



**Hewlett Packard  
Enterprise**

# IN-MEMORY COMPUTING WITH MEMRISTOR CONTENT ADDRESSABLE MEMORY CIRCUITS

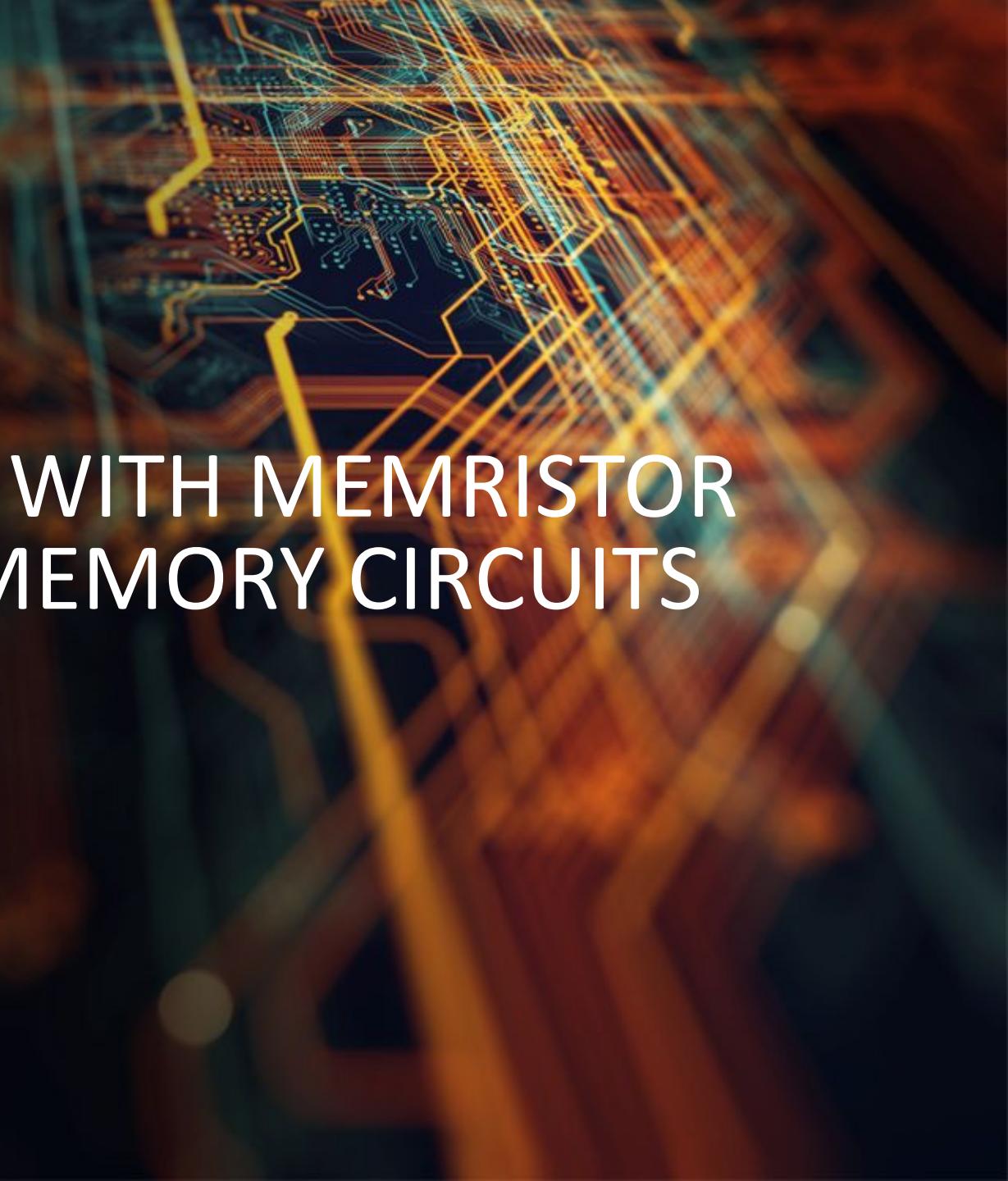
---

Cat Graves

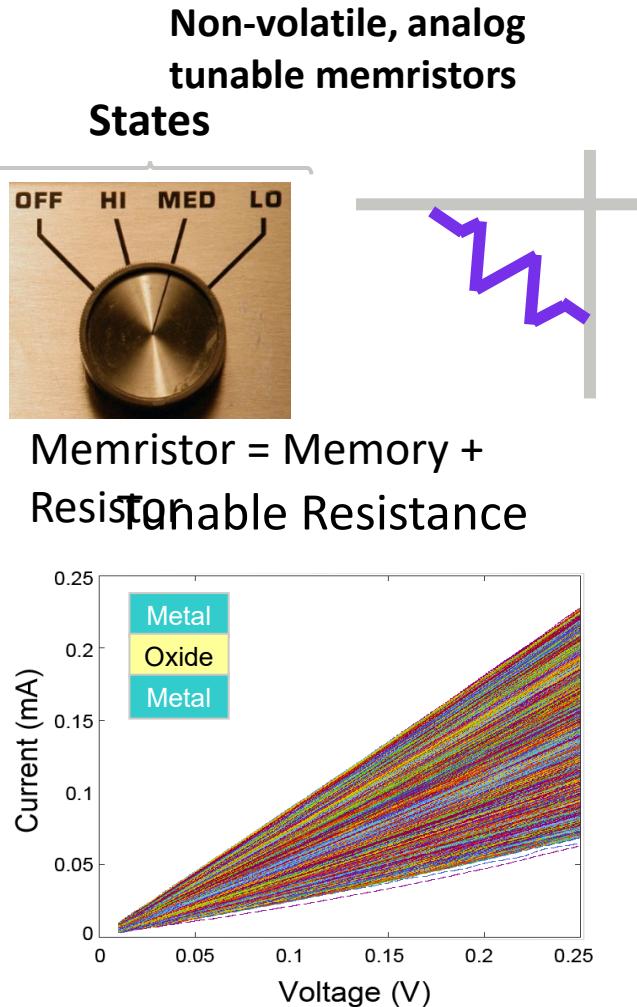
*Hewlett Packard Labs, HPE*

*January 28<sup>th</sup>, 2021*

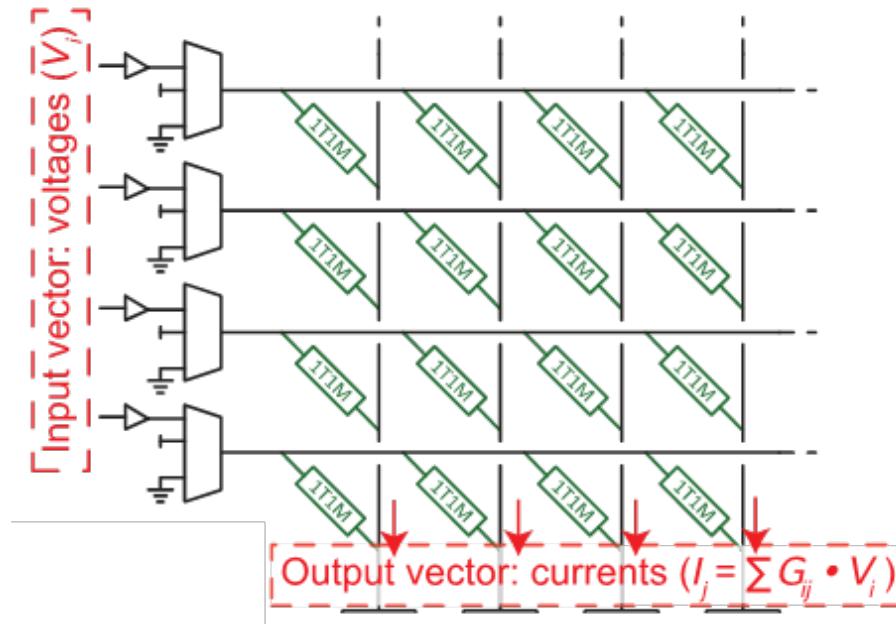
*CRNCH SUMMIT 2021*



# CROSSBARS ACCELERATE MATRIX OPS WITH IN-MEMORY COMPUTE



## Computing: Crossbar circuit of memristors



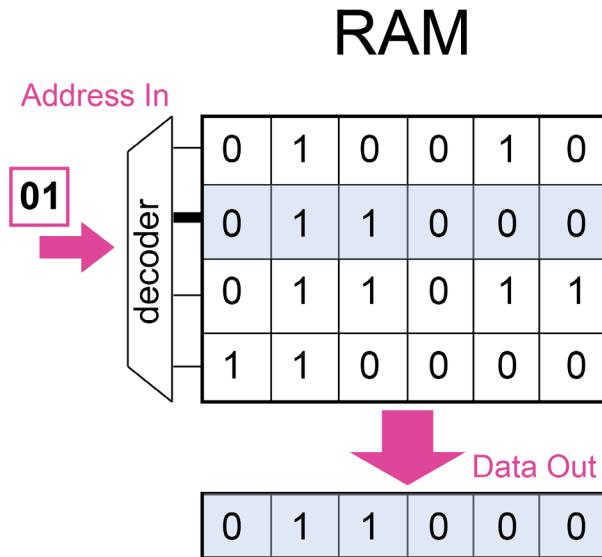
Note: many non-volatile technologies may be used  
(PCM, STT-MRAM, Ferroelectric, Flash)

## Function: vector-matrix multiplication

- Machine Learning
  - Fully connected layers in neural networks
- Scientific computing
  - PDE solver
  - Linear equation solver
- Signal processing
  - FFT, DCT, compression
- Optimization
  - Core for Hopfield NN

# ANOTHER IN-MEMORY COMPUTING PRIMITIVE:

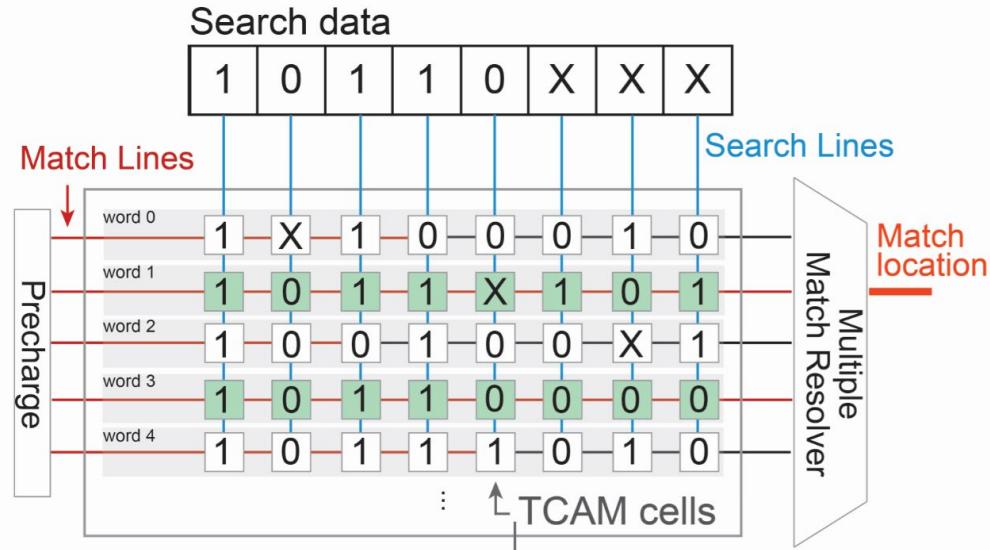
## Content Addressable Memories (CAMS)



### Traditional Random Access Memory (RAM)

- **Input:** Address, **Output:** Contents at this address
