



COSIME: FeFET based Associative Memory for In-Memory Cosine Similarity Search

Che-Kai Liu^{1,2,3}, Haobang Chen¹, Mohsen Imani³, Kai Ni⁴, Arman Kazemi², Ann Franchesca Laguna⁵, Michael Niemier², Xiaobo Sharon Hu², Liang Zhao¹, Cheng Zhuo¹, and Xunzhao Yin¹

¹Zhejiang University, Hangzhou, China, ²University of Notre Dame, USA
³University of California, Irvine, USA, ⁴Rochester Institute of Technology, USA,
⁵De La Salle University, Manilla, Philippines

Von-Neumman Paradigm

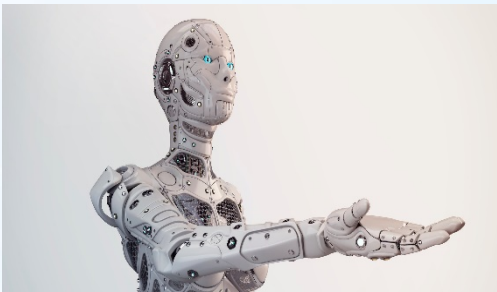
Self Driving Cars



Healthcare



Smart Robots



Finance



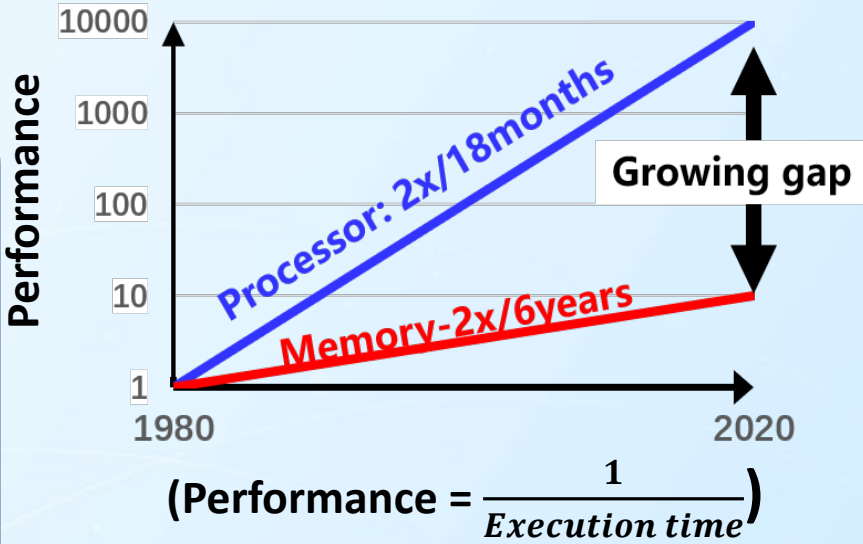
Machine Translation



Gaming



Memory scaling is low

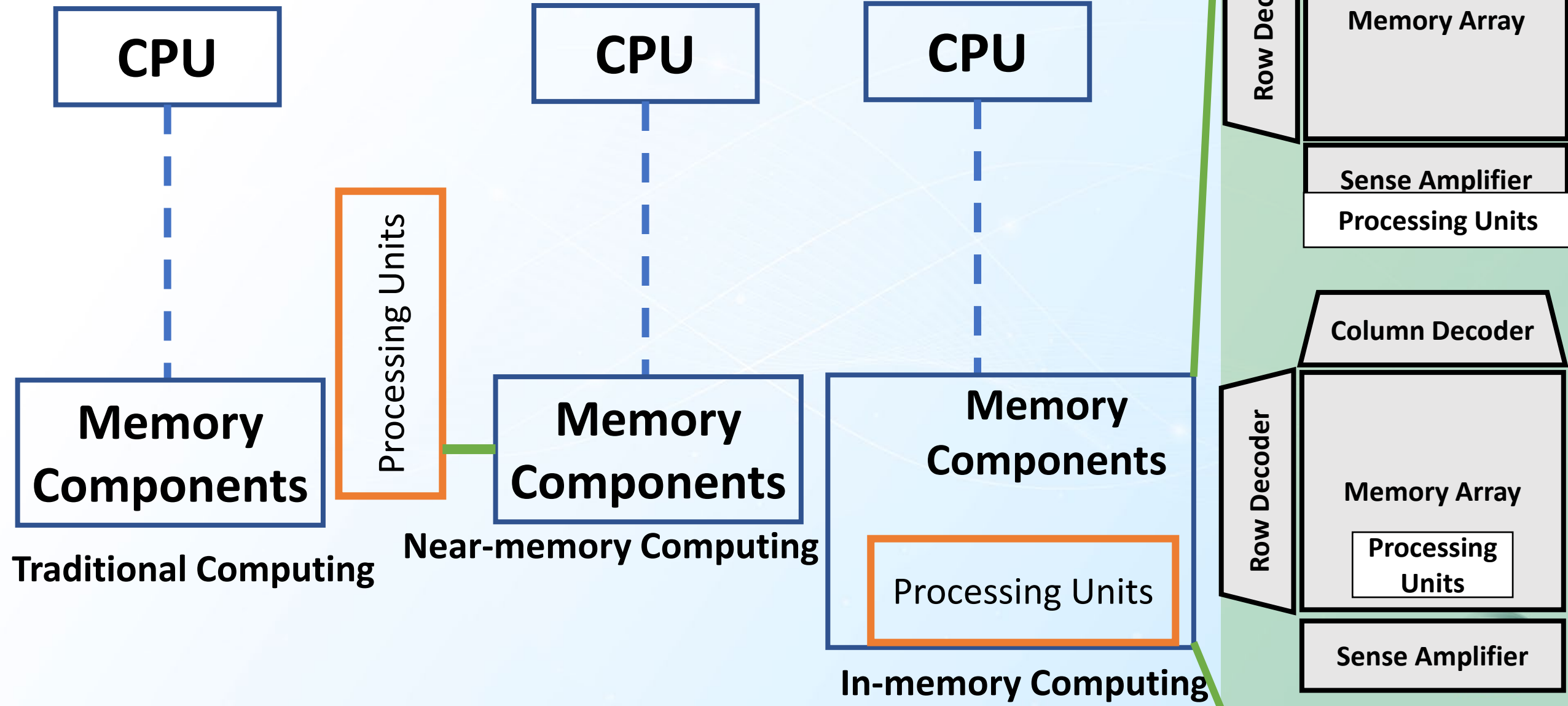


- Concerns
 - Latency
 - Energy
 - Security
 - ...

Operation	Energy (pJ)	Relative Cost
32-bit int ADD	0.1	1
32-bit int MULT	3.1	31
32-bit 32KB SRAM	5	50
32-bit DRAM	640	6400



