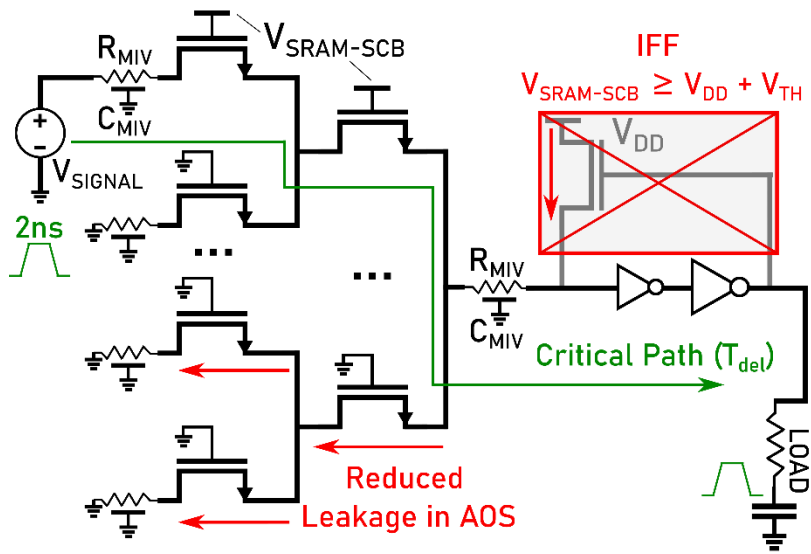
Critical Path Reduction



Reduced footprint = reduced wiring load (RC). Reduced tile level CPD (5.8%) at baseline, but PGs in SB/CB are slower!
 Overdriving (Wide BG = Less Reliability Concern) improves this further (>15%)