- Exact and precise storage of information

### Ternary Content Addressable Memory (TCAM)



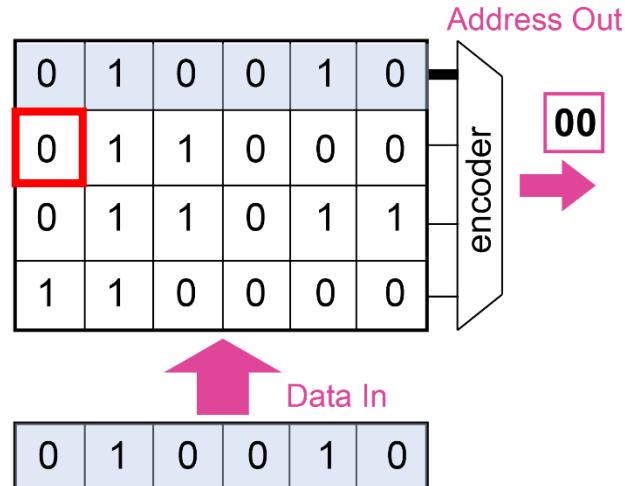
### Content Addressable Memories (CAM) and Ternary CAMs

- **Input:** Search word (content), **Output:** address of a match
- Ternary CAMs store and search 0,1, and X = 'don't care'
- High performance: search word compared to all entries simultaneously
  - Rapid search time, No collisions
  - Used in network routing, highly associative caches, translation lookaside buffers (TLBs), microarchitectural queues etc.

# USE NON-VOLATILE MEMORY INSTEAD OF SRAM IN TCAM CELL

TCAMs cost **30 times more \$ / bit** than DDR SRAMs, and consume **150 times more power / bit**

**CAM/TCAM cell**

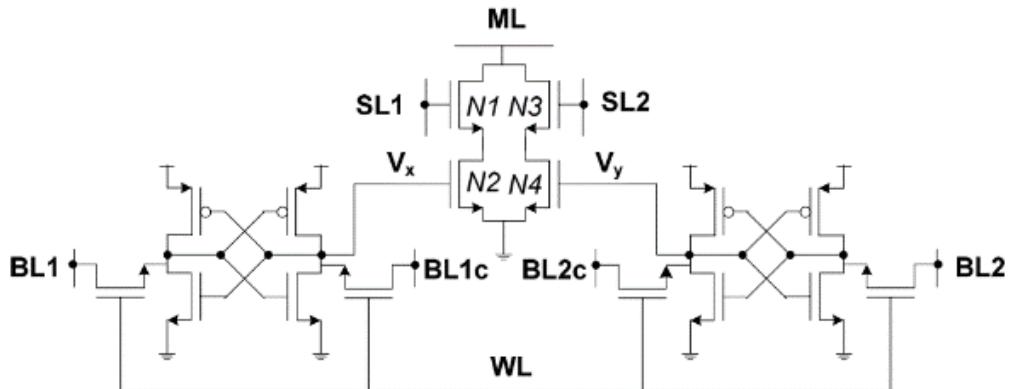


Y. Tsukamoto, et al. Symp. on VLSI Circuits, 2015

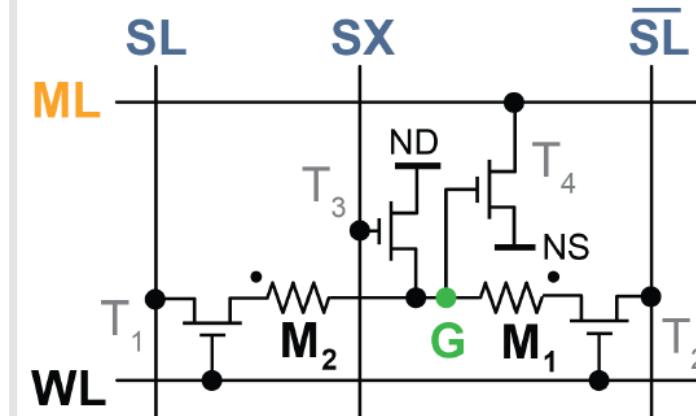
- Utilize non-volatile memristors with no leakage power
  - 2-3.5x lower area than SRAM-TCAM
  - 35x lower power than SRAM-TCAM
  - However, much slower re-programming time (us rather than ns)
- Good match for in-memory computing

