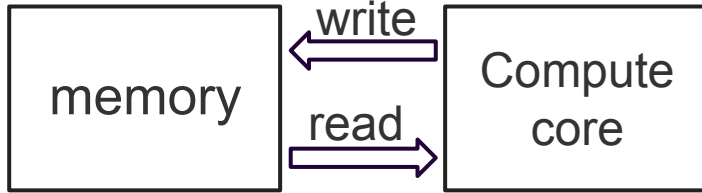# Lecture 13:
# New Computation Paradigms

# Recap

- FovealNet: Advancing AI-Driven Gaze Tracking Solutions for Efficient Foveated Rendering in Virtual Reality
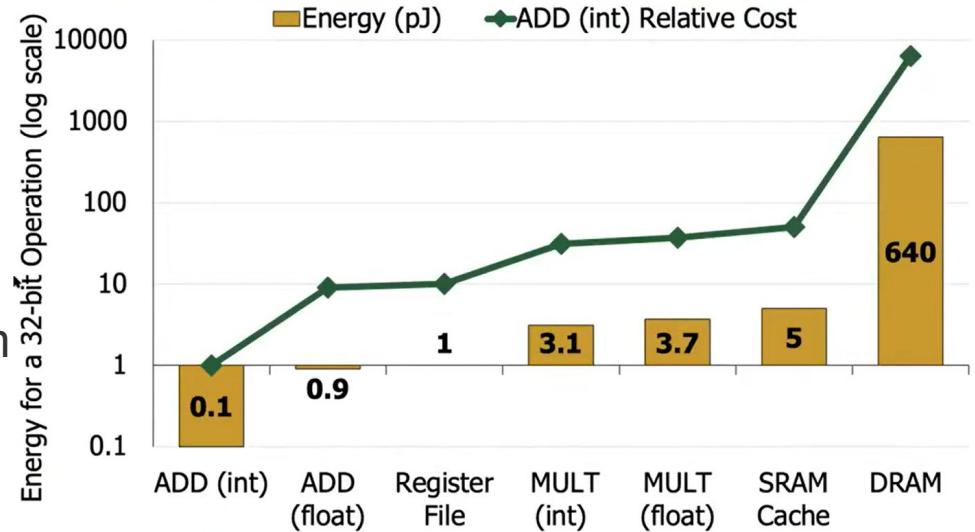- FovealSeg: Efficient Gaze-driven Instance Segmentation for Augmented Reality