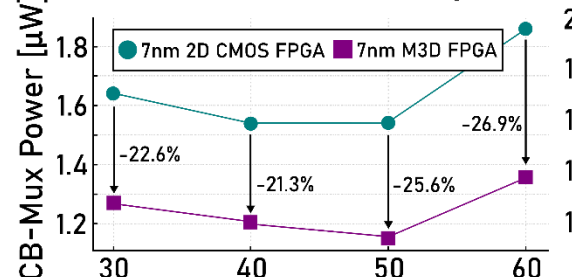


On Routing Power Reduction

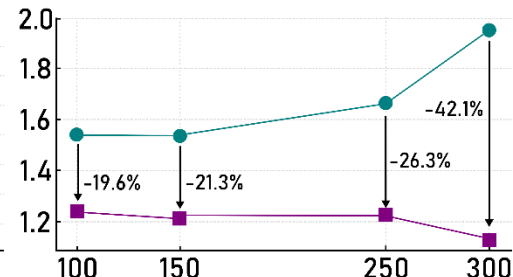


DUT: Example 2-Level M3D MUX @ 250MHz

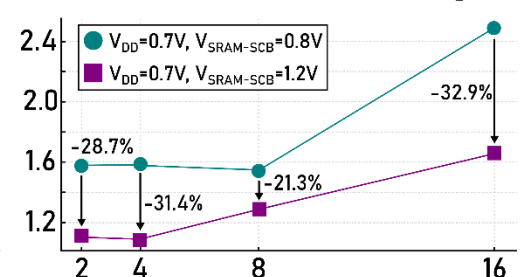
CB-Mux Power vs CLB Input Size



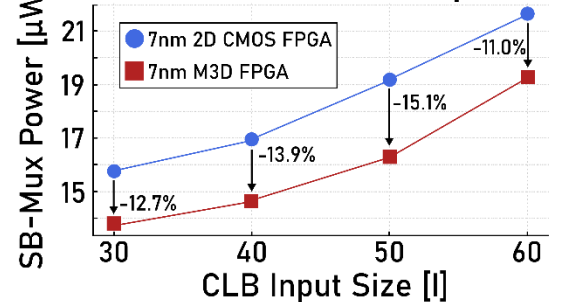
CB-Mux Power vs Bus Width



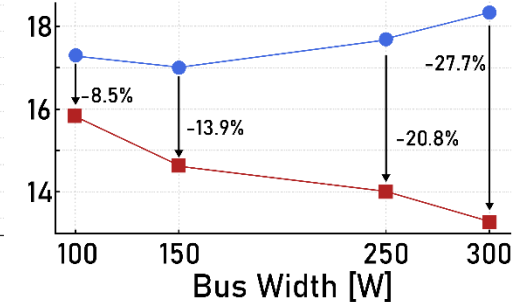
CB-Mux Power vs Bus Length



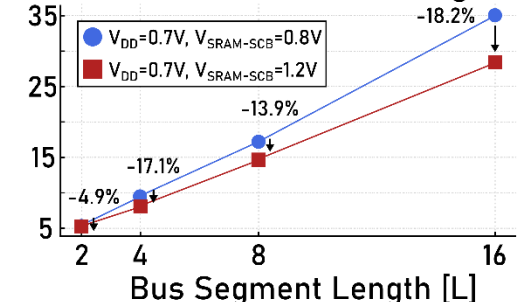
SB-Mux Power vs CLB Input Size



SB-Mux Power vs Bus Width



SB-Mux Power vs Bus Length



Wide BG in AOS – remove level-restorer from PG mux. Reduce total power in reconfigurable mesh

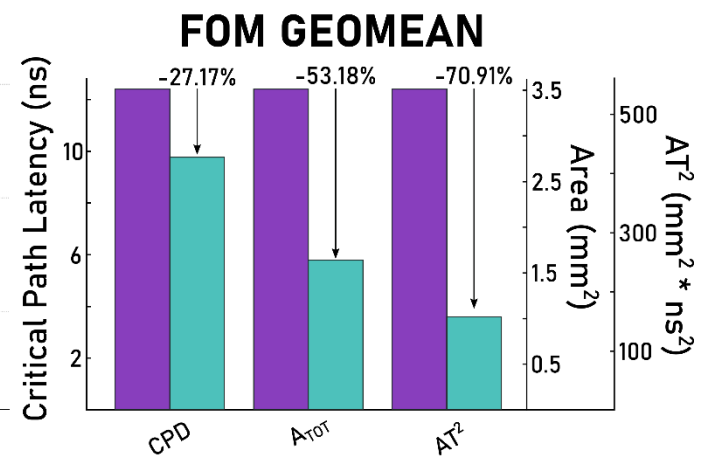
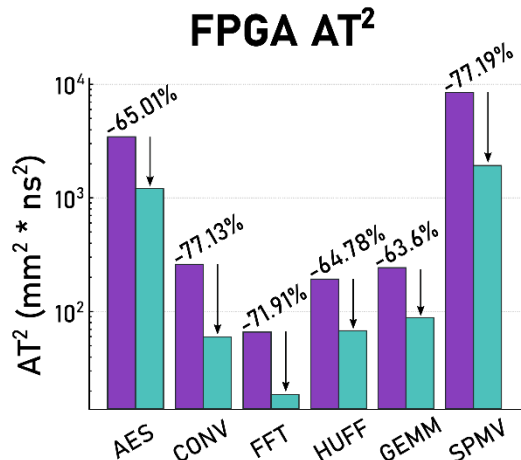
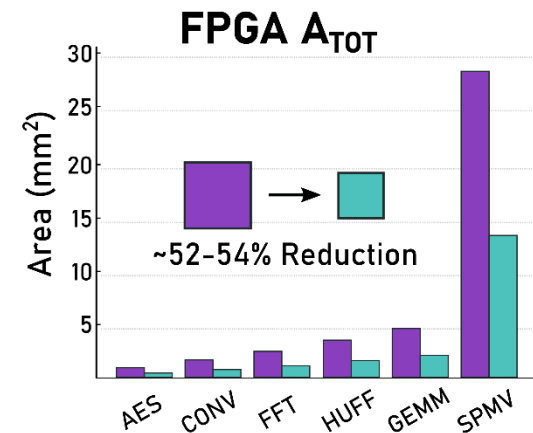
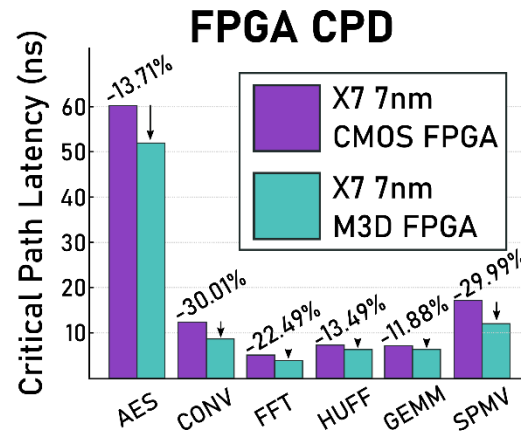
Reductions in SB-Mux power 17.6%, CB-Mux power 13.8% on average