C. E. Graves, et al, NANOARCH (2018);  
C. E. Graves, et al, ICRC (2018);  
C. E. Graves, et al, IEEE TNano, (2019)

**16T SRAM-based TCAM cell**



**4T-2Memristor TCAM cell**

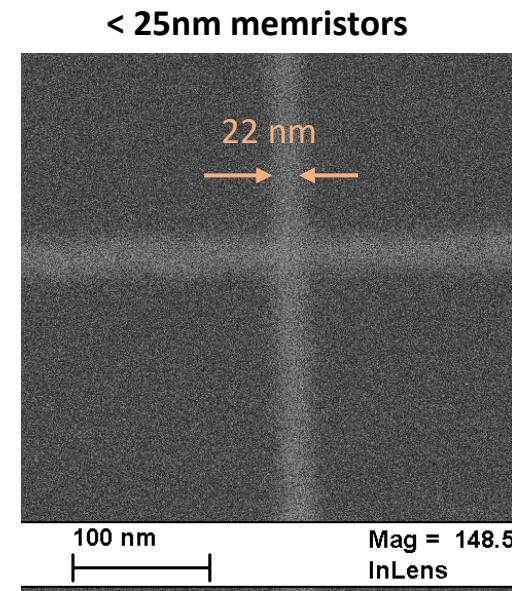
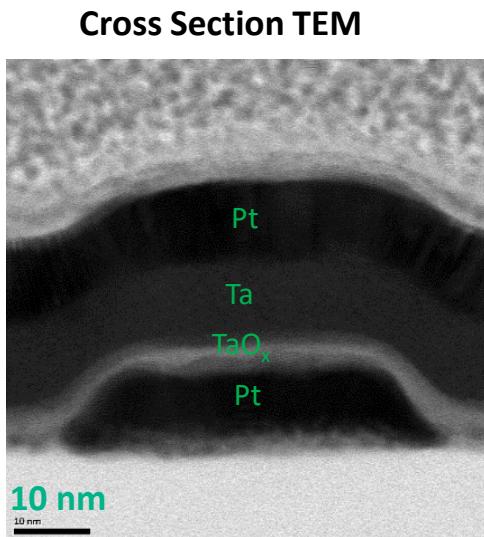
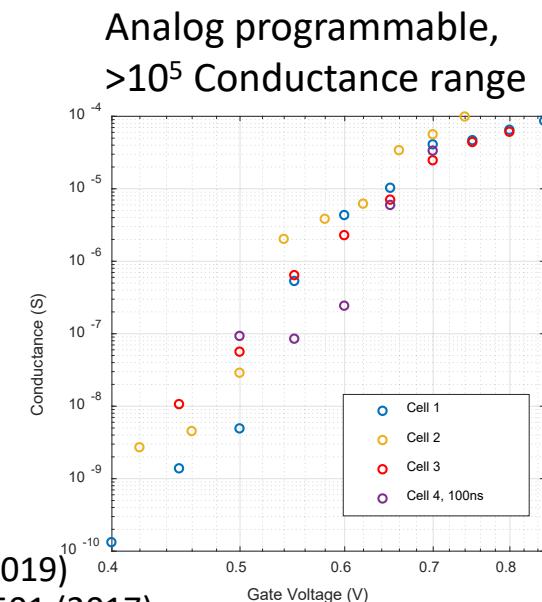
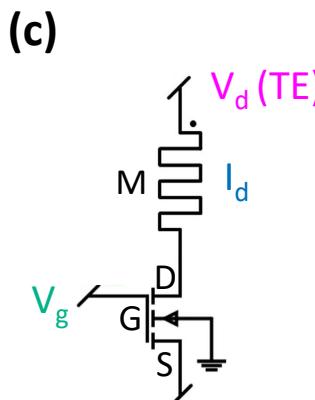
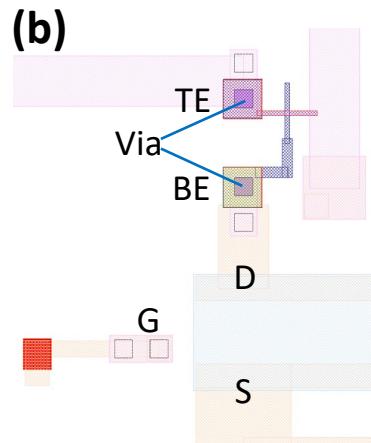
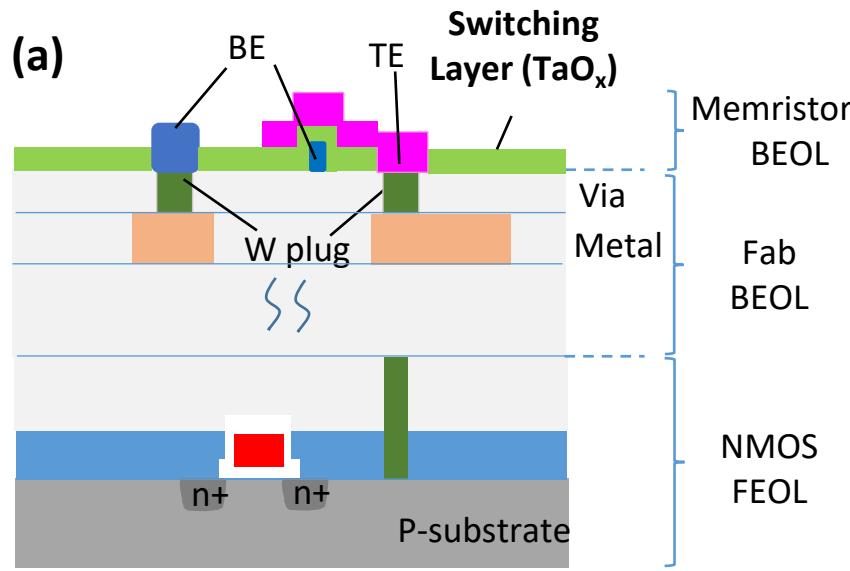


TCAM content	(M <sub>1</sub> , M <sub>2</sub> )
0	(0,1)
1	(1,0)
X	(0,0)

0 = HRS, 1 = LRS

ML = Match Line  
WL = Word Line  
SL/SL-bar = Search Line

# NANOSCALE MEMRISTOR INTEGRATION WITH CMOS



How many bits? ( $2\sigma_{max}$  spacing)

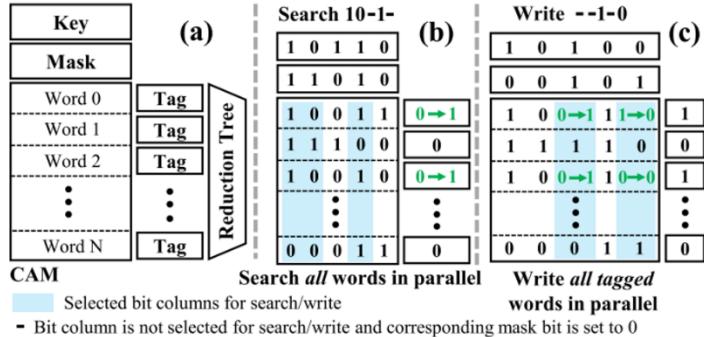
G Range (S)	$\sigma_{max}$ (S)	Number of States	Number of Bits
$10^{-7}-10^{-6}$	$2.5 \times 10^{-7}$	4	2
$10^{-6}-10^{-5}$	$3 \times 10^{-7}$	20	>4
$10^{-5}-3 \times 10^{-4}$	$1 \times 10^{-6}$	150	>7
$3 \times 10^{-4}-6.2 \times 10^{-4}$	$3 \times 10^{-7}$	532	>9