# Topics

- In-memory computing
- Stochastic computing

NYU SAI LAB

# Data Movement Cost vs. Computation Cost



write
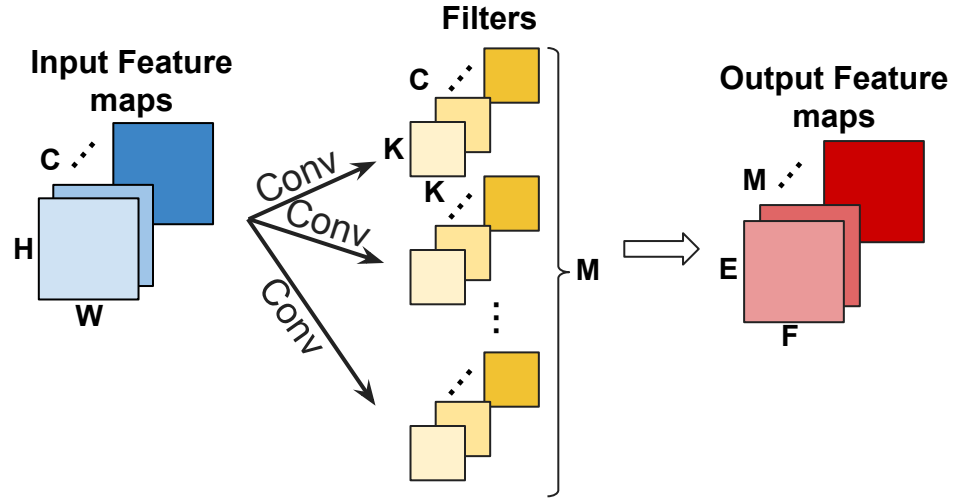
read

memory

Compute core

- Retrieving a single element from memory is more costly than computing it.
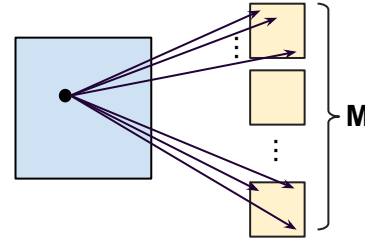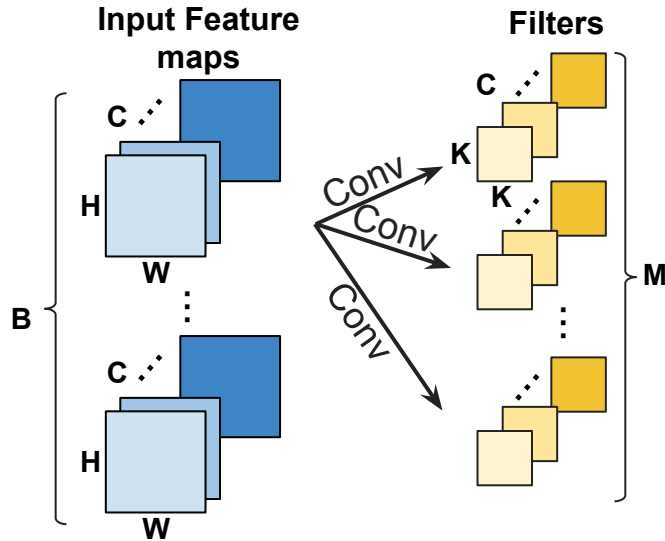
# Data Movement Cost vs. Computation Cost

- **Arithmetic intensity**: the ratio of total floating-point operations to total data movement (bytes)

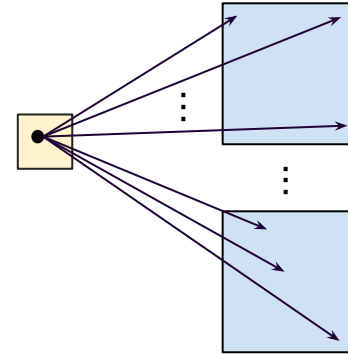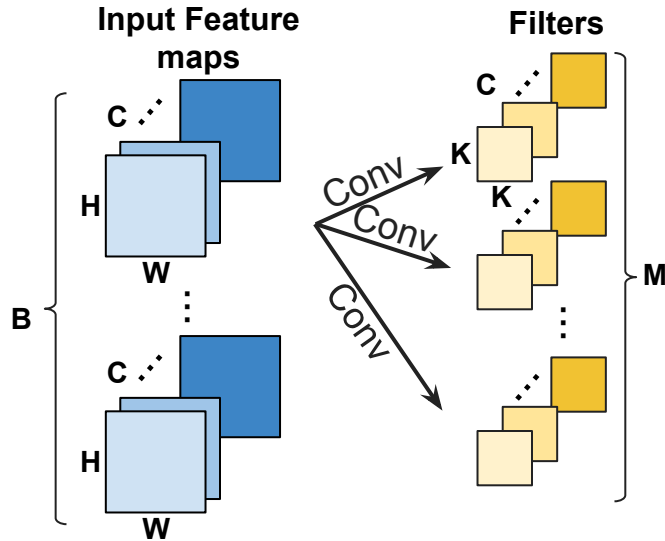$$\frac{\text{Total FLOPs}}{\text{Total data movement}}$$
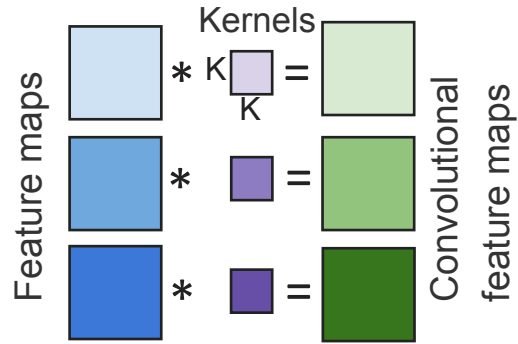
# Data Movement Cost vs. Computation Cost



- For each single element within the input feature maps, the maximum amount of reuse = $K^2M$.
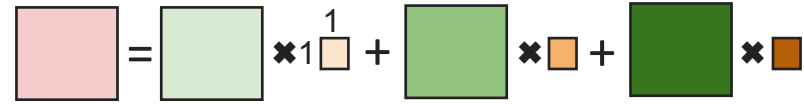
# Data Movement Cost vs. Computation Cost



- For each single element within the weight kernel, the maximum amount of reuse = BHW.
- For standard convolution, the arithmetic intensity is high.

