

# Accelerating DNN Applications with Emerging Memory Technologies



**Tzofnat Greenberg-Toledo**  
**Advisor: Prof. Shahar Kvatinsky**

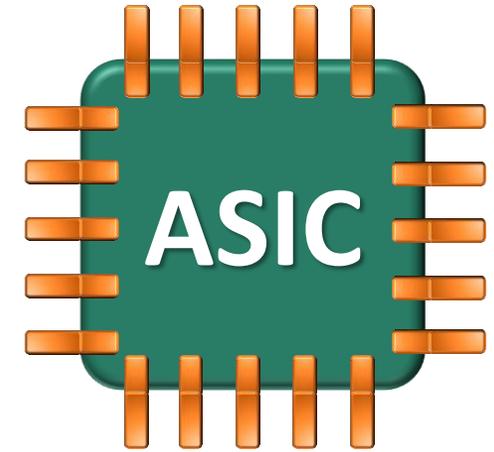
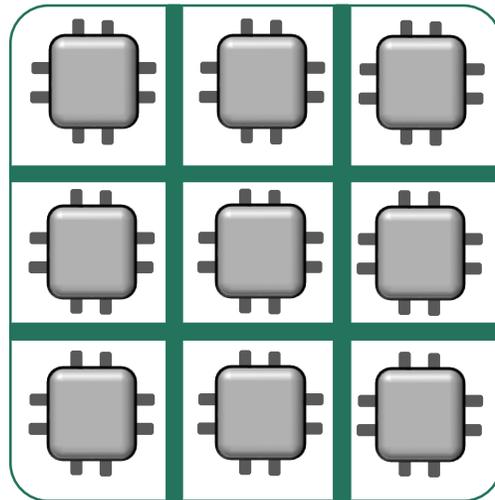
Viterby Faculty of Electrical Engineering  
Technion – Israel Institute of Technology

