# Leveraging Monolithic 3D Technology for Data-Centric Applications

**Lightening Talk at the EES2 Technical Workshop
September 14, 2022**

**Emre Salman and Ayse K. Coskun (Boston University)**

**Department of Electrical and Computer Engineering
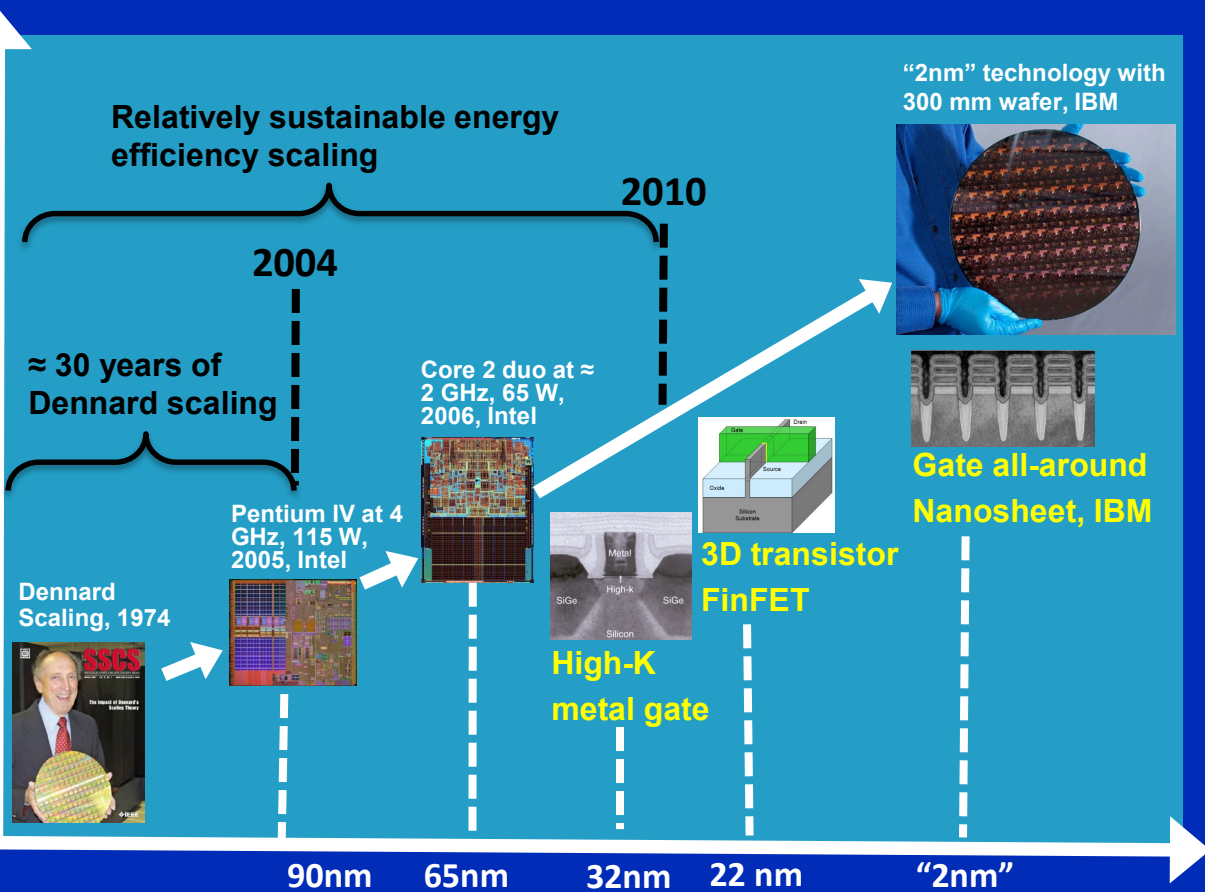Stony Brook University (SUNY), Stony Brook, New York**