# Data Movement Cost vs. Computation Cost

Kernels

Feature maps $*$ K $\square$ K $=$ 

Convolutional feature maps

$*$ $=$

$*$ $=$

**Step 1 Depthwise Convolution**

$=$ $\times 1$ $\overset{1}{\square}$ $+$ $\times \square$ $+$ $\times \square$

**Step 2 Pointwise Convolution**

Conv

Dconv

- For each single element within the input activation, the maximum amount of reuse = $K^2$ for Dconv.

# Breakdown on Computational Cost

**Latency Breakdown**

Matmul | Normalization | Softmax | Others
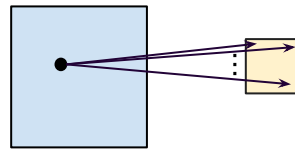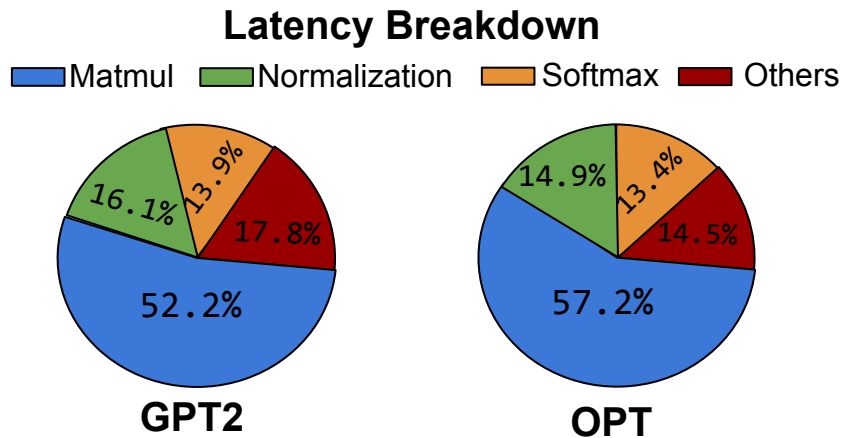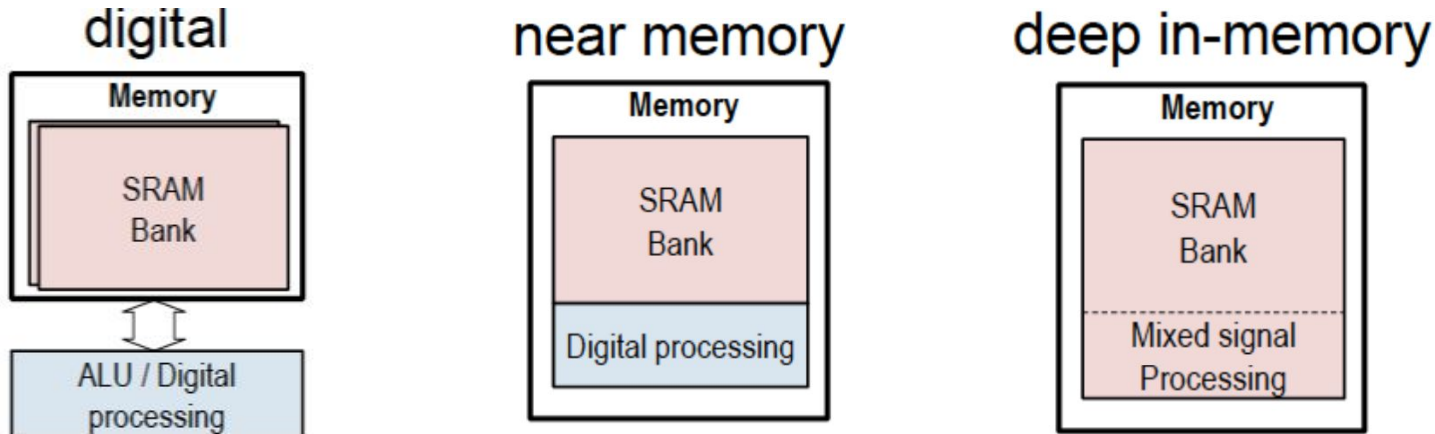