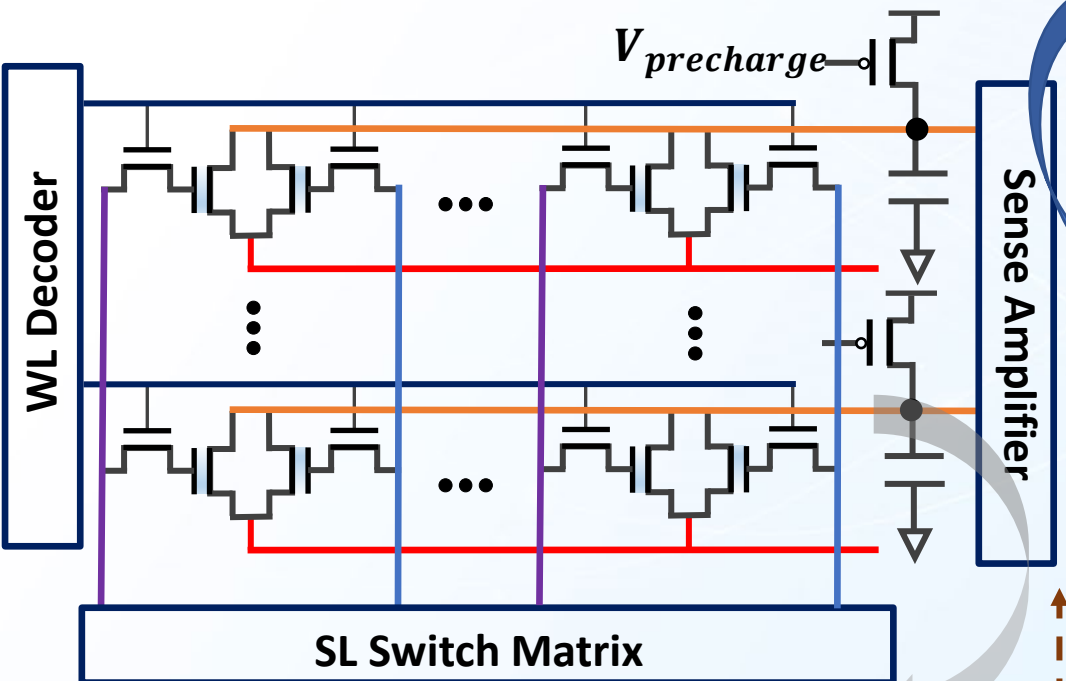
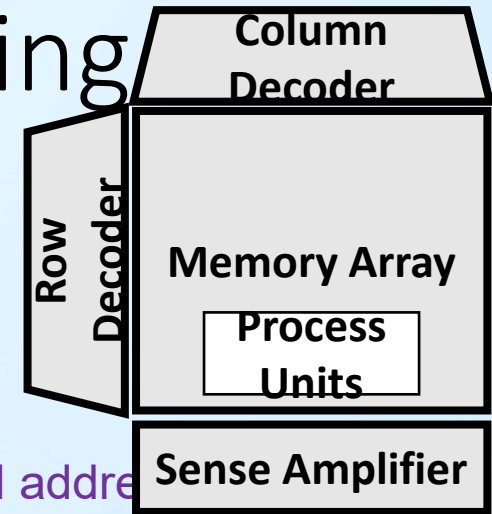
Computing Paradigms



Content Addressable Memory - Hamming

Step2: parallel search
(e.g., 192.168.1.1)

0 1 1 0 0



TCAM address

Don't care

RAM address

00	1	0	1	X	X
01	0	1	1	0	X
10	0	1	1	X	X
11	1	0	0	1	1

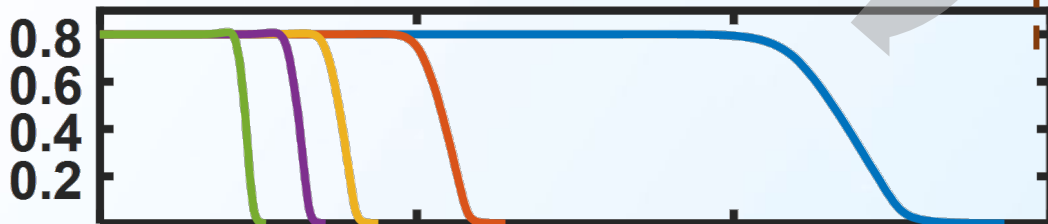
01
Search result

00	Port = A
01	Port = B
10	Port = C
11	Port = D

Output
Port=B

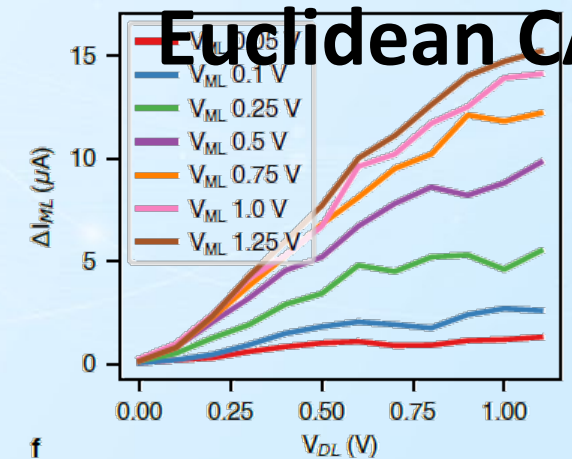
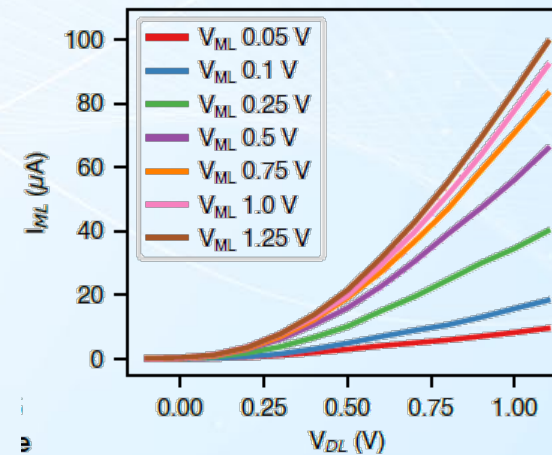
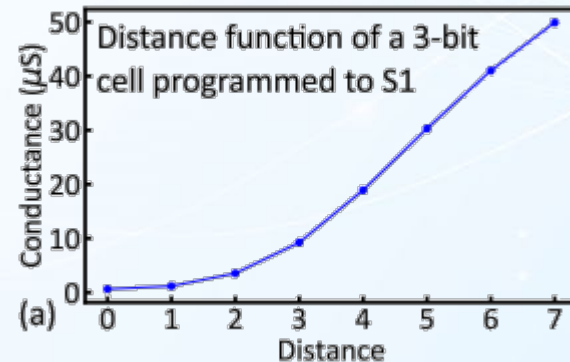
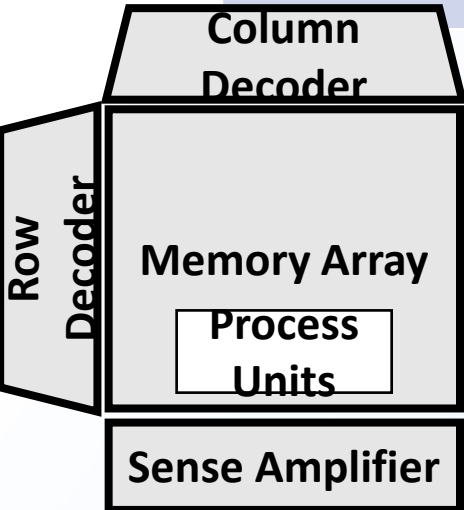
Step1: write 1/0/X(don't care)

Hamming distance!



Content Addressable Memory – Sigmoid/L2

Region	Conditions	Expression	Distance
Cut-off	$V_{gs} \leq V_{th}$	$I_{ds} = 0$	NA
Linear (short)	$V_{gs} > V_{th}, V_{ds} \leq V_{gs} - V_{th}$	$I_{ds} \propto (V_{gs} - V_{th}) V_{ds} - \frac{V_{ds}^2}{2}$	Sigmoid [1]
Saturation (long)	$V_{gs} > V_{th}, V_{ds} > V_{gs} - V_{th}$	$I_{ds} \propto (V_{gs} - V_{th})^2$	Squared Euclidean [2]