Trade-off between #bits and current levels

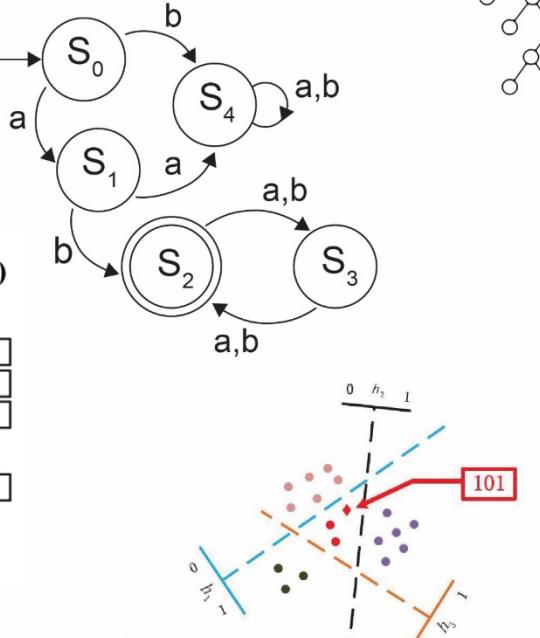
# IN-MEMORY COMPUTING OPPORTUNITIES WITH CAMS

## Lower power and flexible NVM designs drive new applications

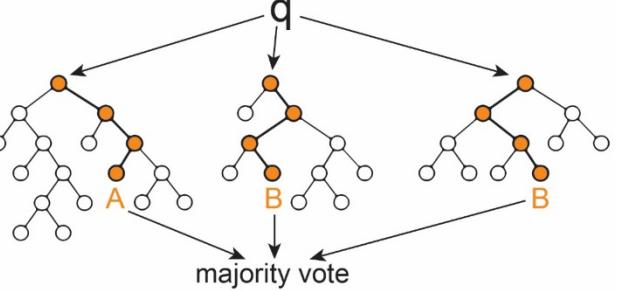
### Associative Computing



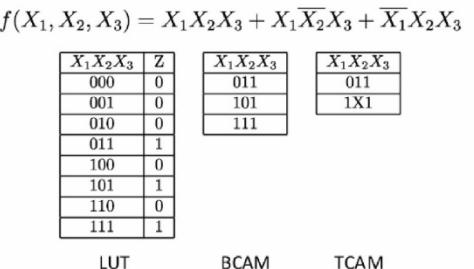
**Finite State Machines**  
(Regular Expression matching for network security, branch prediction and cache controllers for HPC)



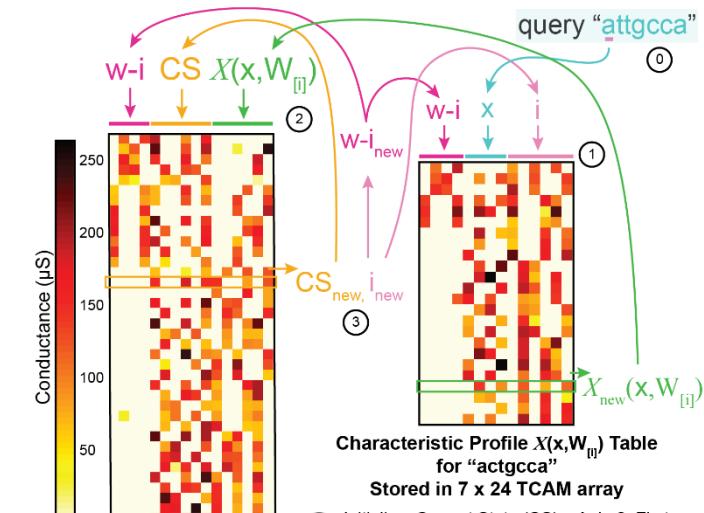
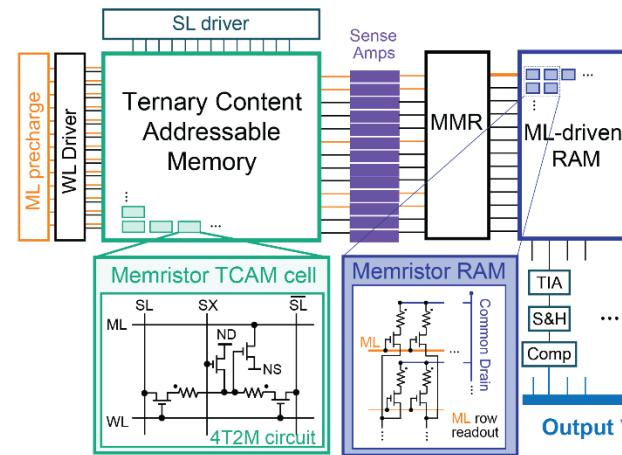
**Data Mining**  
(Locality sensitive hashing (LSH), TF-IDF, data-intensive workloads)



**Tree-based Models**  
(Random Forests, Business Rule Management Systems)



**Reconfigurable Computing**  
(Arbitrary logic functions, associative computing)



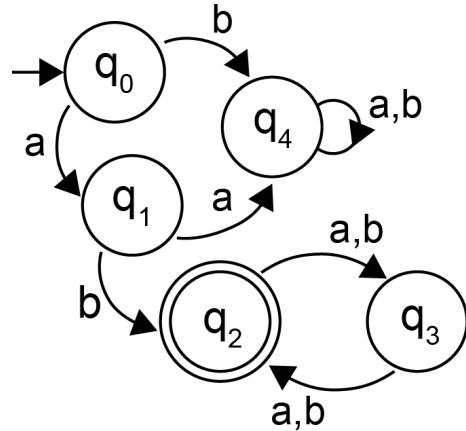
- Levenshtein Automaton**  
can match inexact patterns for genomics read mapping
- |          |                  |
|----------|------------------|
| actgcca  | LEV distance = 1 |
| attgcca  | ✓                |
| actcca   | ✓                |
| actgacca | ✓                |
| atcgaca  | x                |
- Characteristics:**
- Update CS and i with result from (2), take next character in query and repeat from (1) until CS identified as accept or fail state

**Edit distance or Levenshtein Automaton**  
(Genomics)

# USING CAMS FOR COMPUTING: FINITE AUTOMATA

Deterministic Finite Automaton  
to accept  $w$  when  $w$  is of even  
length and begins with 'ab'

Transition Diagram



OR

State Transition Table

Current State	Input	Next State
* $q_0$	a	$q_1$
* $q_0$	b	$q_4$
$q_1$	a	$q_4$
$q_1$	b	$q_2$
$q_2$	a	$q_3$
$q_2$	b	$q_3$
$q_3$	a	$q_2$
$q_3$	b	$q_2$
$q_4$	a	$q_4$
$q_4$	b	$q_4$

Search word =  
Current State + Input

Returns  
Next State

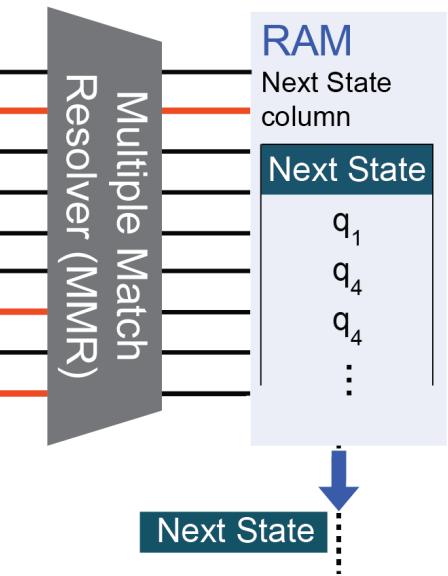
State Transition Table directly  
encoded in TCAM + RAM hardware