GPT2: 52.2% Matmul, 16.1% Normalization, 13.9% Softmax, 17.8% Others

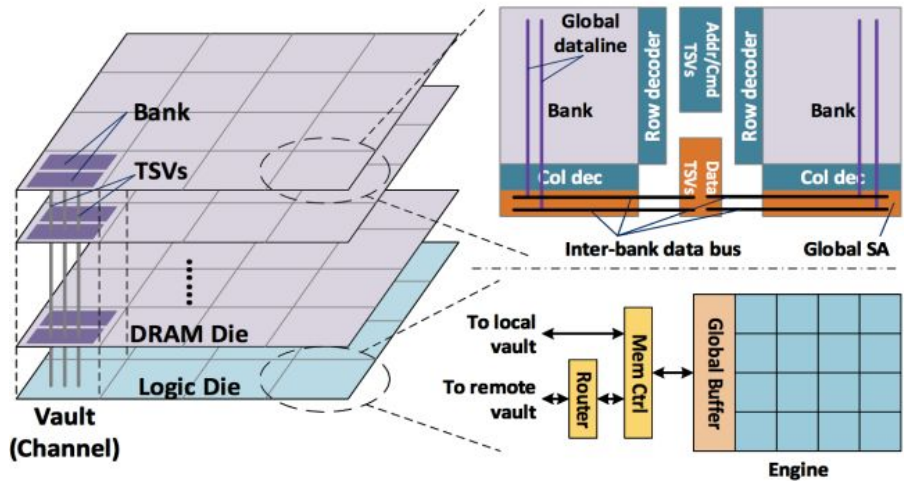OPT: 57.2% Matmul, 14.9% Normalization, 13.4% Softmax, 14.5% Others

- Matmul still contributes to majority of the overall latency.
- Nonlinear operations are not negligible.
- Also other operations (e.g., transposition, reshape) also contributes to a great portion of the overall latency.
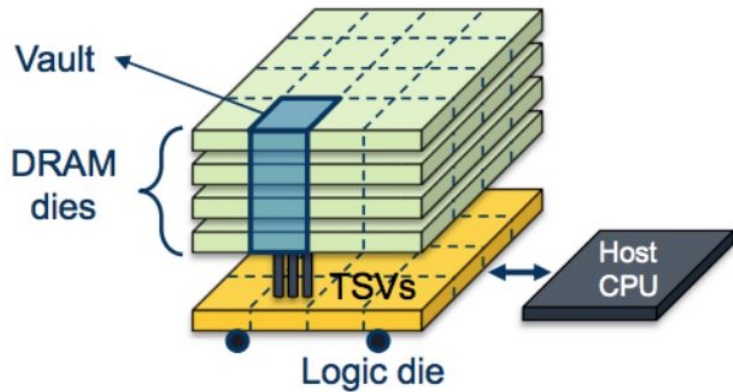
# Near/In-Memory Processing



- Near memory computing has a higher BW, and analog in-memory computing integrate the computation with the memory access.
- Analog PIM brings compute closer to the memory.