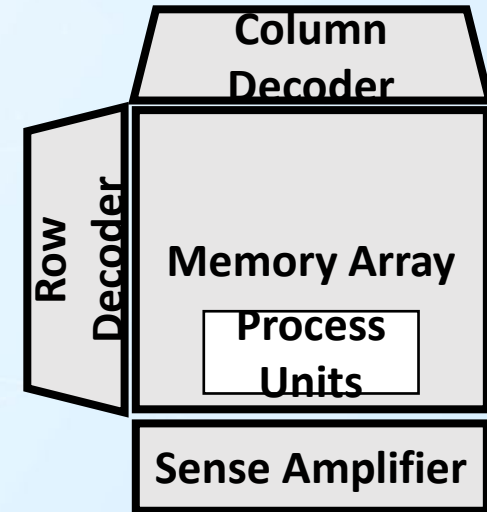
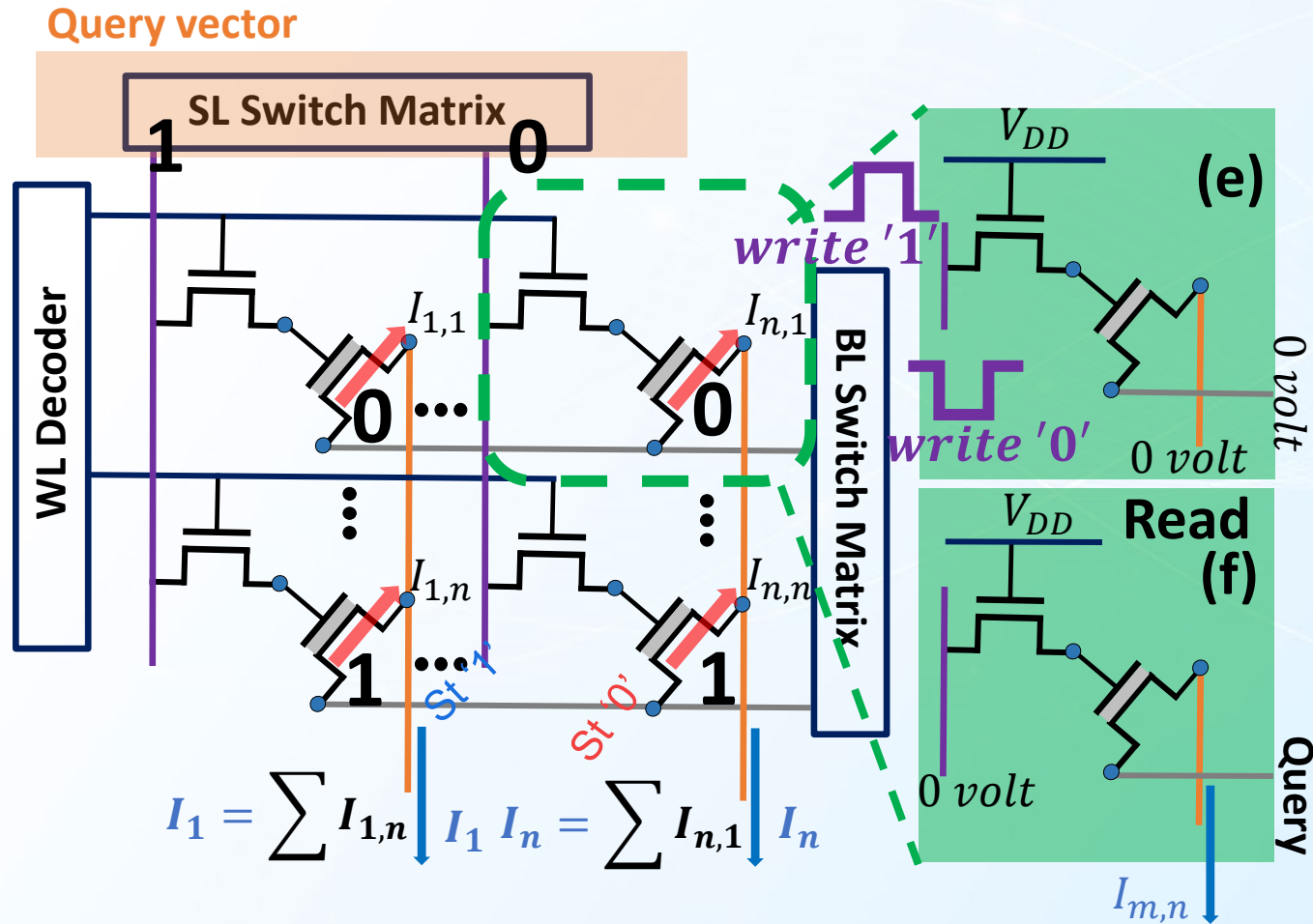


Euclidean CAM!

[1] FeFET Multi-Bit Content-Addressable Memories for In-Memory Nearest Neighbor Search, A. Kazemi et al., IEEE TC

[2] Software-equivalent accuracy in-memory hyperdimensional computing with ferroelectric devices, A. Kazemi et al., To appear in Scientific Report

Crossbar



Inner product!

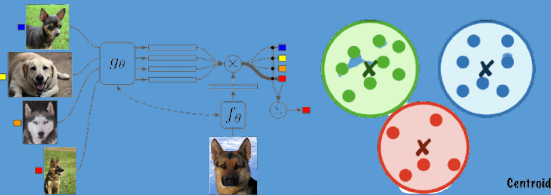
Where is Cosine?

A cross-layer design approach

Bottom-up

Hyperdimensional Computing (HDC)

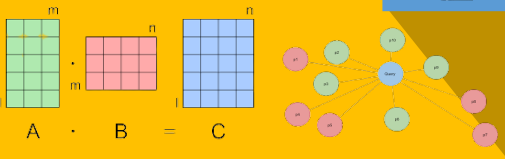
Applications (HDC, few-shot learning, clustering, etc.)



Application driven

Search

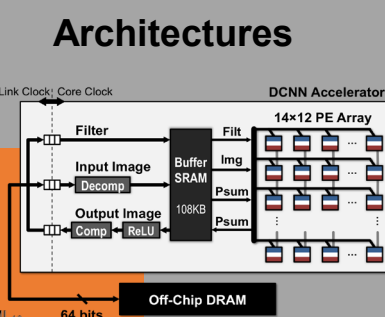
Algorithms (VMMs, search, etc.)



Search is a bottleneck

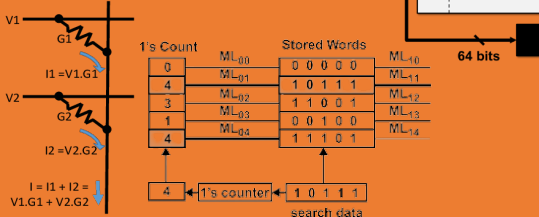
xBar/CAM-based architectures

Architectures



xBars, CAMs

Circuits

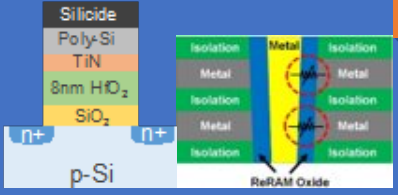


Architectures for long vectors

Nearest Neighbor Search accuracy

FeFETs

Devices (FeFETs, RRAMs, etc.)



Device non-idealities

Top-down

Possible solutions – AM for cosine similarity

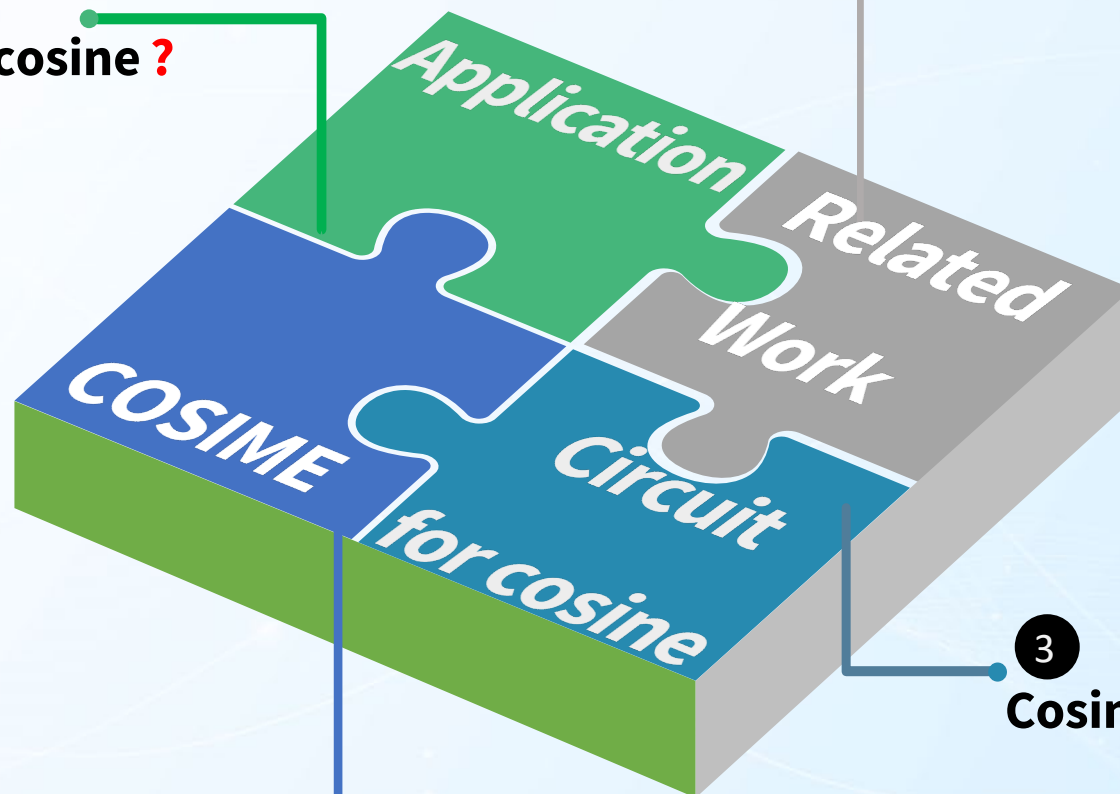
- X-MANN, DAC'19: $\frac{a \cdot b}{||a||_1 + ||b||_1}$
- Nature communication, 2021: Quasi-orthogonal of hyper vectors
- Normalize the data first then cosine? GPU takes millisecond
- Bipolar vector eliminating division?