# DNN Hardware

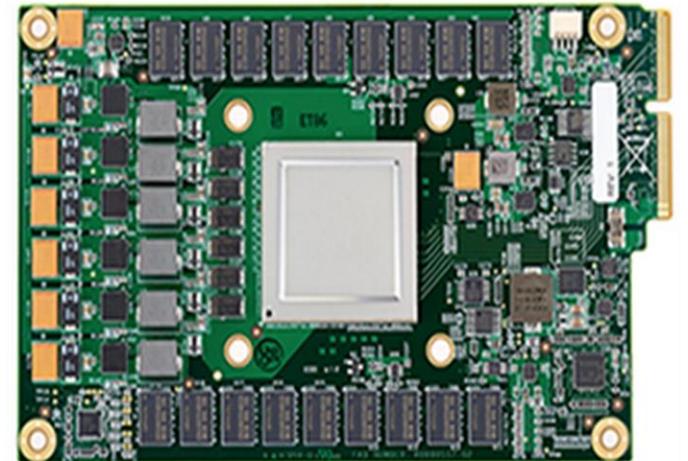
**GPU**



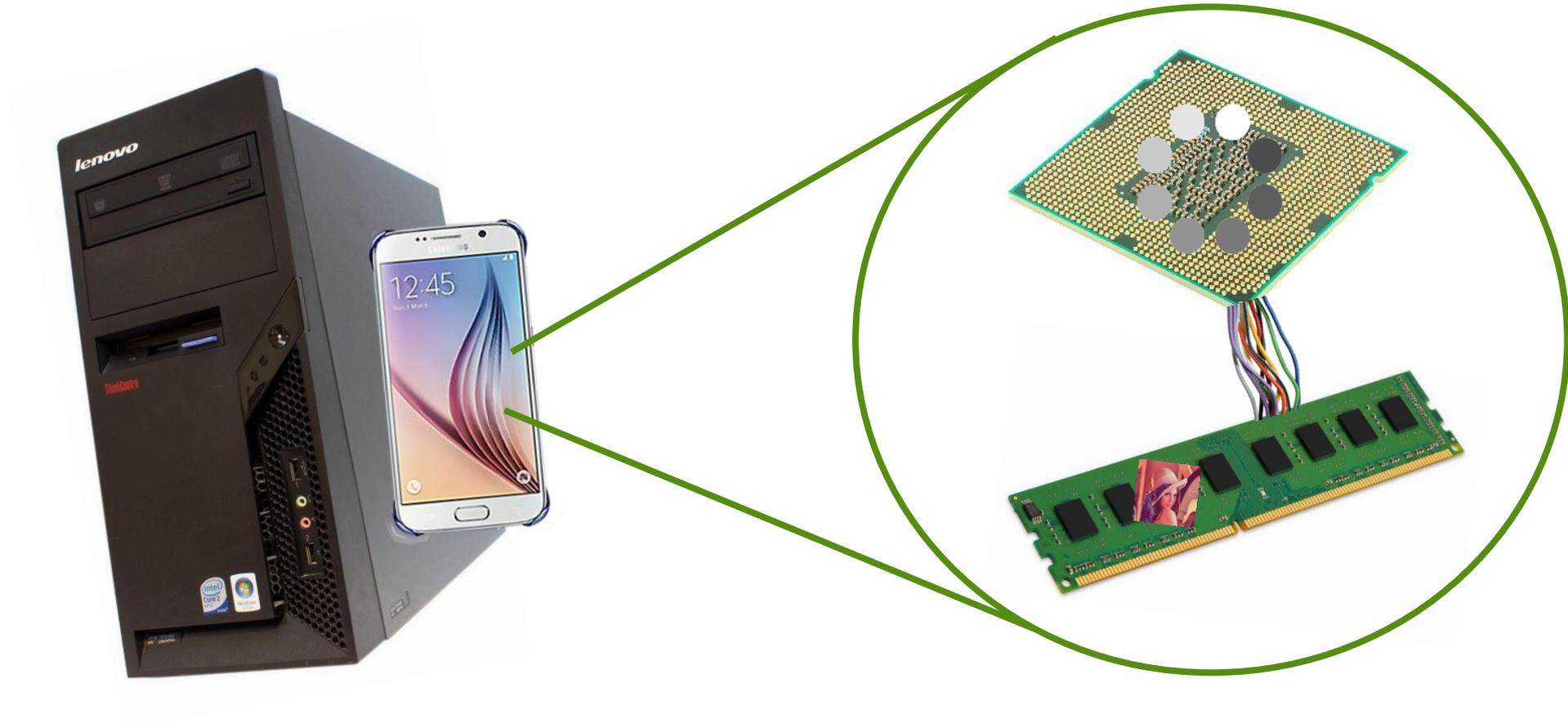
**FPGA**



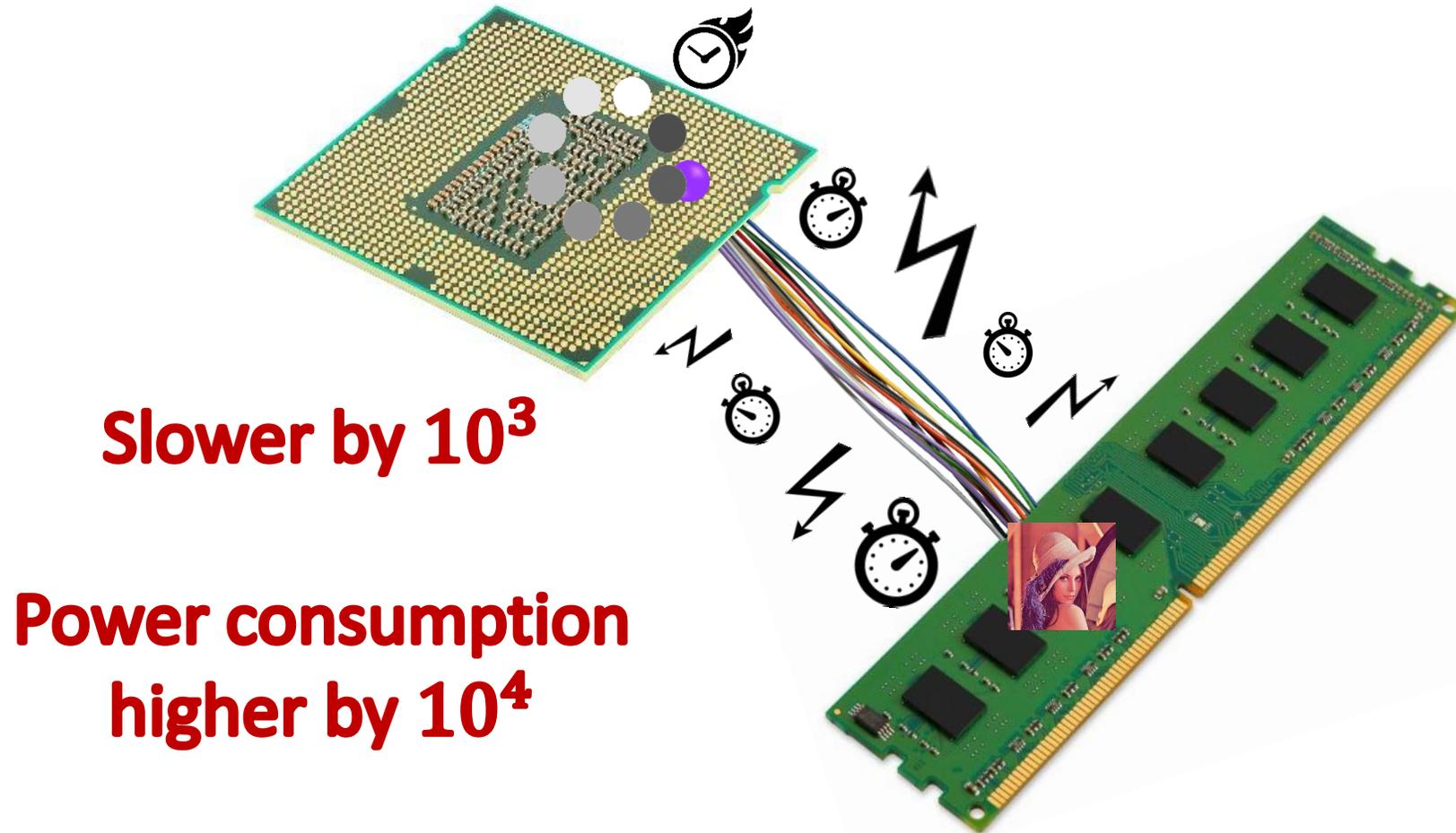
**Google TPU**



# von Neumann