Stony Brook University
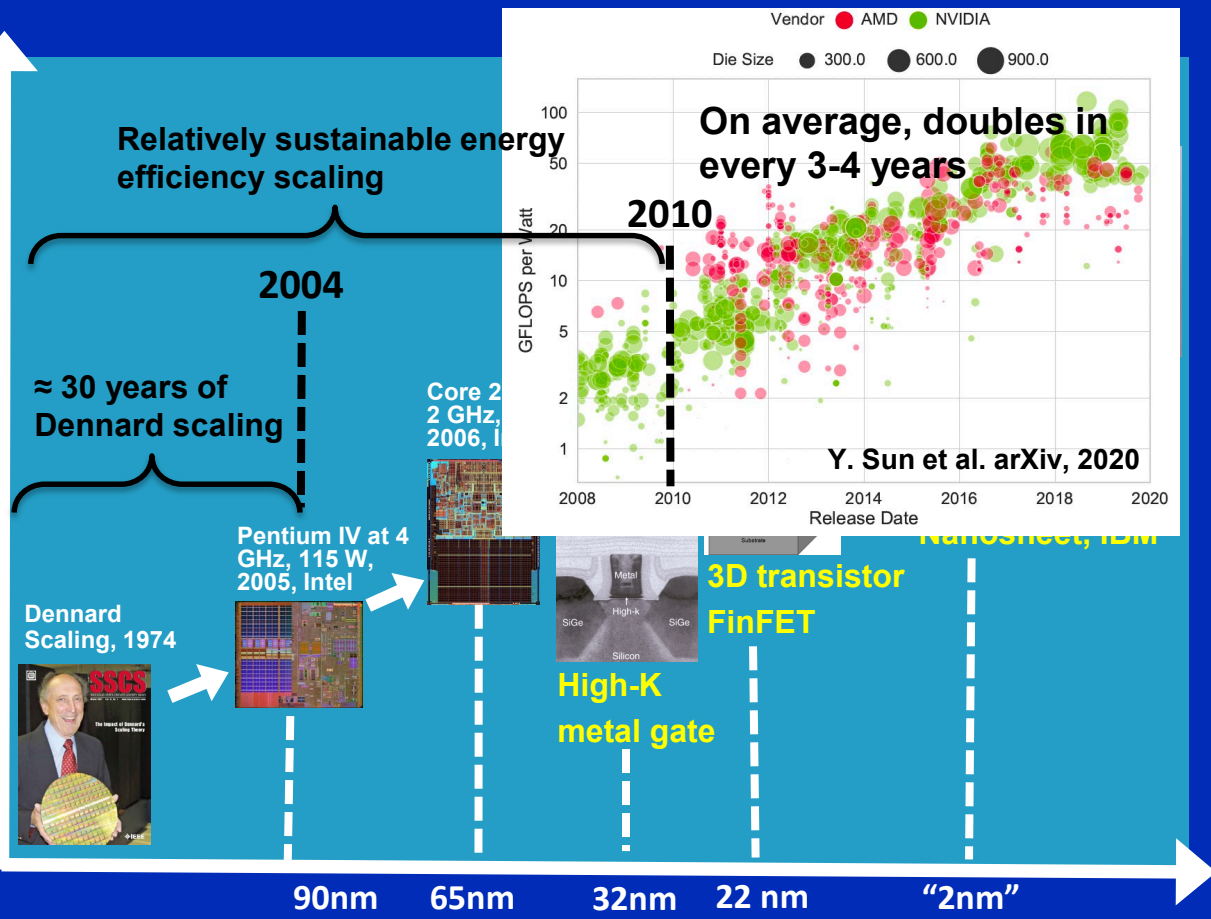
**https://nanocas.ece.stonybrook.edu**