1

- CAM-friendly: Hamming ✓
- Few-shot: L2 on CAM ✓
- HDC, text-related app: cosine ?

2

Approximation of cosine
(related work)



3

Cosine-friendly analog circuit design

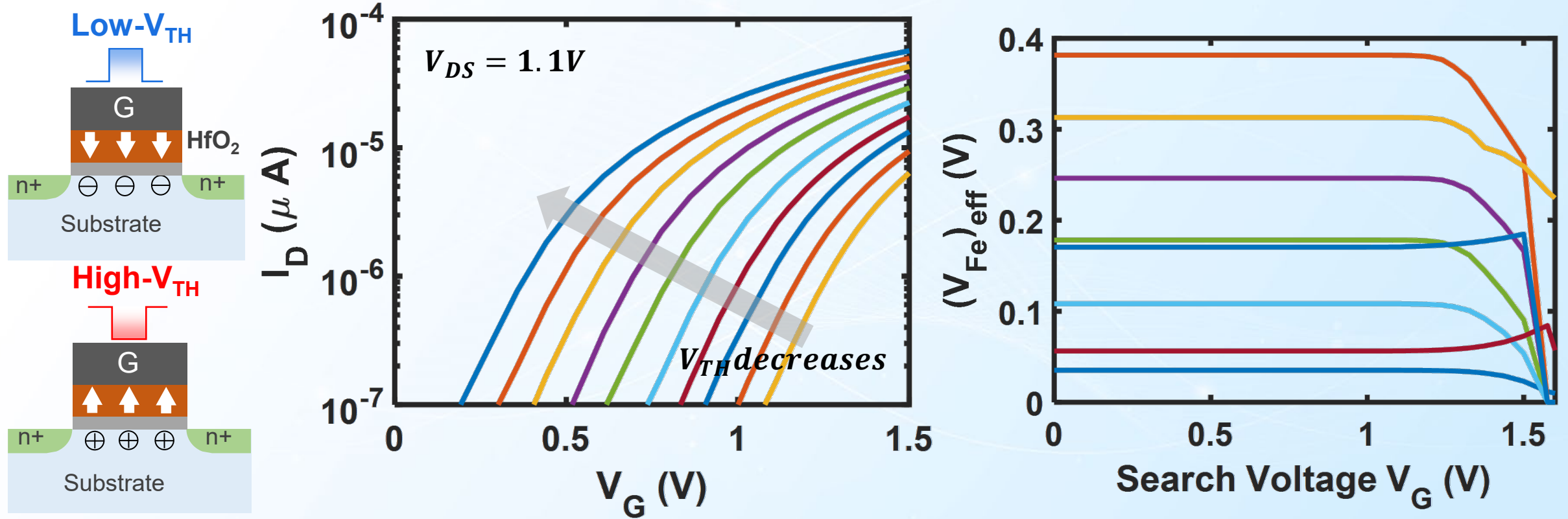
4

Direct implementation of
cosine similarity (this work, COSIME)

$$\frac{a \cdot b}{||a||_2 \cdot ||b||_2} \rightarrow \frac{(a \cdot b)^2}{(\cancel{||a||_2} \cdot ||b||_2)^2} \rightarrow \frac{X^2}{Y}$$

L_2 norm of binary vector is counting "1"s

Ferroelectric field-effect transistor (FeFET)



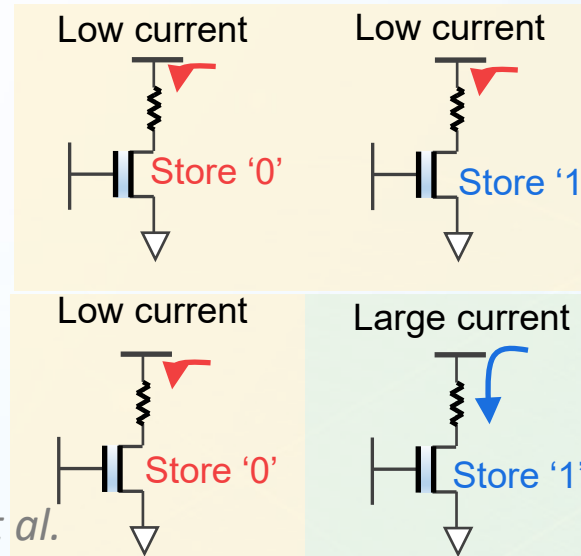
I thank Dr. X. Sharon Hu from the University of Notre Dame for providing this 22nm FeFET model

Breaking down the design

1FeFET1R Cell

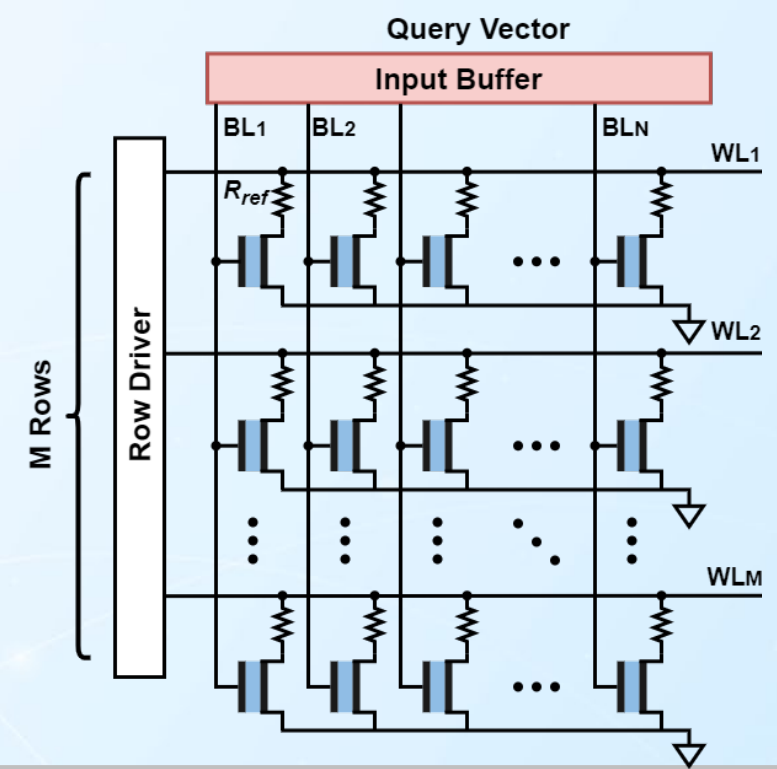
- Mitigate variations
- Scalable array design