TCAM array

Current State and Input columns

Current State	Input
$q_0$	a
$q_0$	b
$q_1$	a
⋮	⋮

Current State Input

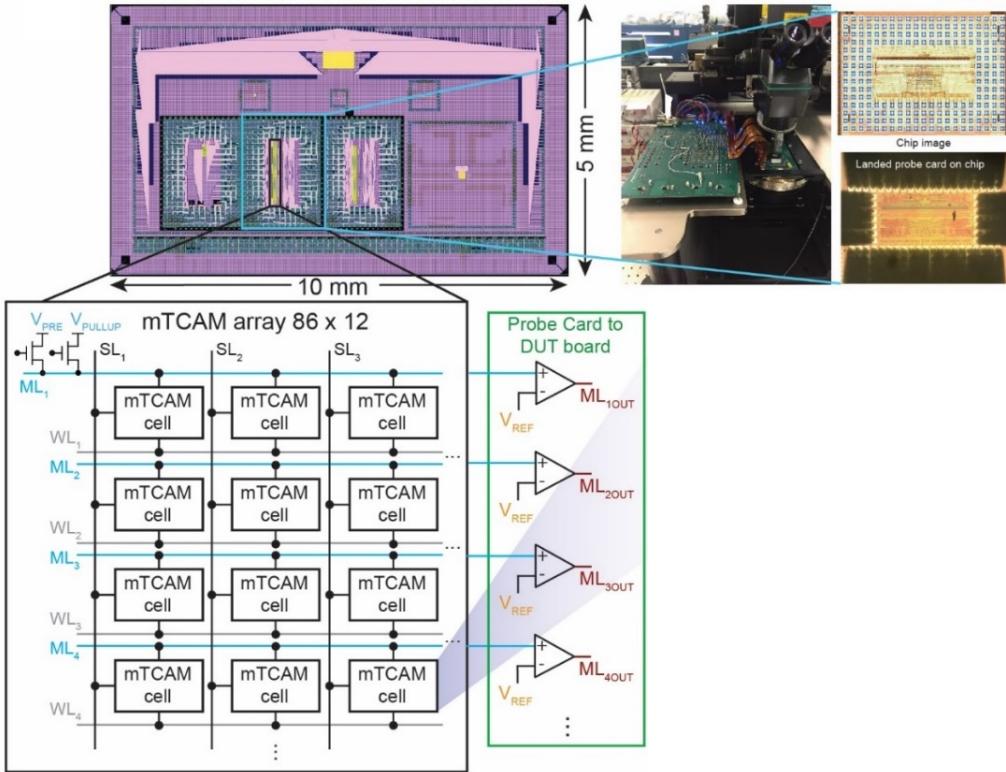


Input String "abaabba..."

Finite Automata have applications in Bio-informatics/Genomics and Security

# EXPERIMENTAL DEMONSTRATIONS FOR NETWORK SECURITY

Non-volatile TCAM arrays taped-out  
180nm CMOS + 50 nm Integrated memristors



C.E. Graves, et al, Advanced Materials, 32, 2003437 (202

Regular Expression matching on chip

RegEx: "sid=[0-9A-F]{32}"

Naive STT: 37 states, 516 rows

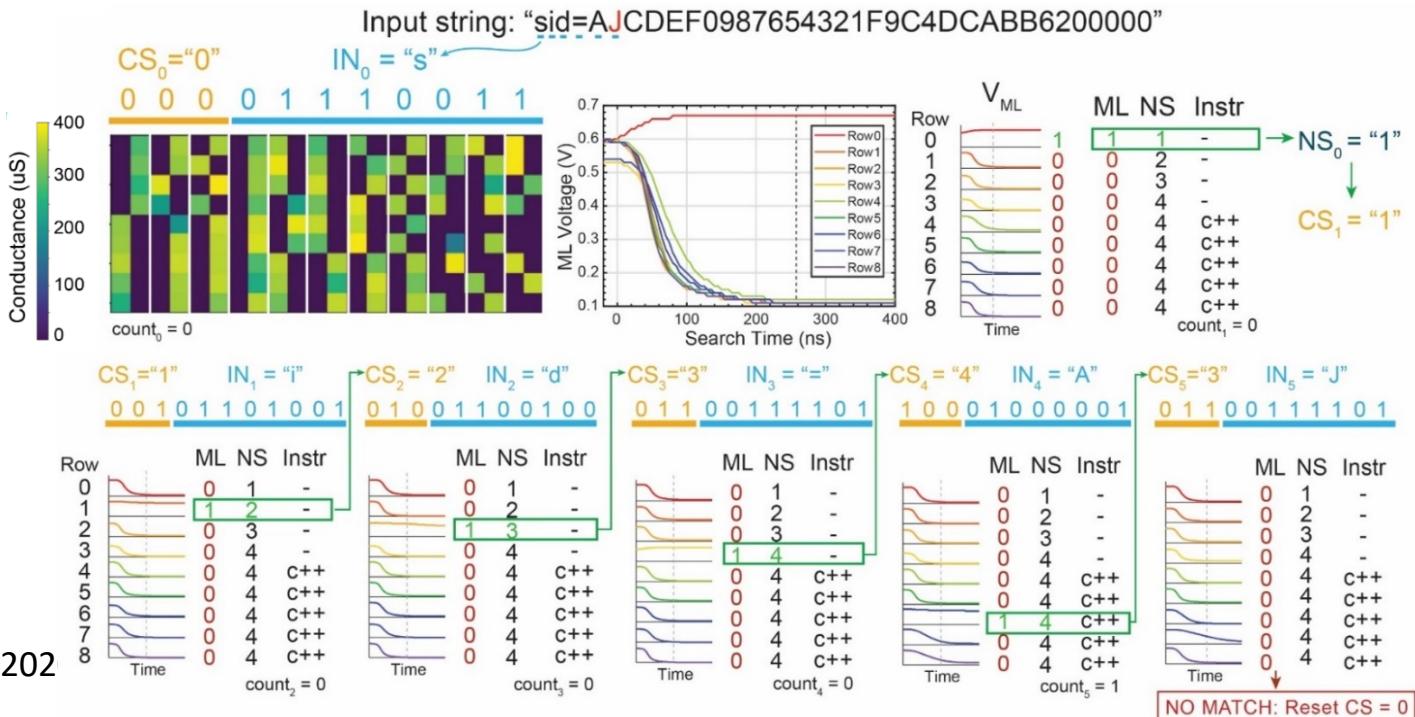
CS	Input	NS
0	01110011 (s)	1
1	01101001 (i)	2
2	01100100 (d)	3
3	00111101 (=)	4
4	[0-9A-F] 16 entries	5
5	[0-9A-F] 16 entries	6
.....	.....	.....
35	[0-9A-F] 16 entries	36

Compressed STT: 5 states, 9 rows

CS	Input	NS	Inst.
000	01110011 (s)	1	
001	01101001 (i)	2	
010	01100100 (d)	3	
011	00111101 (=)	4	Count++
100	00110XXX	4	Count++
100	0011100X	4	Count++
100	01000X1	4	Count++
100	01000X10	4	Count++
100	0100010X	4	Count++