**emre.salman@stonybrook.edu**

# Sustainable Energy Efficiency Scaling



Number of devices

Relatively sustainable energy efficiency scaling

≈ 30 years of Dennard scaling

2004

2010

"2nm" technology with 300 mm wafer, IBM

Core 2 duo at ≈ 2 GHz, 65 W, 2006, Intel

Pentium IV at 4 GHz, 115 W, 2005, Intel

Dennard Scaling, 1974

High-K metal gate

3D transistor FinFET

Gate all-around Nanosheet, IBM

90nm   65nm   32nm   22 nm   "2nm"

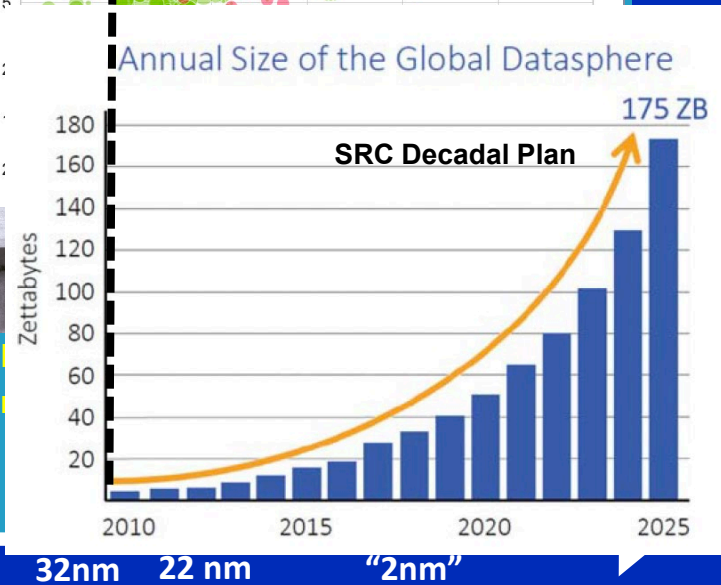# Sustainable Energy Efficiency Scaling

**Number of devices**

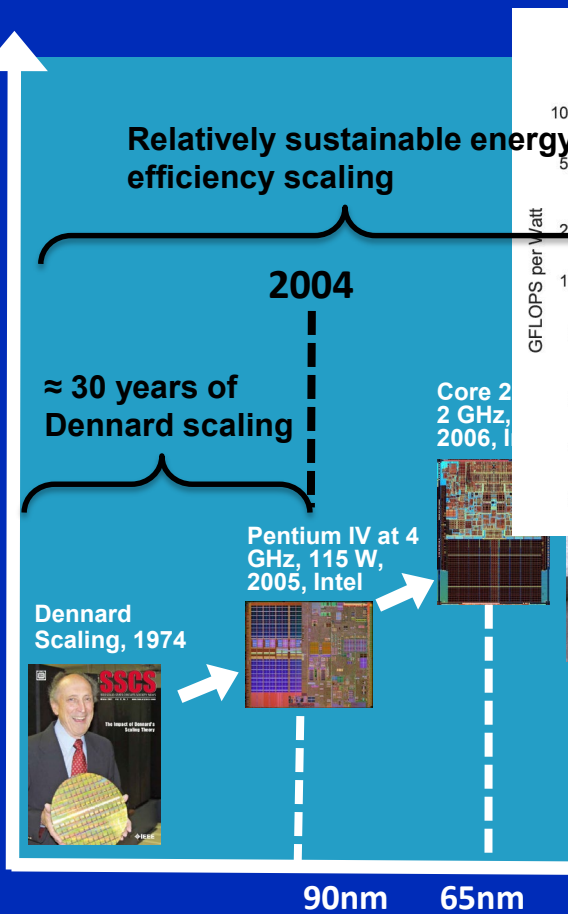**Relatively sustainable energy efficiency scaling**