Benchmarking on Modern Workloads

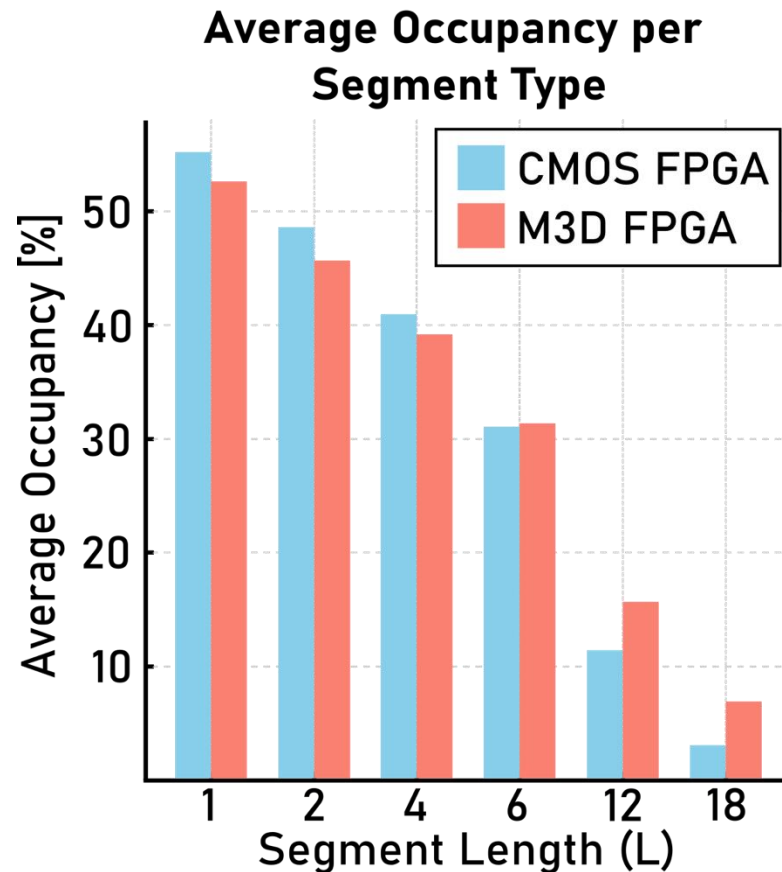
Key Kernel Benchmarking Parameters

Benchmark	Parameters	Input Type
FFT	Input Data Size $n = 64$	32b FP
Convolution	Input Size = $128 \times 228 \times 228$ Kernel Size = $128 \times 128 \times 5 \times 5$ Output Size = $128 \times 224 \times 224$ 5x5 Convolution Unit Number = 16	8b INT
GEMM	Matrix Size = 512×512 Systolic Array Size = 8×8	8b INT
SPMV	Single HiSparse [39] SPMV Cluster	32b FP
AES Enc. & Dec.	Key Size = 256b	128b INT
Huffman Encoder	Input Symbol Size = 256 Max Tree-Construction Depth = 64 Max Rebalanced Tree Depth = 27	10b UINT 32b UINT

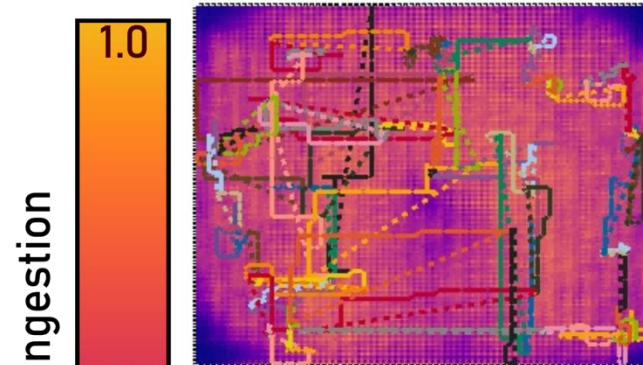


Reductions cumulatively are 27% lower CPD, 53% lower logic area utilization, and 70% AT² reduction