# Near Memory Processing



**Tetris**

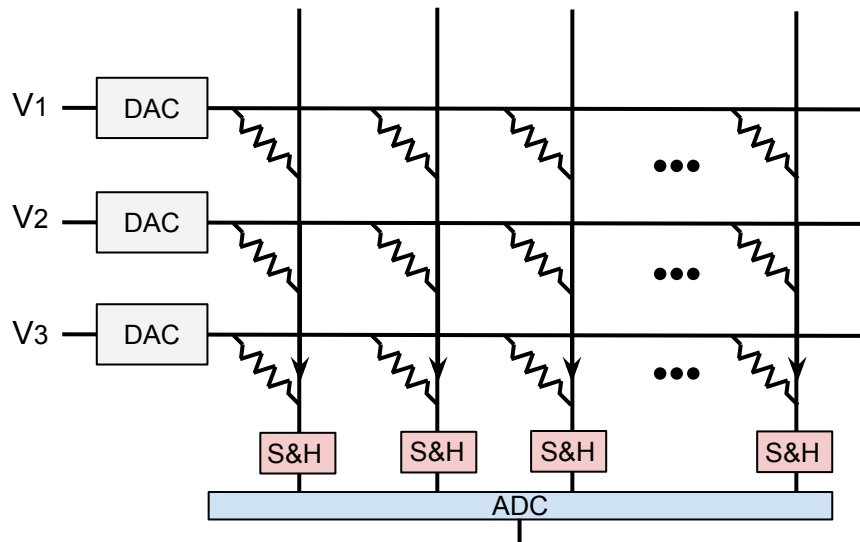**NeuroCube**

Gao, Mingyu, et al. "Tetris: Scalable and efficient neural network acceleration with 3d memory." *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*. 2017.
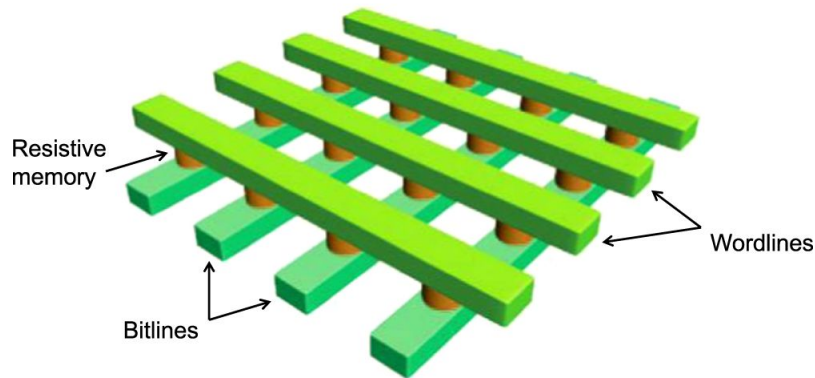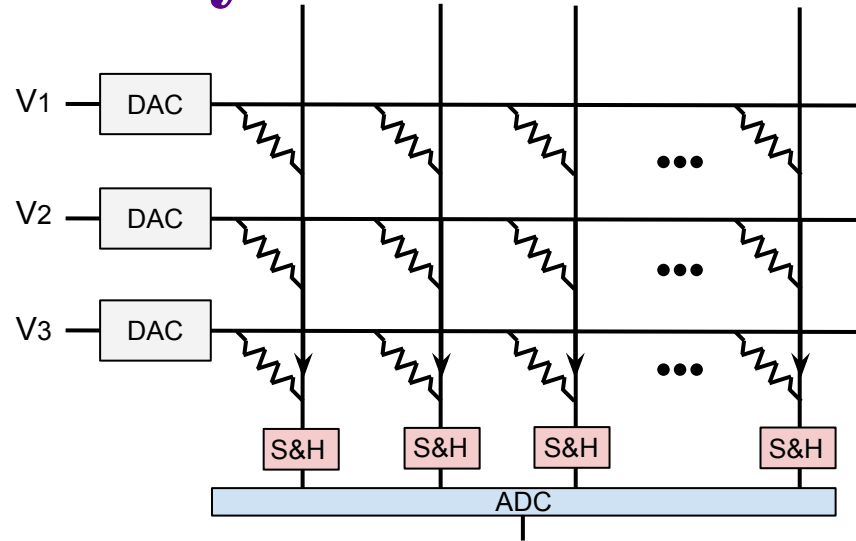
# Resistive Memory



- Resistive RAM (ReRAM or RRAM) is a type of non-volatile RAM that works by changing the resistance across a dielectric solid-state material, often referred to as a memristor.