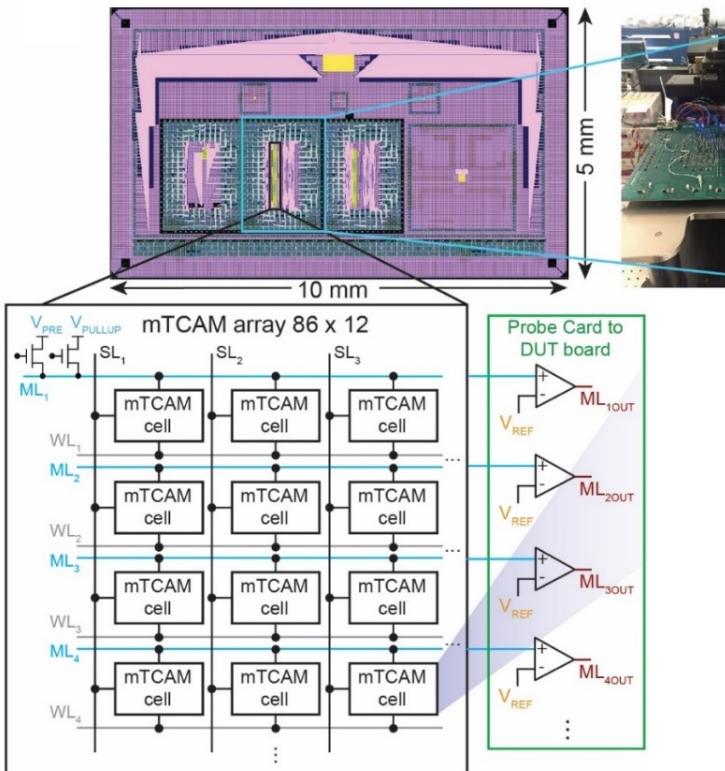
Compressed STT for TCAM array

0	0	0	0	1	1	1	0	0	1
0	0	1	0	1	1	0	1	0	0
0	1	0	0	1	1	0	0	1	0
0	1	1	0	0	1	1	1	0	1
1	0	0	0	1	1	0	X	X	X
1	0	0	0	1	1	1	0	X	X
1	0	0	0	1	1	1	0	X	X
1	0	0	1	0	0	0	X	1	
1	0	0	1	0	0	0	X	10	
1	0	0	1	0	0	0	1	0	X