$$\frac{(\frac{I_x}{N} \times N)^2}{\frac{I_y}{N} \times N} = I_z$$

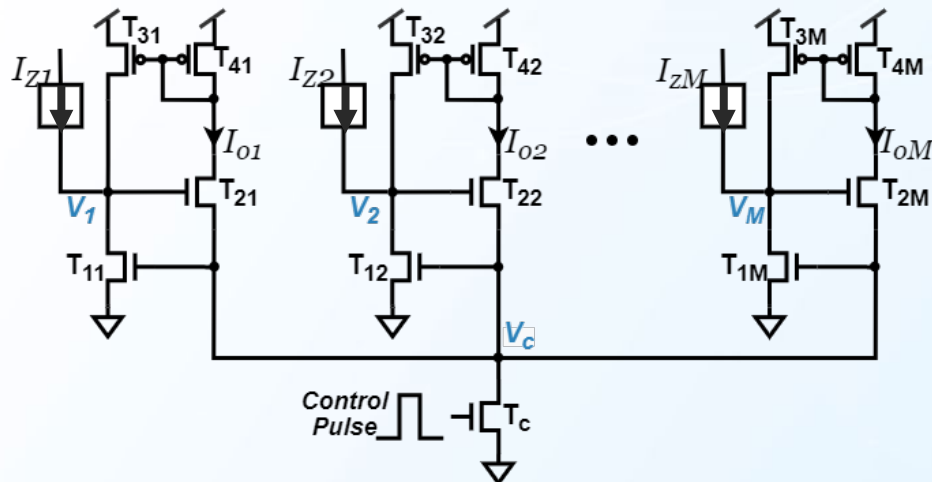


Process is demonstrated
in Symp. on VLSI, 2021, D. Saito et al.

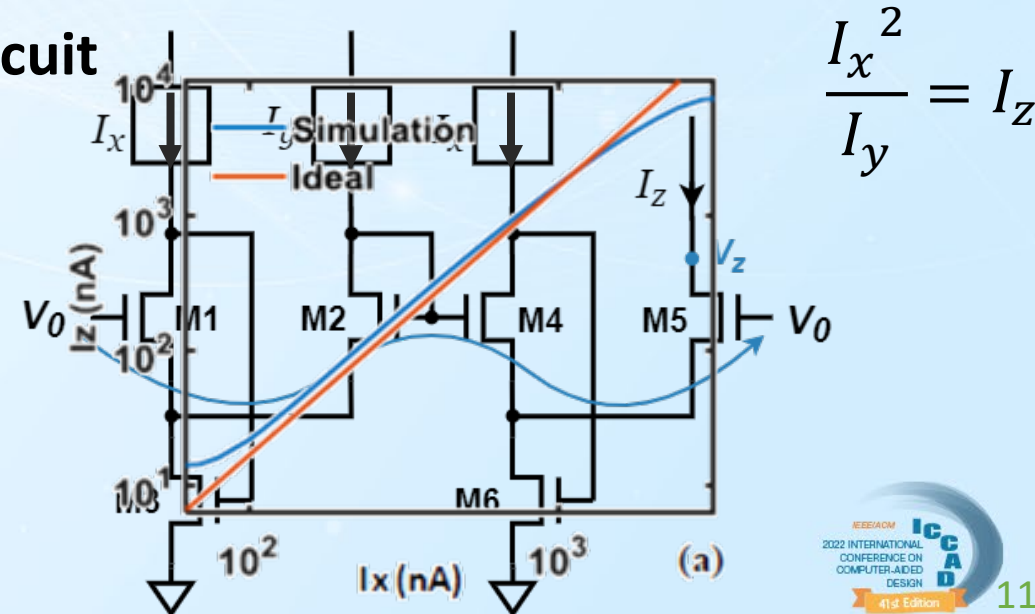
Crossbar

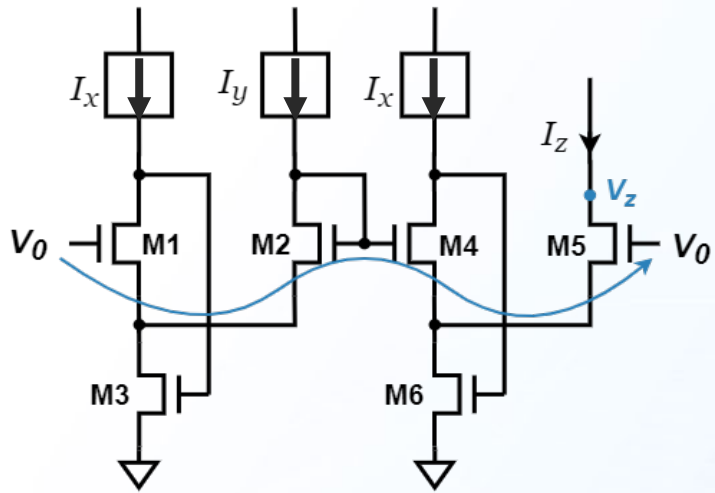


Winner-take-all circuit



$\frac{x^2}{y}$ circuit





1

$$I_{DS} \approx I_0 \frac{W}{L} e^{\frac{V_{GS}}{\eta T}}$$

2

$$V_{GS} = V_T \eta \ln\left(\frac{I_{DS} L}{I_0 W}\right)$$

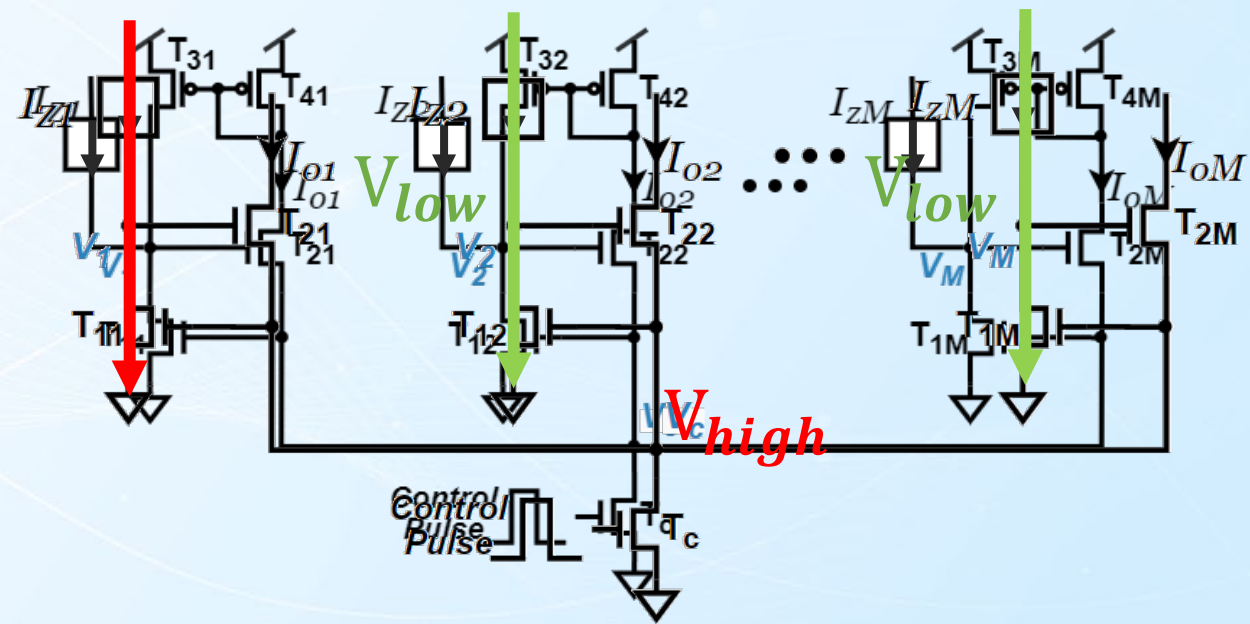
3

$$\sum_{\text{ClockWise}} V_{GS} = \sum_{\text{counter CW}} V_{GS}$$

(M1, M2, M4, M5)