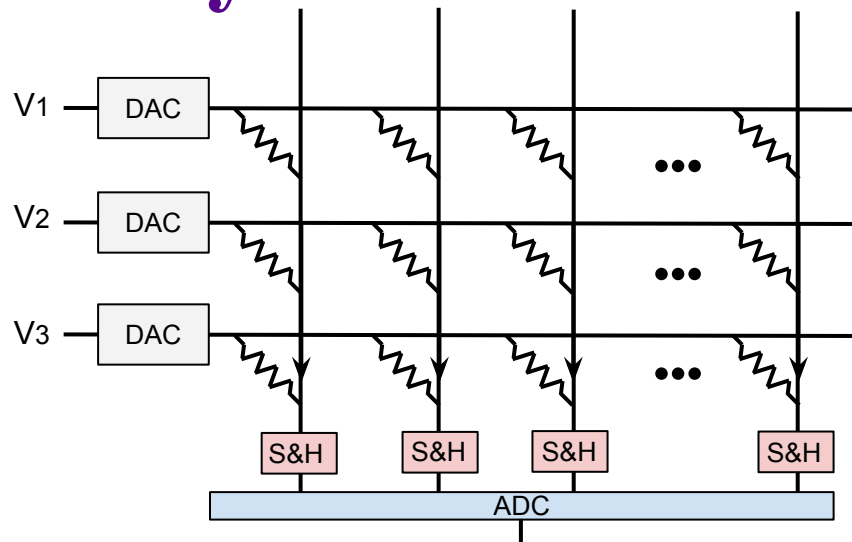Resistive memory

Wordlines

Bitlines

# Processing In Memory

$$I = V_1G_1 + V_2G_2$$



- The digital input are first passed to the DAC and converted to the analog input voltages.
- The voltages are applied to each of the rows in the crossbar array.

# Processing In Memory

$$I = V_1G_1 + V_2G_2$$



- The output current accumulated at the bottom of each column is the dot product between the voltages and the conductances across the rows.
- A sample-and-hold (S&H) circuit receives the bitline current and feeds it to a shared ADC unit

14

# Processing In Memory

**Original**

V1= 3
(11)

DAC

G1= 1

V2= 2
(10)

DAC

G2= 3

S&H

ADC

Result = 9

T = 1

1

DAC

G1= 1

0

DAC

G2= 3

S&H

ADC

Result = 1

T = 2

1

DAC

G1= 1

1

DAC

G2= 3

S&H

ADC
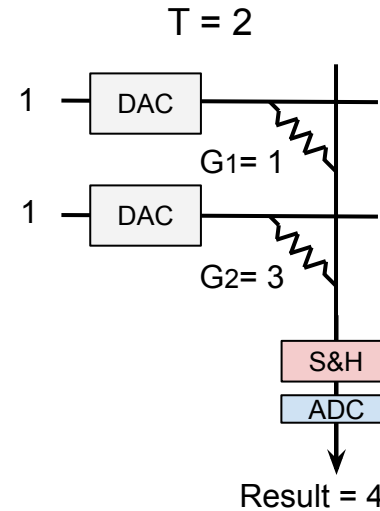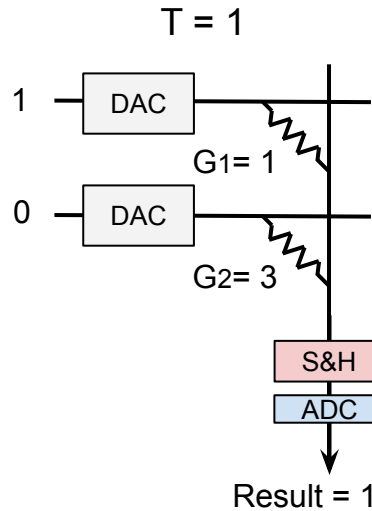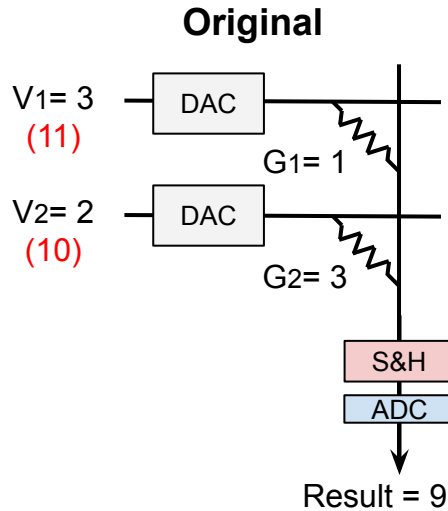
Result = 4

$4 \times 2^1 + 1 \times 2^0 = 9$

- Assume both inputs and weights are 16 bits, we need a 16-bit DAC to provide input voltage, $2^{16}$ resistance levels in each cell, and an ADC which can handle over 16 bits, which leads to a significant overhead.