Architecture

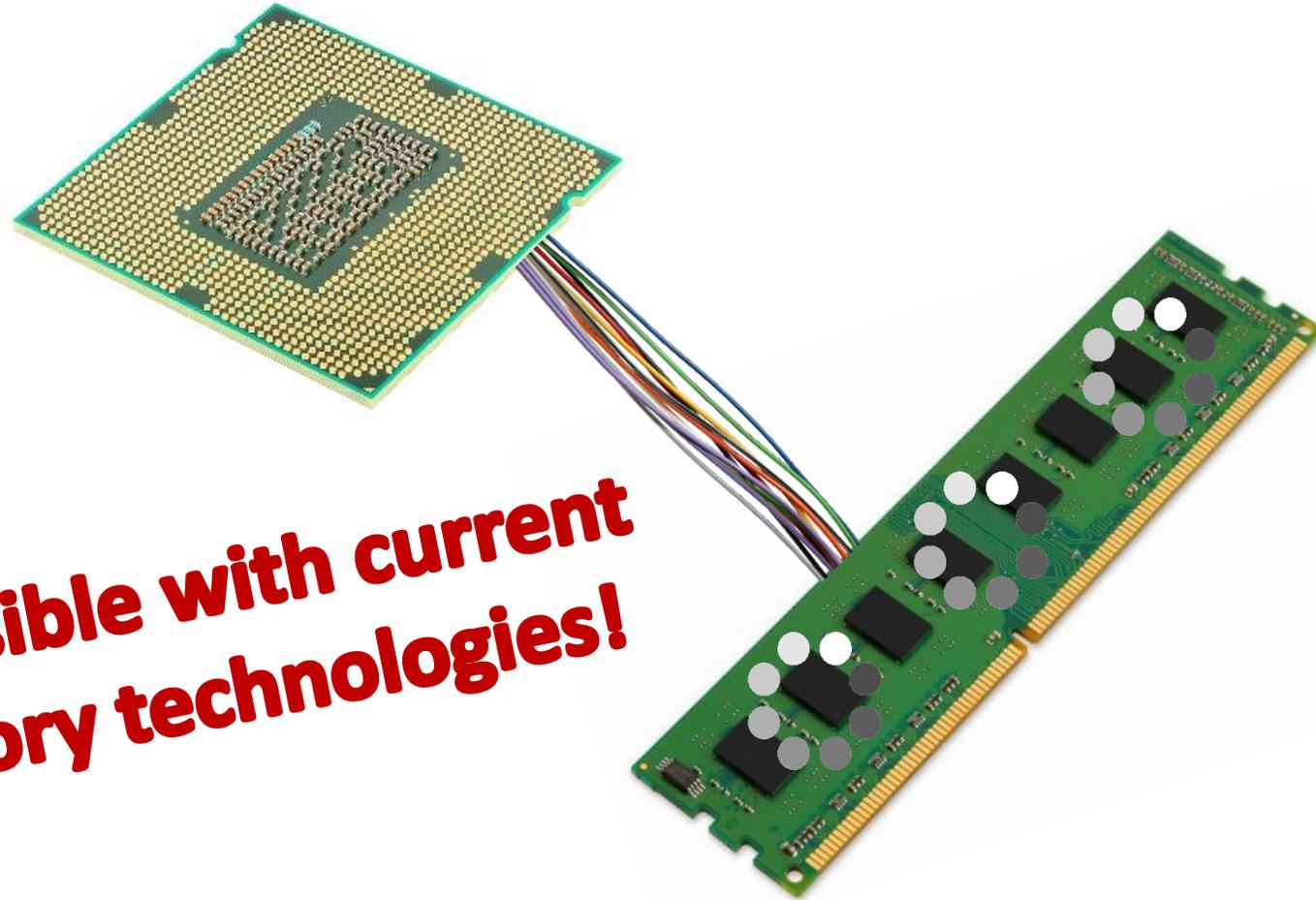


# Processing Data in von Neumann Architecture



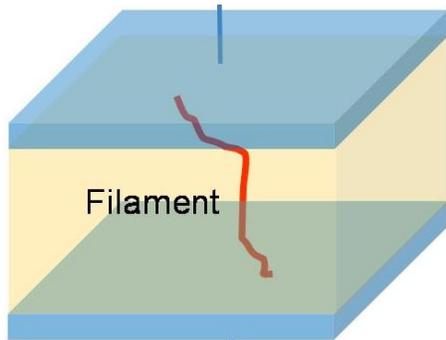
# Solving the Bottleneck: Processing Data Within Memory

**Impossible with current  
memory technologies!**

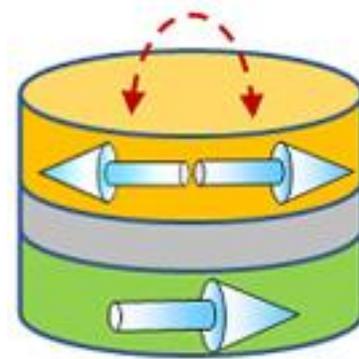


# Solution: Novel Memory Technologies

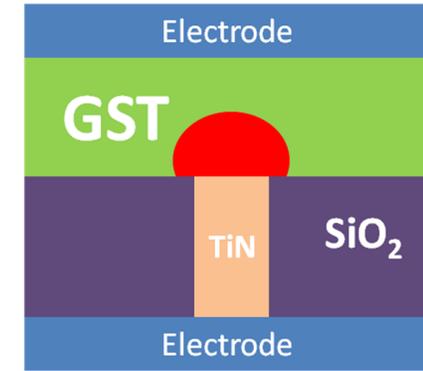
## Resistive RAM (RRAM)