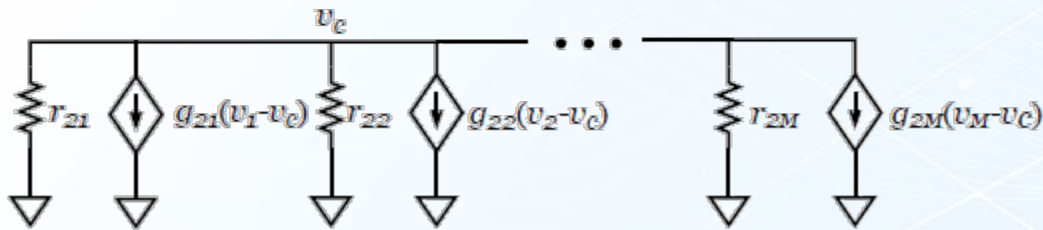
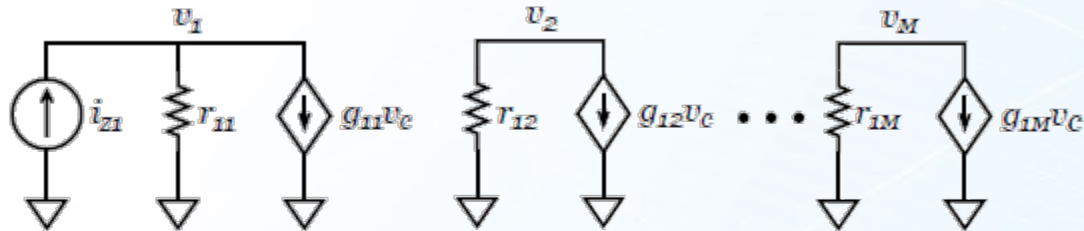
4

$$\frac{I_x^2}{I_y} = I_z$$

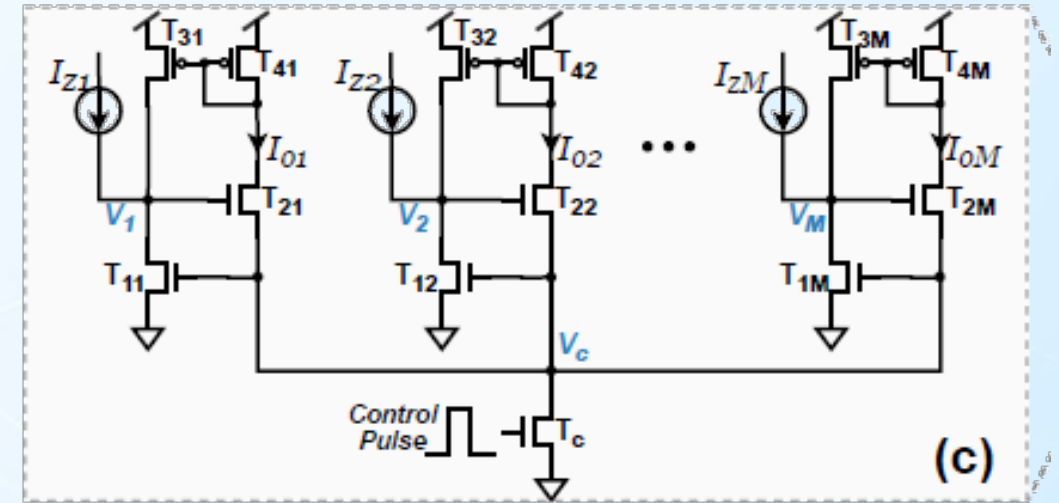


Scalability Analysis

$$g_{11} = \frac{I_{Z1}}{V_T} \text{ and } r_{11} = \frac{V_A}{I_{Z1}} \text{ etc.}$$



$$\begin{cases} i_{z1} = v_1 \frac{I_{Z1}}{V_A} + v_c \frac{I_{Z1}}{V_T} \\ v_j \frac{I_{Zj}}{V_A} = -v_c \frac{I_{Zj}}{V_T}, \quad \forall j \in [2, M] \\ \sum_{i=1}^M \left[\frac{I_{oi}}{V_T} (v_i - v_c) + v_c \frac{I_{oi}}{V_A} \right] \end{cases}$$



$$\frac{dV_1}{dI_{Z1}} = \frac{M-1}{M} \frac{V_A}{I_{Z1}}$$

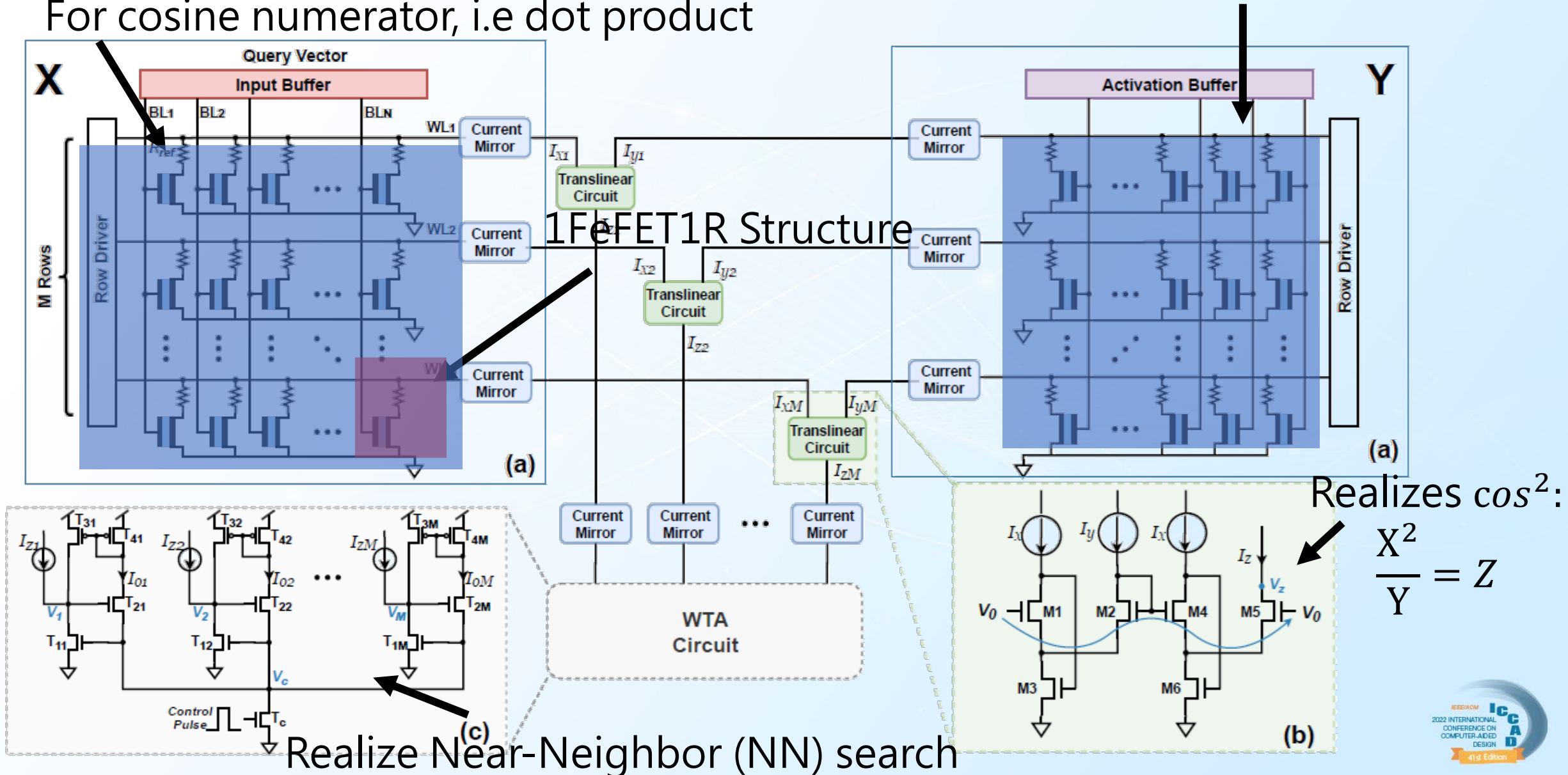
$$\frac{dV_2}{dI_{Z1}} = \frac{-1}{M} \frac{V_A}{I_{Z1}}$$

Overall Design of "COSIME"