**2004**

**≈ 30 years of Dennard scaling**

**Core 2 2 GHz, 2006, I**

**Pentium IV at 4 GHz, 115 W, 2005, Intel**

**Dennard Scaling, 1974**

Vendor ● AMD ● NVIDIA

Die Size ● 300.0 ● 600.0 ● 900.0

**On average, doubles in every 3-4 years**

**2010**

GFLOPS per Watt

100
50

20

10

5

2

1

**Y. Sun et al. arXiv, 2020**

2008    2010    2012    2014    2016    2018    2020

Release Date

**NanoSheet, IBM**

**3D transistor FinFET**

**High-K metal gate**

**90nm    65nm    32nm    22 nm    "2nm"**
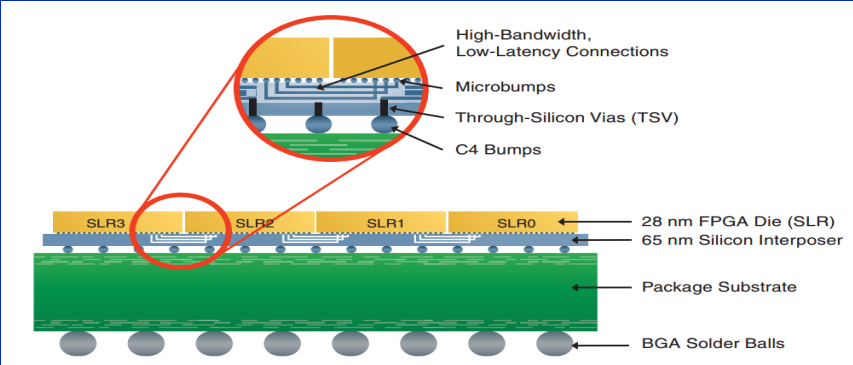
# Sustainable Energy Efficiency Scaling



**Number of devices**

**Relatively sustainable energy efficiency scaling**

**2004**

**≈ 30 years of Dennard scaling**

**Dennard Scaling, 1974**

**Pentium IV at 4 GHz, 115 W, 2005, Intel**

**Core 2, 2 GHz, 2006, I...**

Vendor ● AMD ● NVIDIA

Die Size ● 300.0 ● 600.0 ● 900.0

**On average, doubles in every 3-4 years**

**2010**

GFLOPS per Watt

**Y. Sun et al. arXiv, 2020**

Annual Size of the Global Datasphere

**175 ZB**

**SRC Decadal Plan**

Zettabytes

180 / 160 / 140 / 120 / 100 / 80 / 60 / 40 / 20

2010   2015   2020   2025

**90nm    65nm    32nm    22 nm    "2nm"**

- **Dominance of data-centric applications**
  - **Integration platforms for mitigating dominance of data on performance/power**
    - **2.5D (chiplet)**
    - **TSV based 3D**
    - **Monolithic 3D**

- **FLOPS / Watt**
  - **Upper bound in power**
  - **Lower bound in performance**
  - **Application specific**
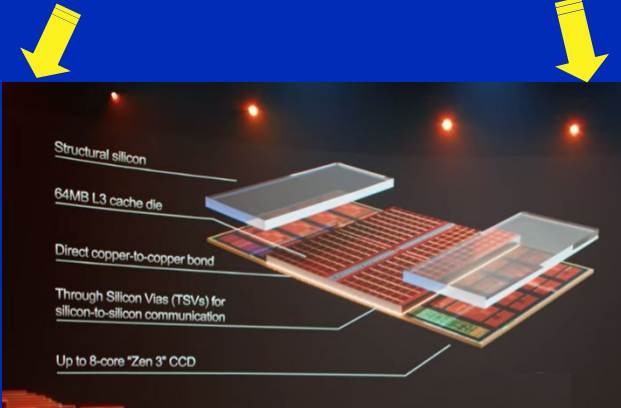    - **Cloud vs. edge**