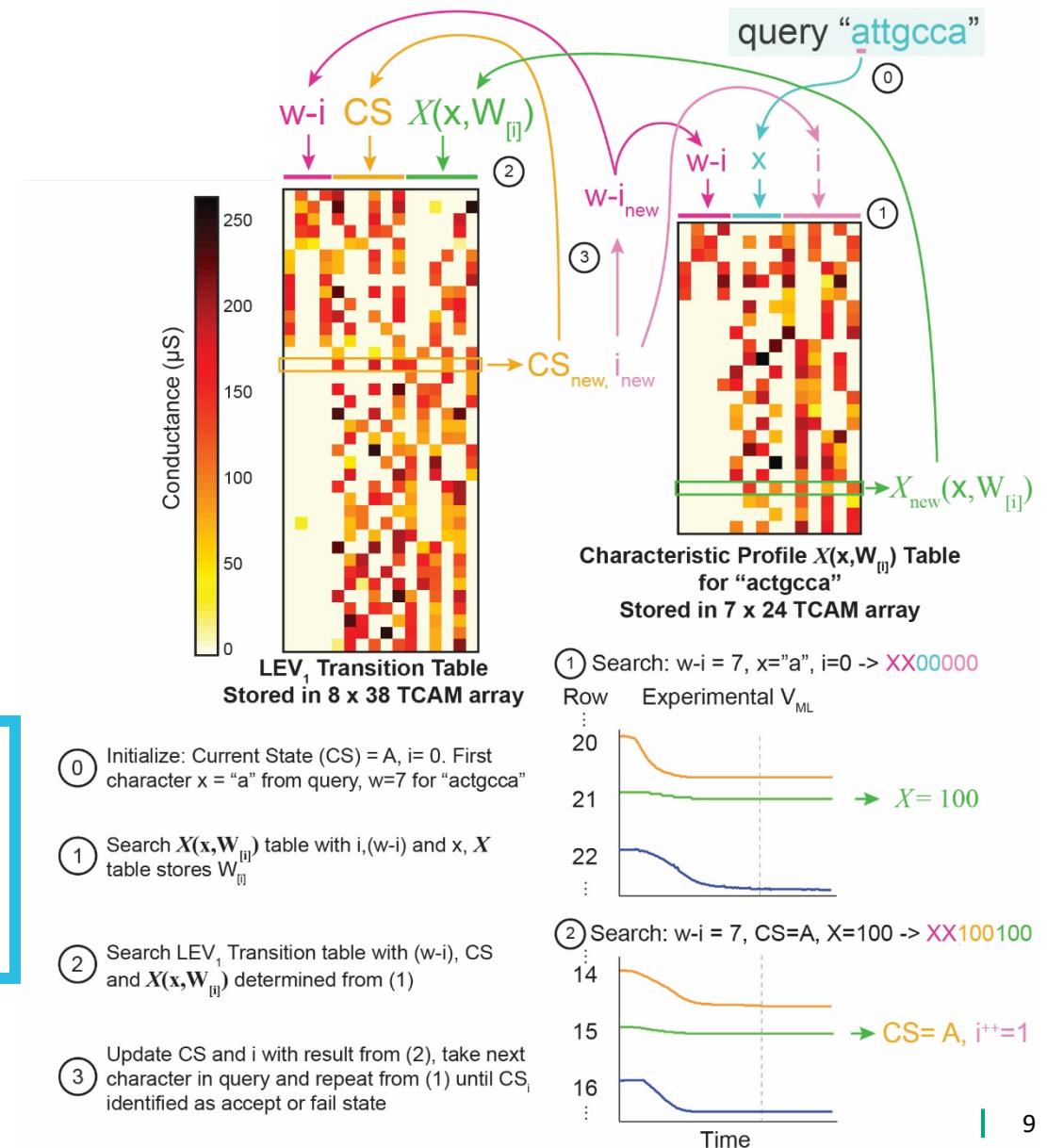
# EXPERIMENTAL DEMONSTRATIONS FOR GENOMICS APPLICATIONS

Non-volatile TCAM arrays taped-out  
180nm CMOS + 50 nm Integrated memristors

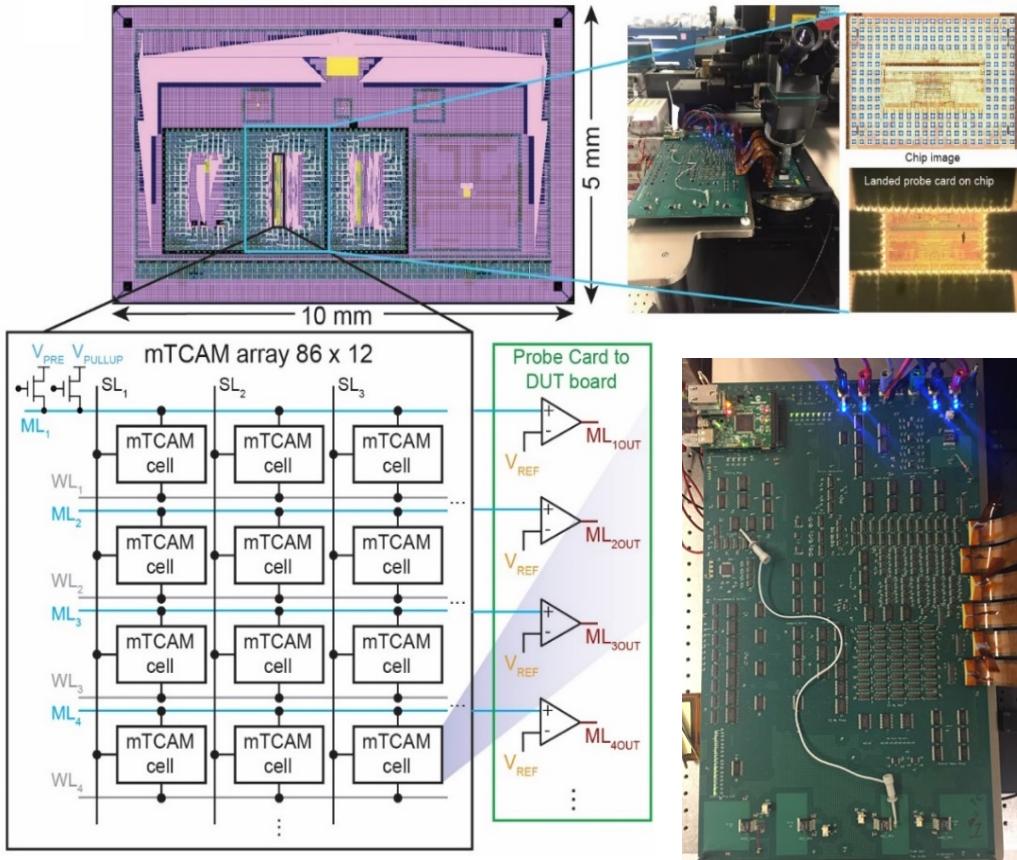


## Edit Distance Automata

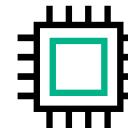
<b>actgcca</b> <b>attgcca</b> <b>actcca</b> <b>actgacca</b> <b>atctgaca</b>	<b>LEV distance = 1</b> <span style="color: green;">✓</span> <span style="color: green;">✓</span> <span style="color: green;">✓</span> <span style="color: red;">X</span>
---	---



# PERFORMANCE: 25X IMPROVEMENT VS STATE-OF-THE-ART FPGA



Performance on a Network Security (“SNORT”) ruleset:

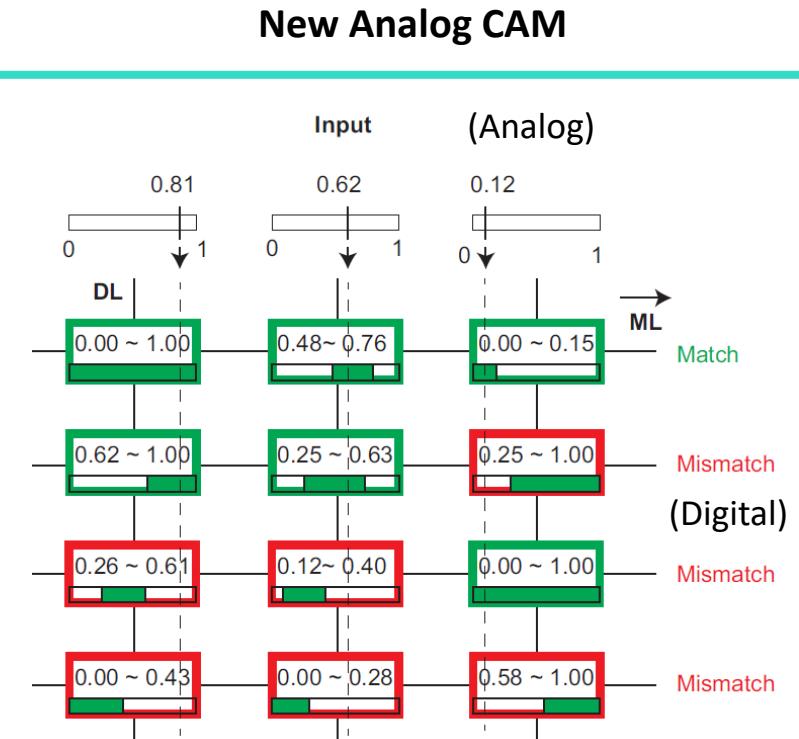
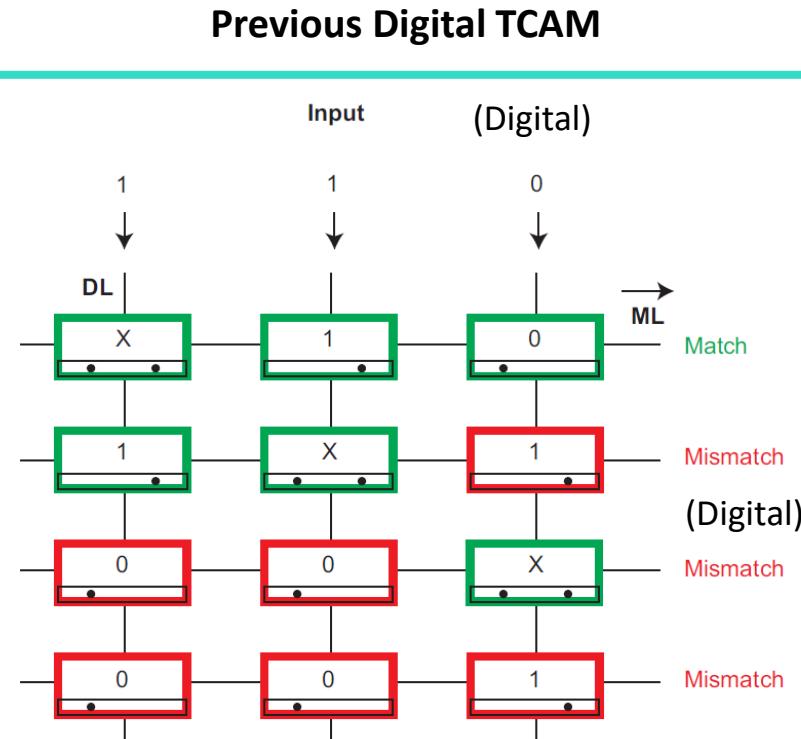