## STT MRAM



## Phase Change Memory (PCM)



## Industry investment (partial list)

**TOSHIBA**

**winbond**



**Micron**

**IBM**



**QUALCOMM**

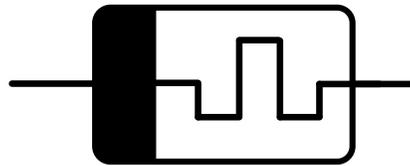
**EVERSPIN TECHNOLOGIES**



# Processing within Memristive Memories

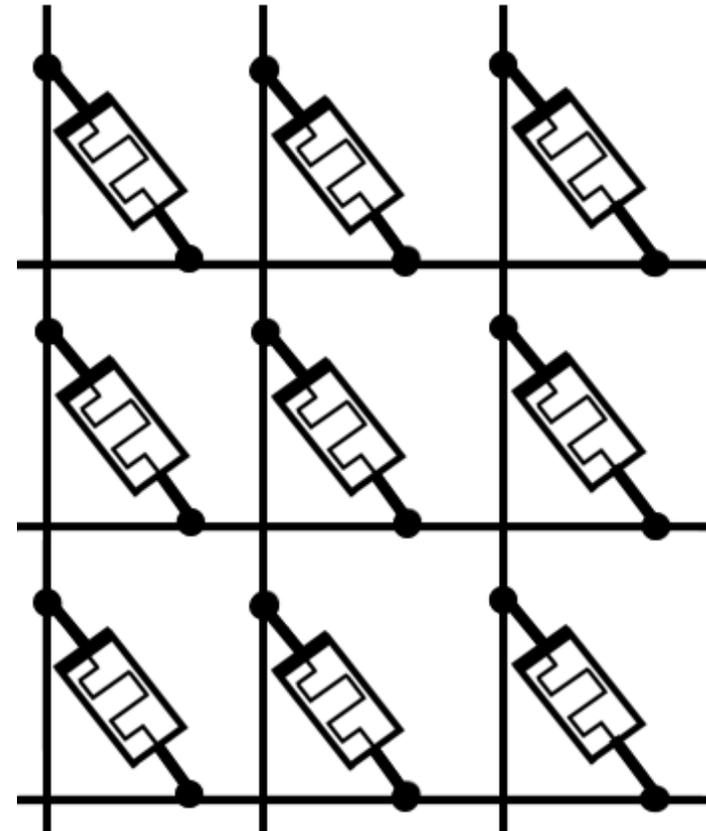


## Basic memory cell

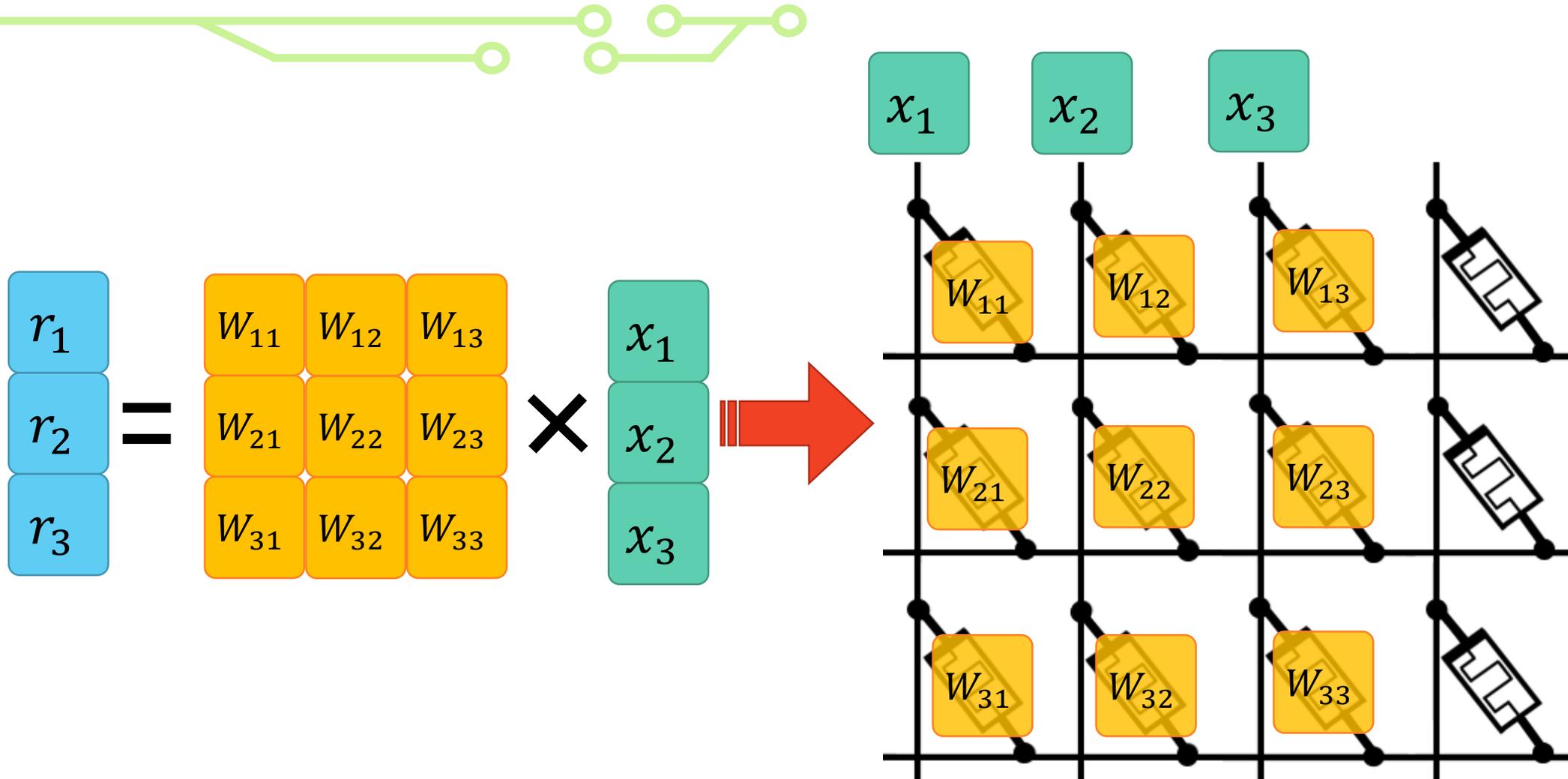


Memory resistor  $\Rightarrow$  **Memristor**

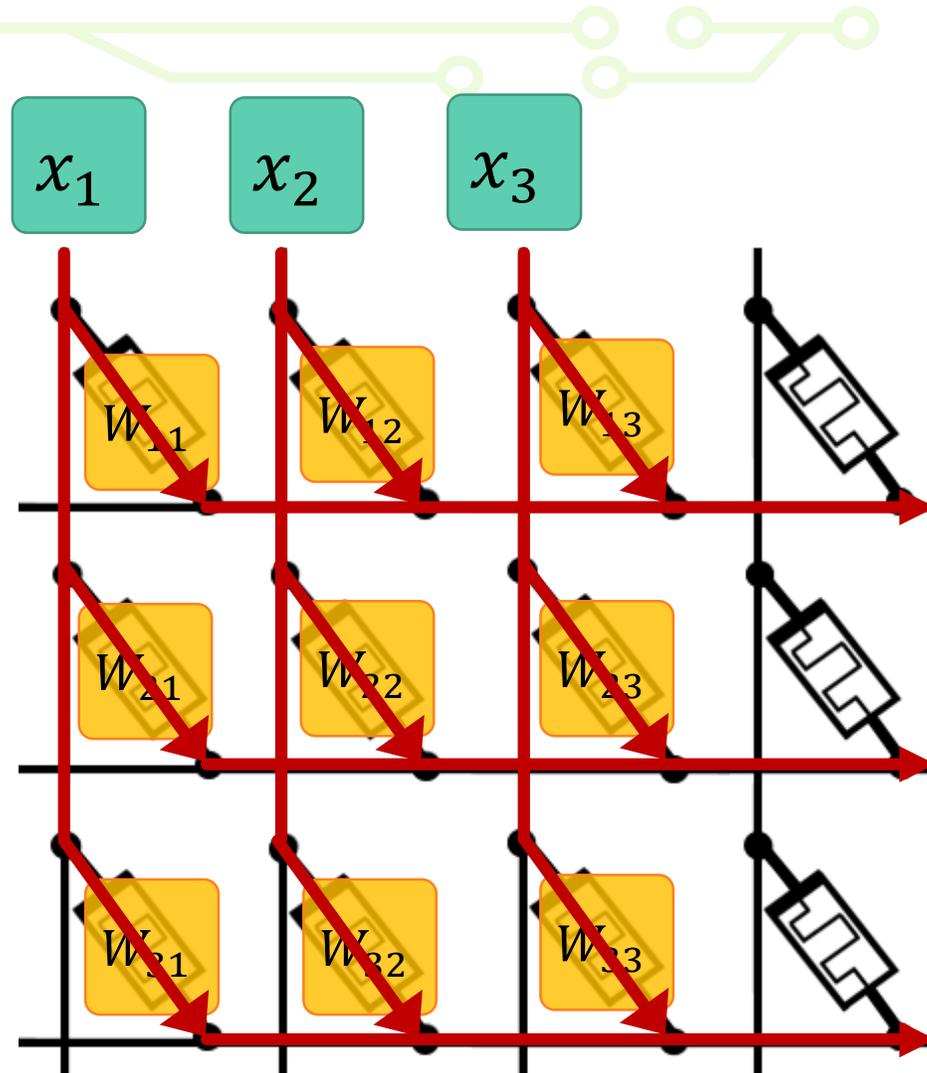