4

# Emerging Integration/Packaging Technologies



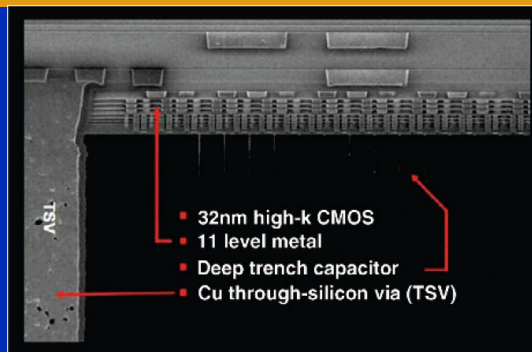Interposer Based Integration (2.5D), Xilinx



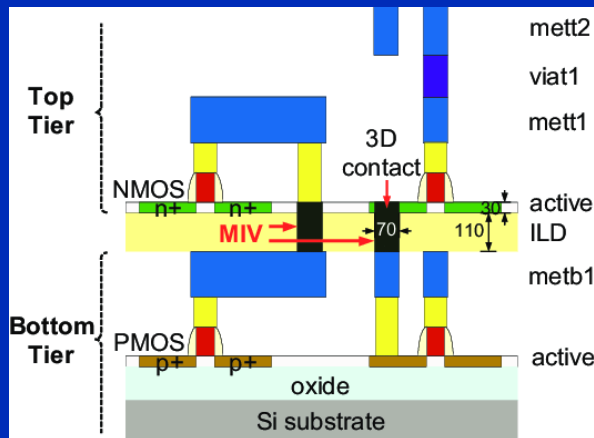Through Silicon Via (TSV), Intel FOVEROS



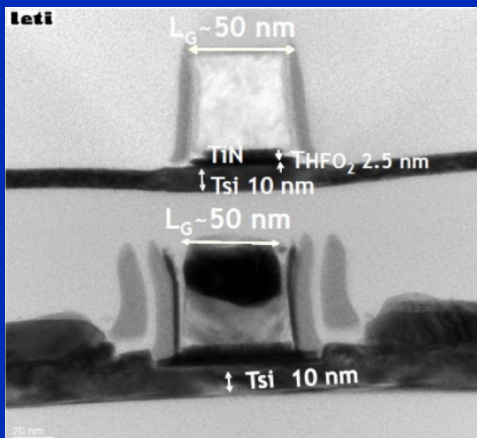Cache-core stack, AMD Ryzen 5000 series

# Monolithic 3D Integration Technology



Courtesy of IBM

- 32nm high-k CMOS
- 11 level metal
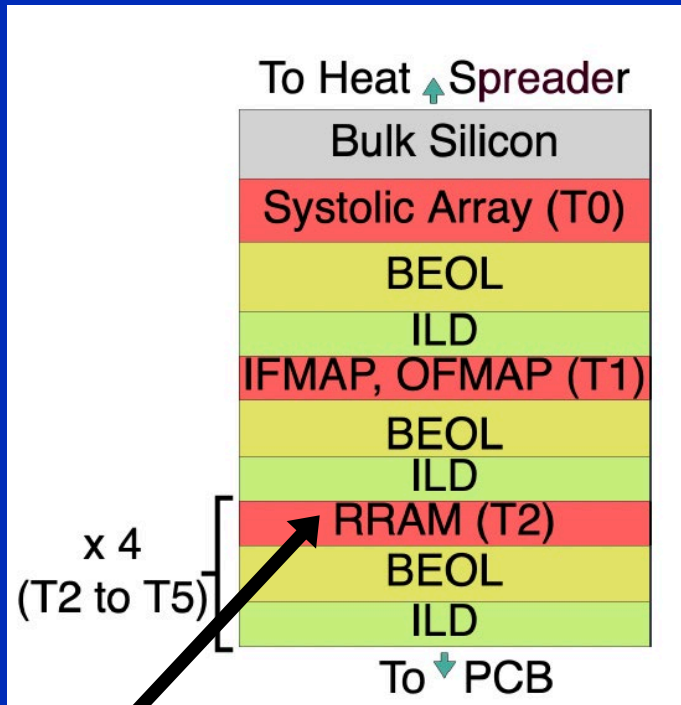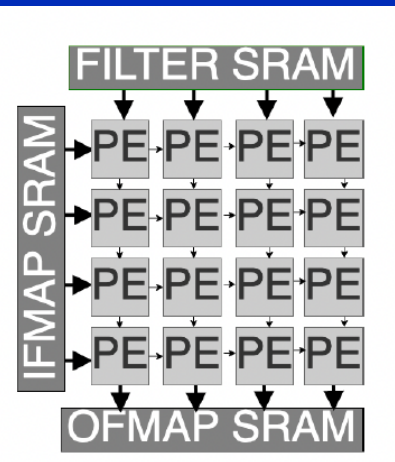- Deep trench capacitor
- Cu through-silicon via (TSV)