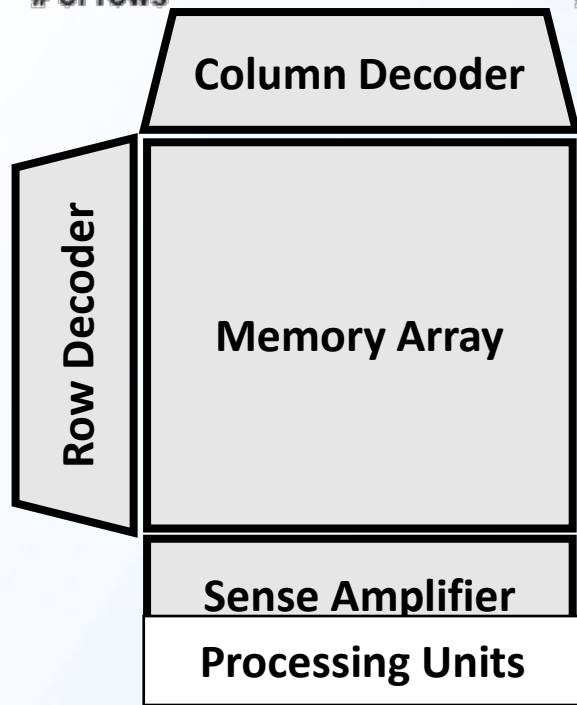
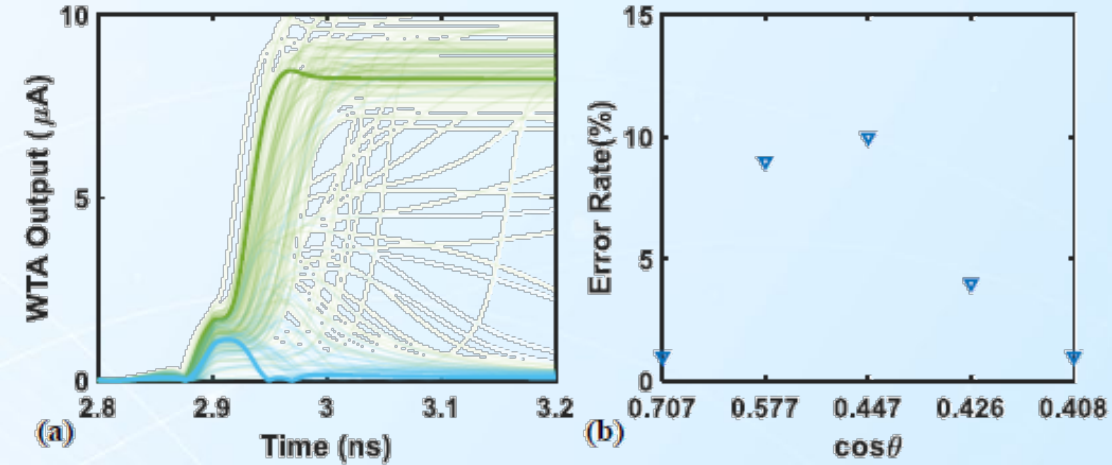
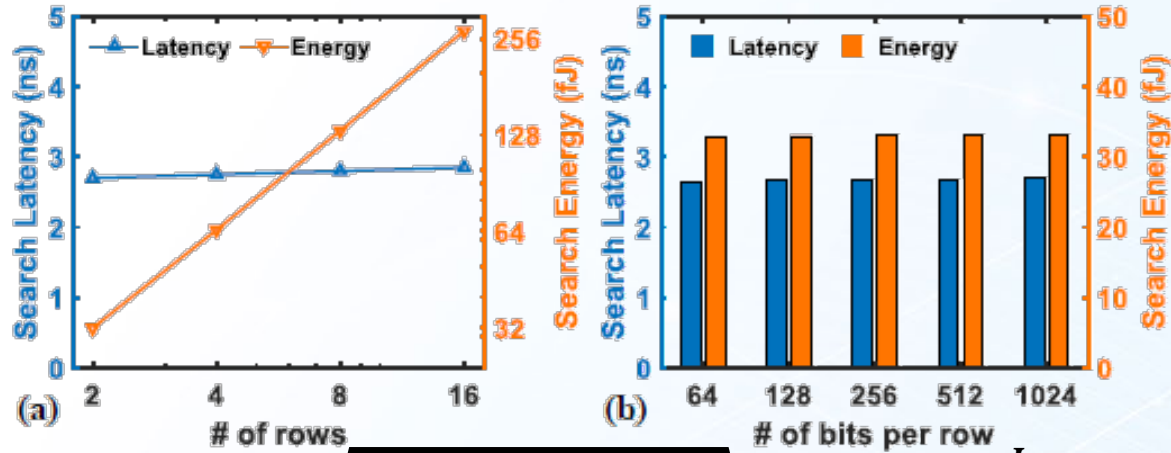
For cosine denominator, i.e. L_2 norm