NYU SAI LAB

# Processing In Memory



**Original**

$V_1 = 3$
(11)

DAC

$G_1 = 1$

$V_2 = 2$
(10)

DAC

$G_2 = 3$

S&H

ADC

Result = 9

**T = 1**

1 — DAC

$G_1 = 1$

0 — DAC

$G_2 = 3$

S&H

ADC

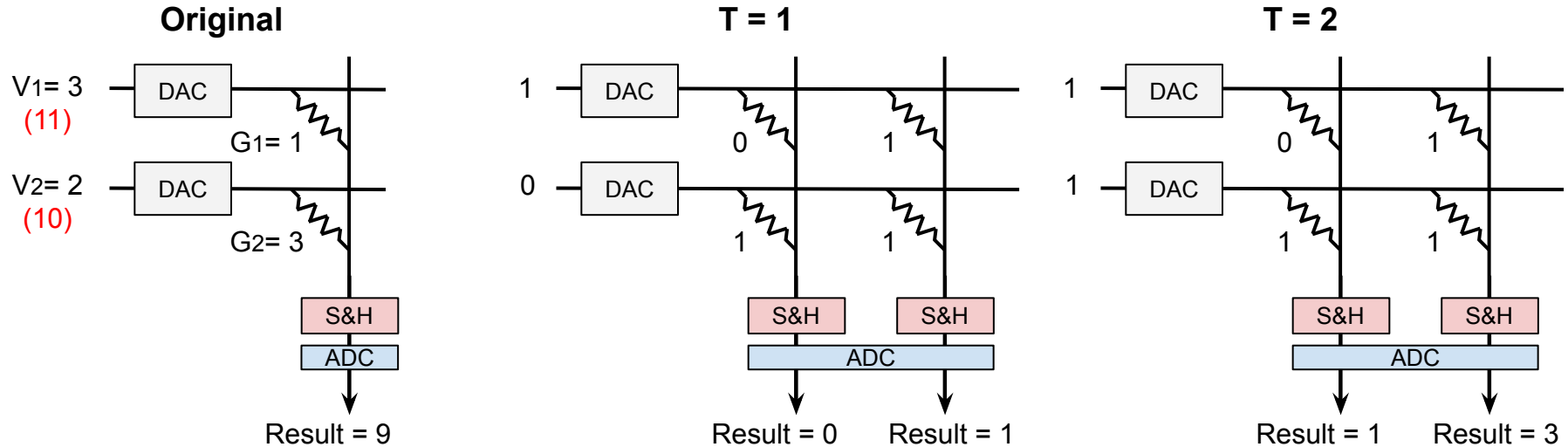Result = 1

**T = 2**

1 — DAC

$G_1 = 1$

1 — DAC

$G_2 = 3$

$4 \times 2^1 + 1 \times 2^0 = 9$

S&H

ADC

Result = 4

- Instead, the digital input enters the crossbar in a bit-serial manner, the intermediate results are buffered in the register. Shift-Add operation is them performed after all the input bits entering the crossbar.

# Processing In Memory



**Original**

$V_1 = 3$
(11)

DAC

$G_1 = 1$

$V_2 = 2$
(10)

DAC

$G_2 = 3$

S&H

ADC

Result = 9

**T = 1**

1 — DAC

0        1

0 — DAC

1        1

S&H        S&H

ADC

Result = 0        Result = 1

**T = 2**

1 — DAC

0        1

1 — DAC

1        1
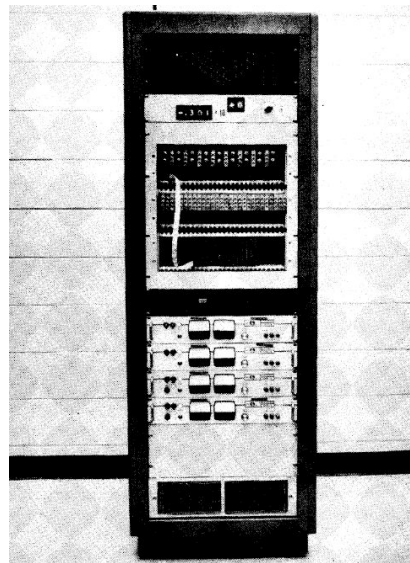
S&H        S&H

ADC

Result = 1        Result = 3

$1 \times 2^0 + 0 \times 2^1 + 1 \times 2^1 + 3 \times 2^1 = 9$

NYU SAI LAB

# Topics

- Processing in memory
- Stochastic Computing
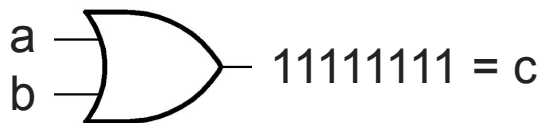
# Stochastic Computing

- Stochastic computing is a computational approach that utilizes random bit streams to perform numerical calculations, offering benefits in power efficiency and hardware simplicity, particularly for error-tolerant applications.
- Introduced by John von Neumann in 1953.