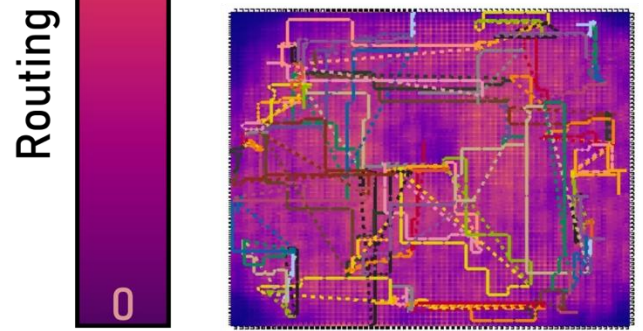
Benchmarking on Modern Workloads



Routing Congestion & Crit. Paths:
AES: 7nm CMOS FPGA

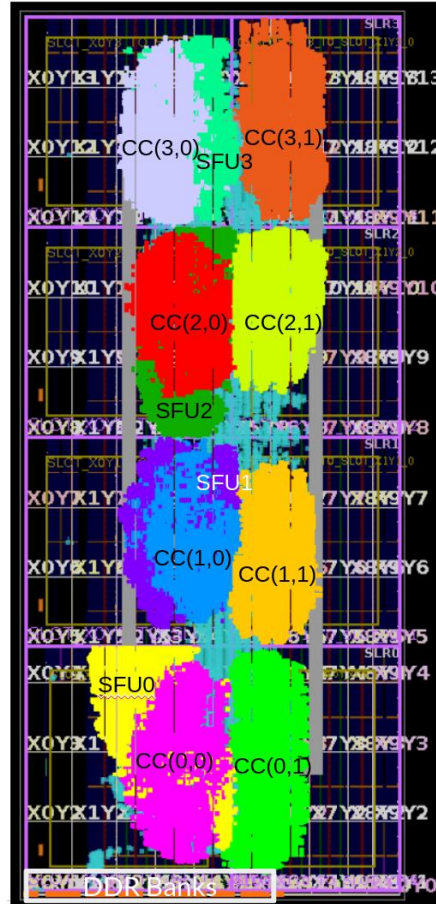


AES: 7nm M3D FPGA



Higher affinity for longer segments in M3D FPGA leading to **lower congestion, less hotspots, more opportunity for hard-block placement in local-proximity**

Benchmarking on GPT-2 Acceleration



Comparison of LLM (GPT-2) Implementation on M3D FPGA

System	Versal VPK 180	7 nm M3D FPGA
Area	736.0 mm ²	423.8 mm ² (-42.4%)
Delay	4.17 ns	3.31 ns (-20.6%)
AT ²	12798.9 mm ² ×ns ²	4643.2 mm ² ×ns ² (-63.7%)



Using Xilinx Versal VPK180 FPGA, Vivado 2023.2, and resource utilization from GEMM, translate findings onto an equivalent M3D FPGA. **Significant A, CPD and AT² reduction on GPT-2 Class LLM**

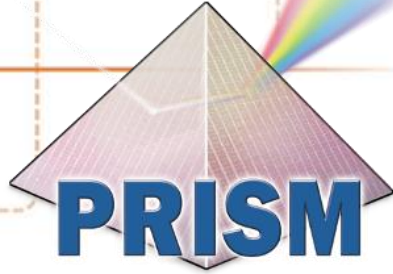
Conclusions

- We present **DSE of an M3D FPGA** based on **AOS configuration memories** and reconfigurable mesh
- The proposed FPGA:
 - i.) Is **programmable at logic-compatible voltages** (an issue in other eNVM designs)
 - ii.) Has up to **30% lower critical path delay**, **54% lower area util**, **77% lower AT^2**
 - iii.) **reduce ASIC/FPGA performance disparity to 2.05-4.37 \times** and area to **4.2 \times**
- AOS SRAMs can **reduce static power by $\sim 60\%$** while **maintaining adequate H-SNM**
- Demonstrate an **average reduction in SB/CB power by $\sim 15\%$**

This Work Was Supported By:



Semiconductor
Research
Corporation



Faaïq Waqar

Ph.D. Student, Georgia Tech



Jiahao Zhang

Ph.D. Student, UCLA



Anni Lu

Senior Engineer, NVIDIA



Zifan He

Ph.D. Student, UCLA



Jason Cong

Volgenau Chair, Distinguished
Professor, UCLA



Shimeng Yu

Dean's Professor, Georgia
Tech

**Thank you for listening
Questions?**

Email: faaiq.waqar@gatech.edu, shimeng.yu@ece.gatech.edu