## Memristive memory



# Process in Memory for DNN



# Process in Memory for DNN



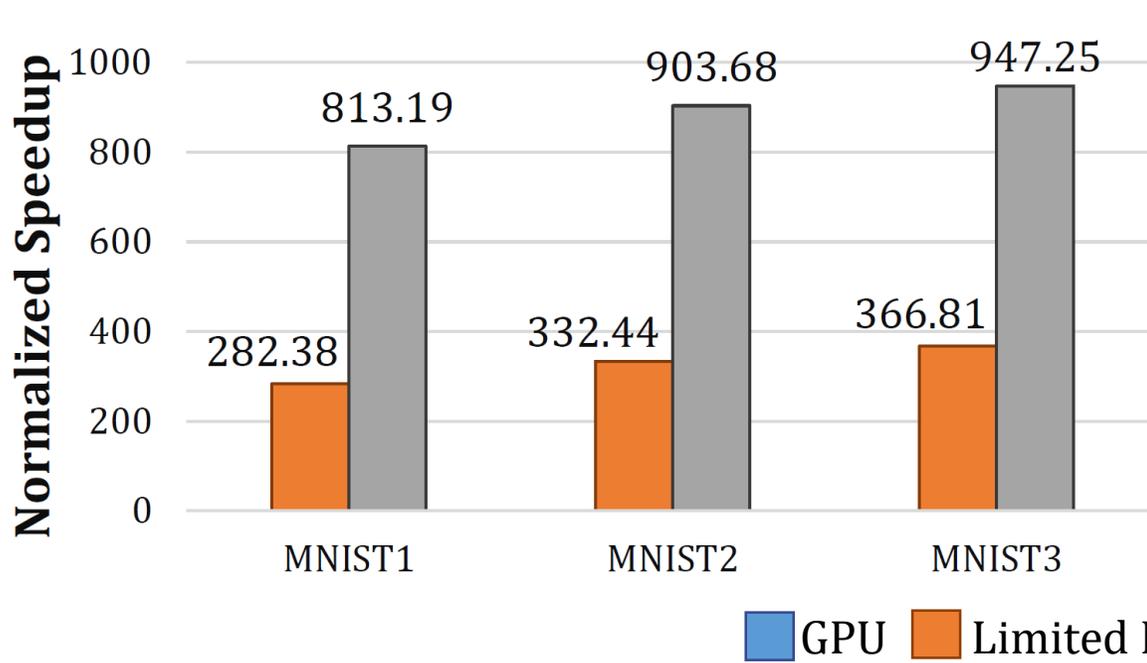
- Reduce data movement
- Highly parallel operation
- Immediate results
- Low power consumption