Monolithic 3D ICs with MIVs (LETI)

- **TSVs are large**
  - Diameter in the 2 to 10 μm range
  - Height in the 8 to 60 μm range
  - Large pitch (30 to 50 μm) and keep-out zone (KOZ) requirements

- **MIVs provide unprecedented interconnect density**
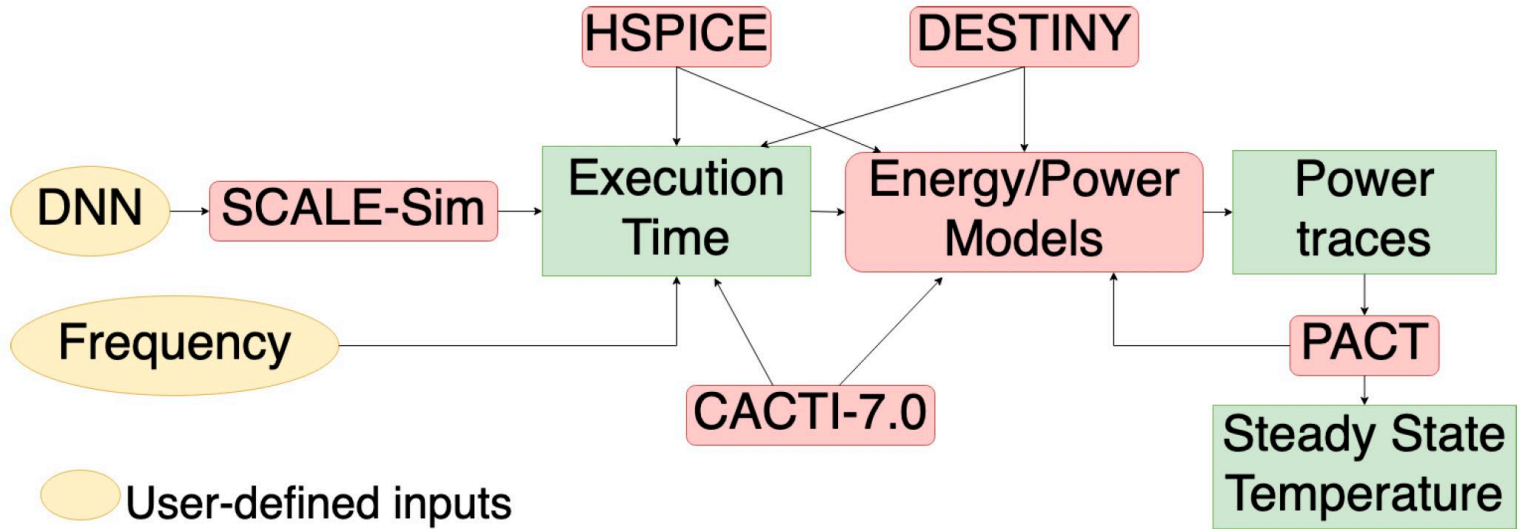  - Diameter ≈ 50 nm (similar to a metal via)
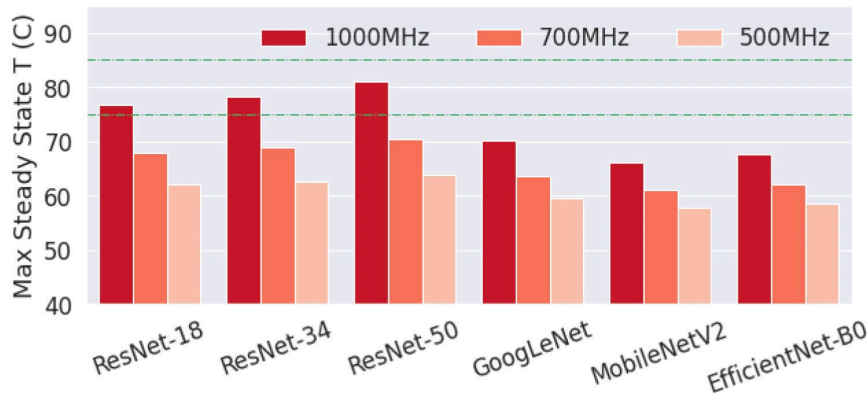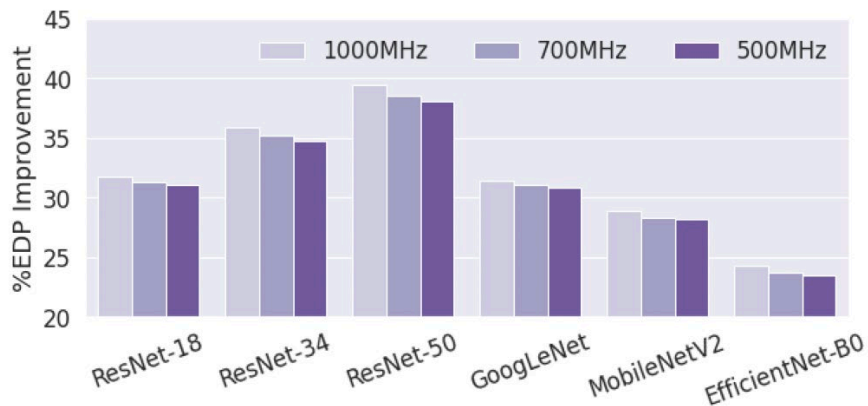
# Mono3D Deep Neural Network Accelerator