For cosine numerator, i.e dot product

1FeFET1R Structure



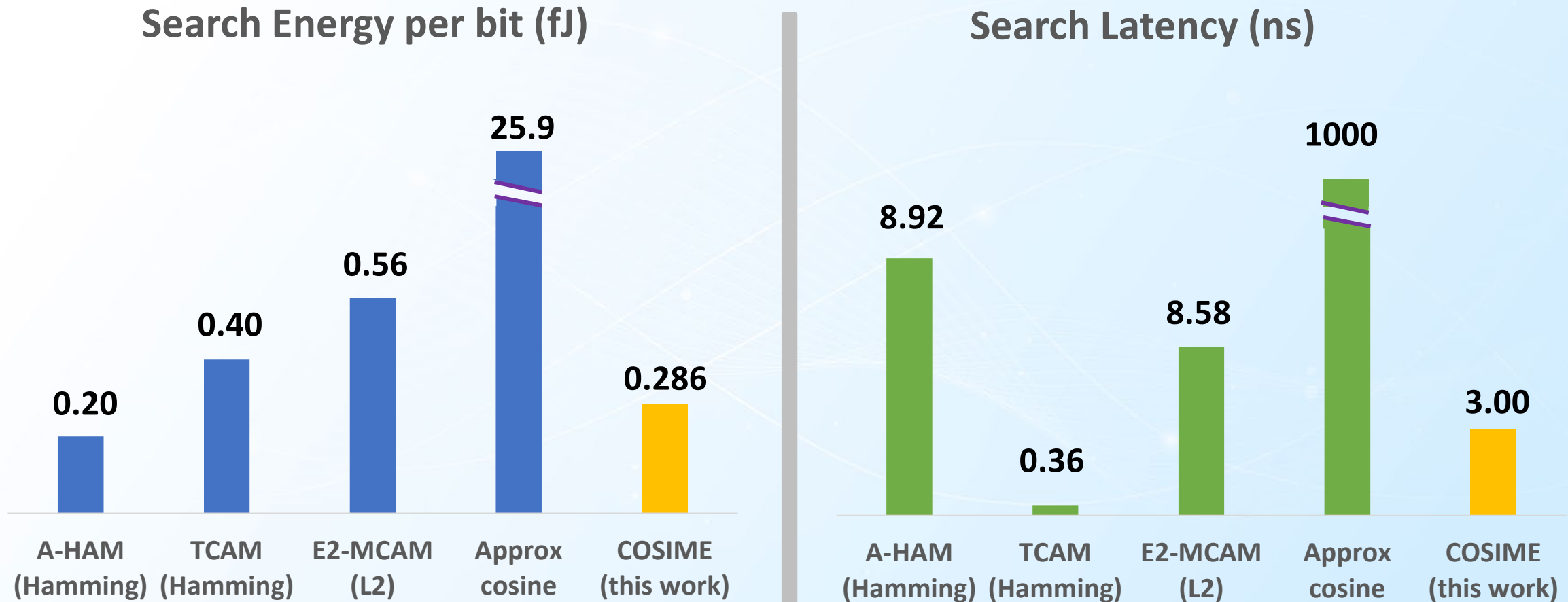
Array-level Evaluation



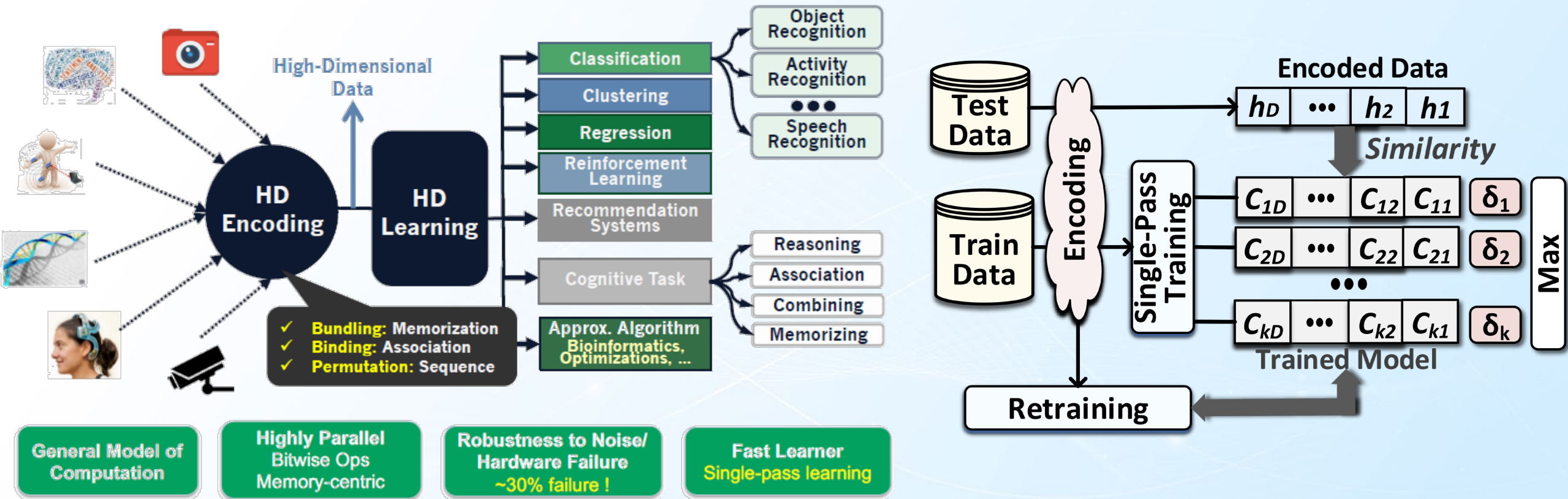
$$\frac{\left(\frac{I_x}{N} \times N\right)^2}{\frac{I_y}{N} \times N} = I_z$$

Variation including:
VDD, FeFET, Resistor in the cell, MOSFET

Existing Associative Memories



HDC Basics



- Limitation: Complex tasks such as Cifar-10 \rightarrow 60%
- Robustness: What if Fp32 MSB get flipped in NN?

Detail operations in Python

- SOTA encoding: Inspired by kernel method [1].

- $\cos(\vec{F} \cdot \vec{B} + \vec{b}) \sin(\vec{F} \cdot \vec{B})$

- Single-pass training:

$$\vec{C}_l = \sum \vec{H}_l$$

- Iterative training: Inspired by Perceptron learning algorithm (PLA).

- MNIST: single-pass 86% Iterative: 93%

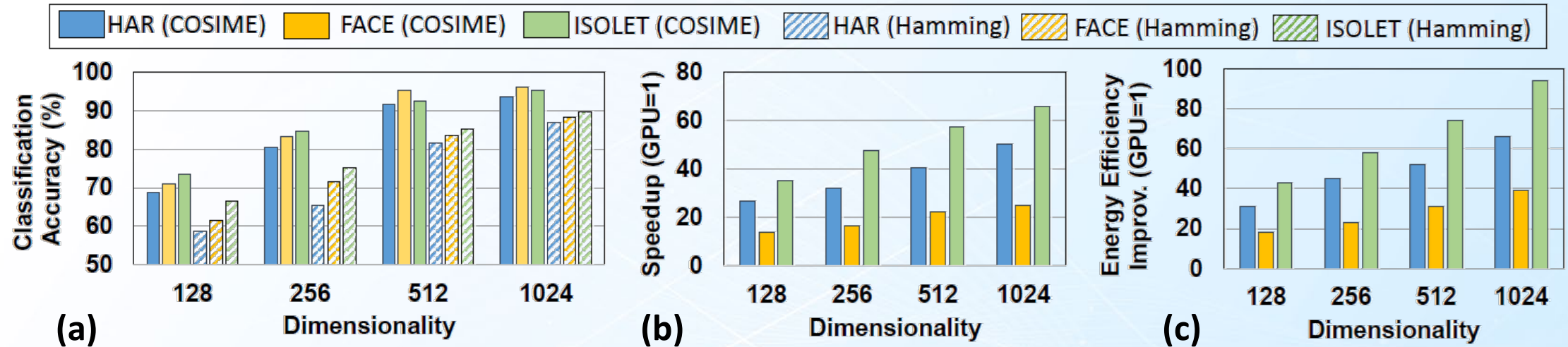
$$\vec{C}_l \leftarrow \vec{C}_l + \eta(1 - \delta)\vec{Q}$$

$$\vec{C}_{l'} \leftarrow \vec{C}_{l'} - \eta(1 - \delta)\vec{Q}$$

- Inference

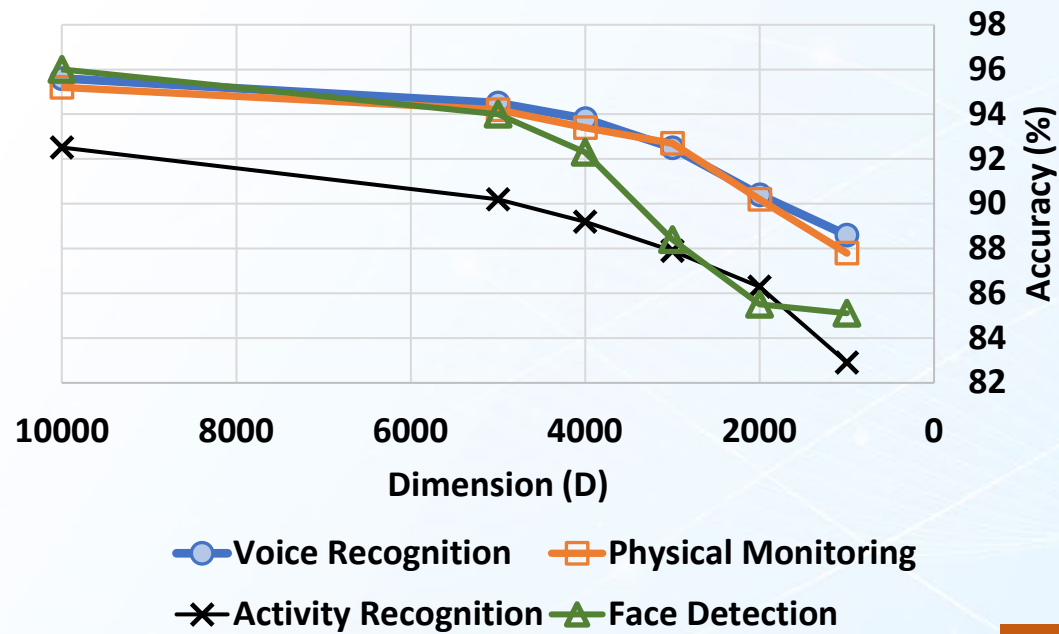
Benchmarking with HDC

	n	K	Train Size	Test Size	Description
UCIHAR	561	12	6,213	1,554	Activity Recognition[40]
FACE	608	2	522,441	2,494	Face Recognition[41]
ISOLET	617	26	6,238	1,559	Voice Recognition [42]



- On average 7% software acc improvement than Hamming distance.
- 47.1x faster/98.5x energy saving compared to 1080 GPU.
- Not so impressive? Because the dimension is low due to hardware.

A view from top-down



HD Computing loses accuracy on lower dimensionality

Challenges and Research Opportunities

- Right memory hierarchy.
- Multibit associative memory.
- Low conducting current. Redesign the translinear operating region.
- Peripherals for NVMs

Conclusion

- A non-approximated design of cosine similarity search in-memory is presented.
- The evaluation and benchmarking results at array and application level show the robustness and superior accuracy to Hamming distance implementation.
- The proposed COSIME design is not limited to FeFET technology, but is rather general.

Thank you for your attention!

cliu29@nd.edu