# Example: Support SGD + Momentum

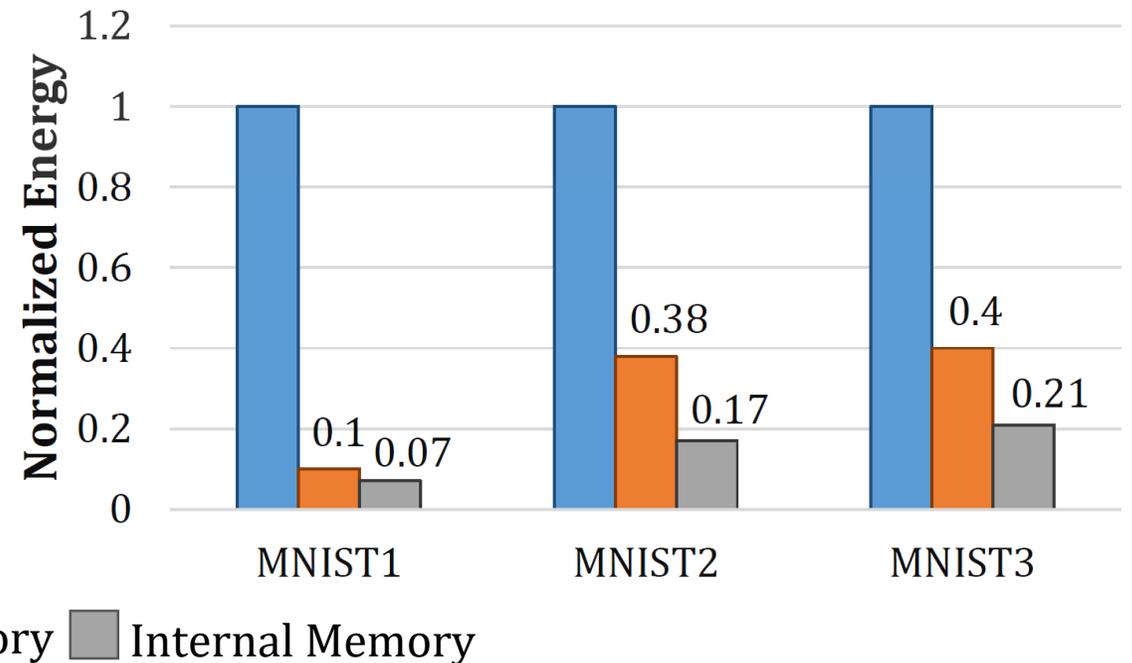
**~1000 × Speedup**

**80% Energy reduction**

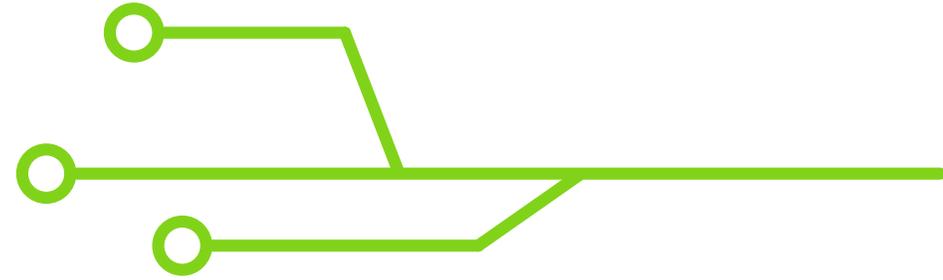
## Speedup Vs. GPU



## Energy consumption Vs. GPU



# Thanks!



**ASIC<sup>2</sup>** ARCHITECTURES  
SYSTEMS  
INTELLIGENT COMPUTING  
INTEGRATED CIRCUITS