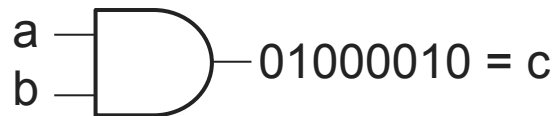The RASCEL stochastic computer, circa 1969

# Stochastic Computing

- a = 0.5, b = 0.5
  - a = 00111100  $p_a(1) = 0.5$
  - b = 11000011  $p_b(1) = 0.5$

$$a \lor b \rightarrow 11111111 = c$$

$$p_c(1) = 1$$

- a = 0.5, b = 0.5
  - a = 11001010  $p_a(1) = 0.5$
  - b = 01010011  $p_b(1) = 0.5$

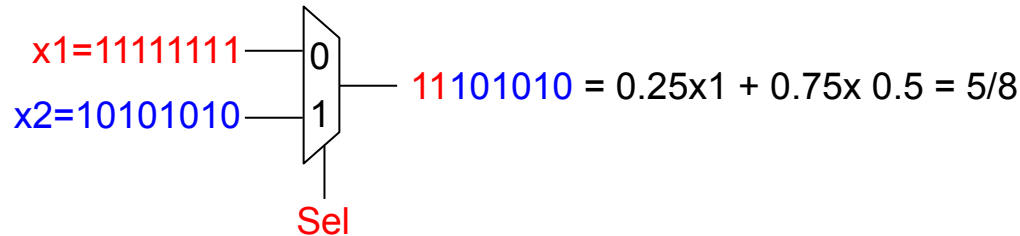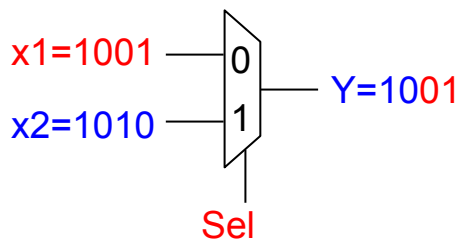$$a \land b \rightarrow 01000010 = c$$

$$p_c(1) = 0.25$$

- As the input stream lengthens, the multiplication process will become more accurate.
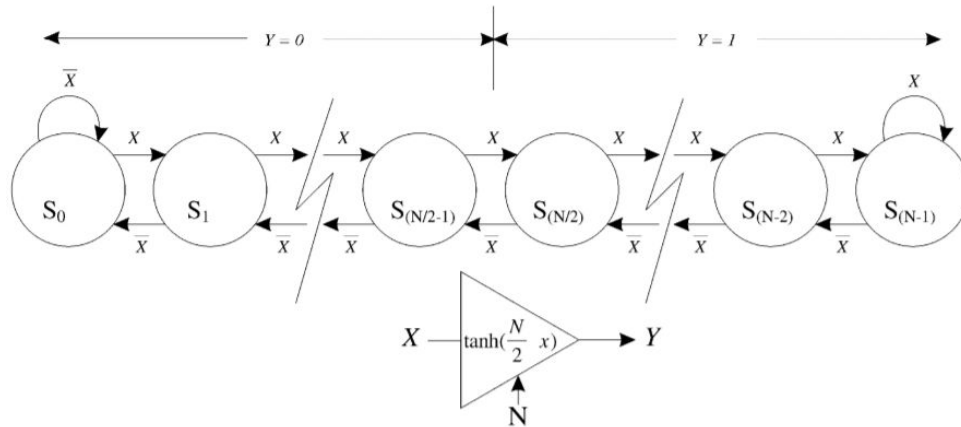
NYU SAI LAB

# Addition with Stochastic Computing

- ## MUX implementation
  - By adjusting Sel over time, the output of the multiplexer will equal to the weighted sum of the input bit streams.
  - The accuracy gets worse when the number of inputs to the MUX is large.



x1=1001 — 0
x2=1010 — 1
Sel
Y=1001

x1=11111111 — 0
x2=10101010 — 1
Sel
11101010 = 0.25x1 + 0.75x 0.5 = 5/8

# Nonlinear Operation with Stochastic Computing

- The tanh function is highly suitable for SC-based implementations because i) it can be easily implemented with a K-state finite state machine (FSM) in the SC domain.



- The major advantage of stochastic computing is the significantly lower hardware cost for a large category of arithmetic calculations.