**Resistive RAM to enable high bandwidth interface**

- **Accelerator hardware in 22nm CMOS**
- **Footprint of 8 mm²**
  - **256X256 array**
  - **2 MB SRAM for IFMAP**
  - **2 MB SRAM for OFMAP**
  - **32 MB RRAM for weights**
    - **RRAM is read only**
    - **Endurance issue is mitigated**
- **Architectural optimizations**
  - **Read all weights into PE array in one cycle**
  - **Multicast all IFMAPs to all the PEs**

# Thermal-Aware Evaluation Framework

# Improvements w.r.t. 2D System



- Up to **40% reduction** in **energy-delay product**

- Up to **81% improvement** in **inference per second per Watt**

- Up to **10X improvement** in **inference per second per Watt per footprint**

# Final Notes

- **Impractical to achieve the same bandwidth with TSV-based 3D**
  - **Footprint would increase to more than 100 mm$^2$**

- **Optimizations to dataflow, architecture, and circuits are critical**
  - **Switching to Mono3D and RRAM alone does not produce desired improvements**
    - **EDP improvements less than 5%**
  - **Importance of co-design**

- **Thermal-awareness will play a key role in future advanced integration/packaging technologies**
  - **Higher power density**
  - **Temperature dependence of emerging nonvolatile memory**
  - **Edge applications with limited cooling capability**

  - **Questions: emre.salman@stonybrook.edu**