**47.2 Gbps at 0.3W**  
Memristor-TCAM system  
7-var stride [1]

**3.9 Gbps at 0.63W**  
FPGA, Best reported  
4-stride [2]

**25x improved Throughput/Watt**  
(with smaller area and greatly improved scalability)

# NEW INVENTION: ANALOG CONTENT ADDRESSABLE MEMORY (ACAM)



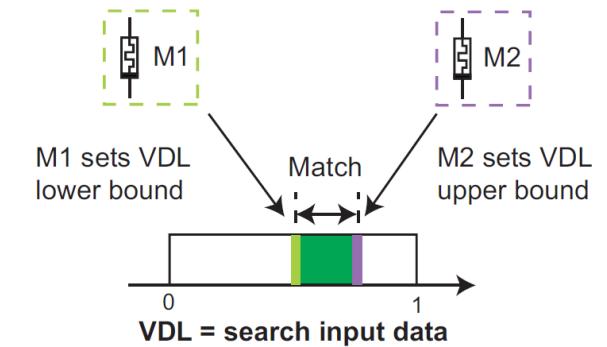
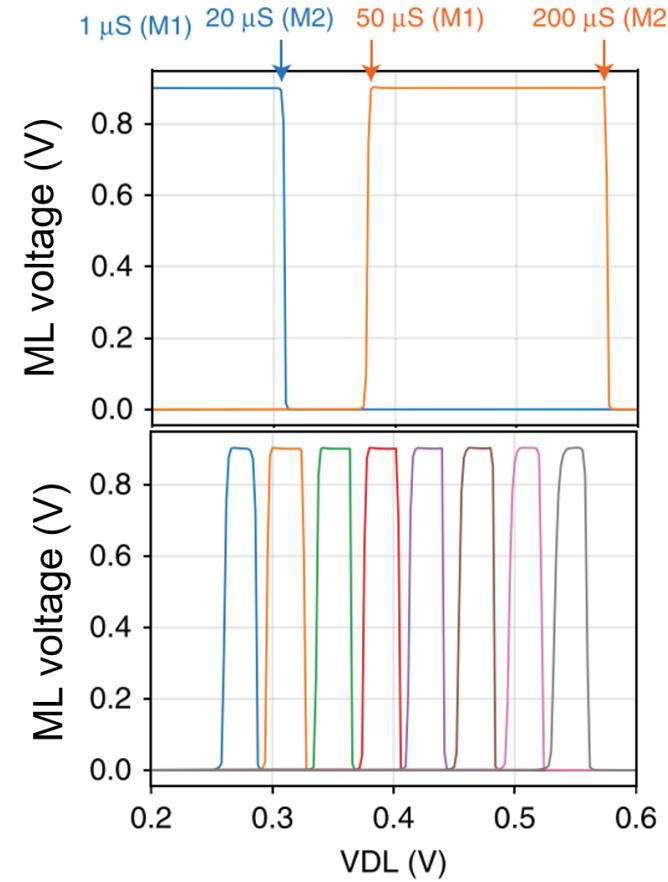
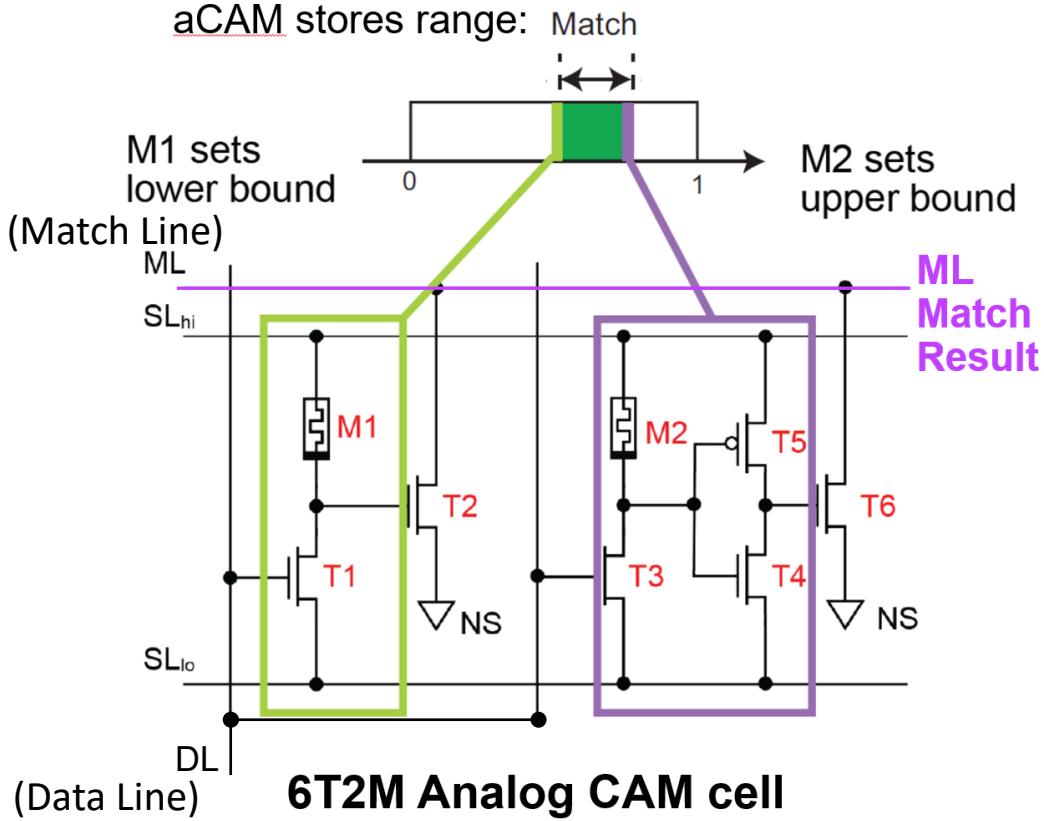
Digital CAM / TCAM

Stores: 0,1,or 'X'  
Input Search: 0,1, or 'X'

Analog CAM

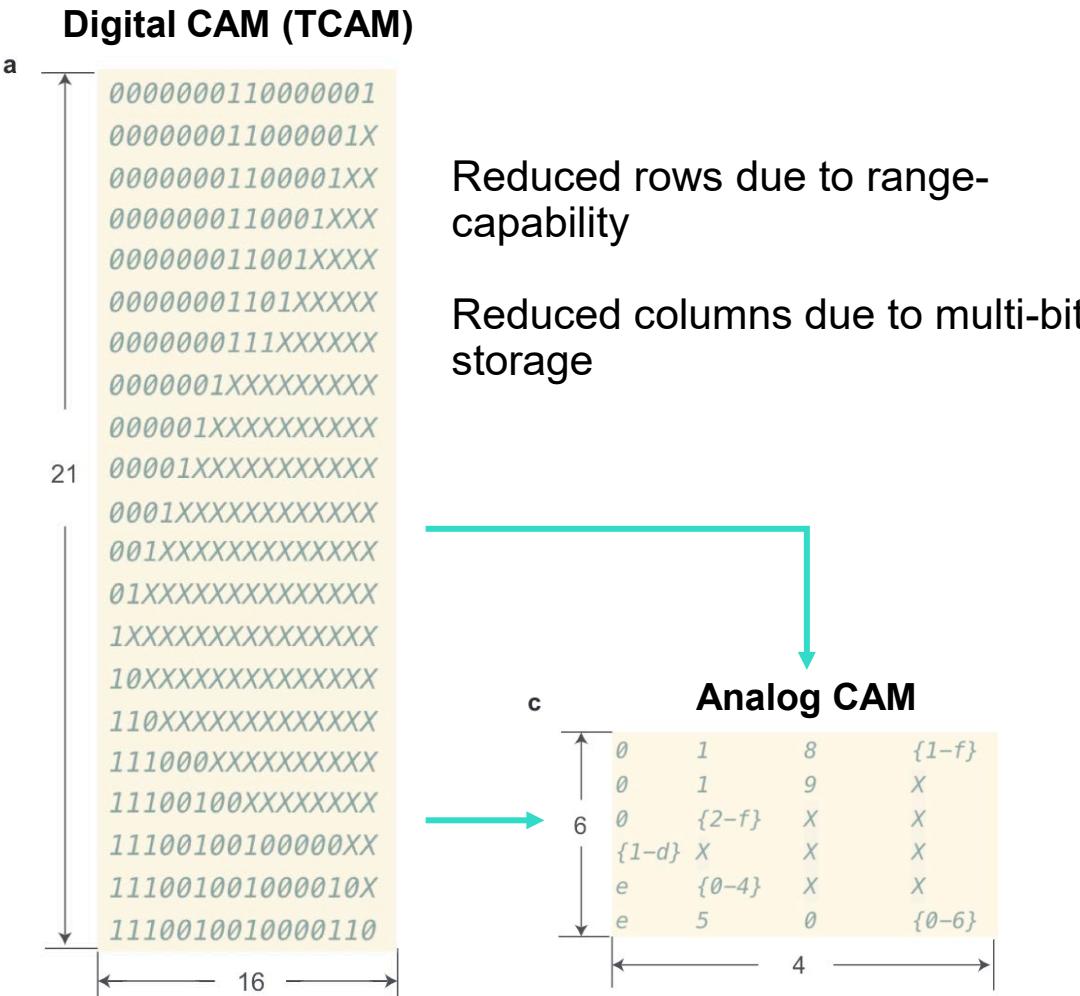
Stores: analog ranges  
Input Search: analog values (multi-bit)

# ANALOG CAM (A-CAM) CIRCUIT USING MEMRISTOR TUNABLE CONDUCTANCE



# AREA/ENERGY COMPARISON TO HIGHLY OPTIMIZED SRAM-TCAM

## Class B IP address routing table



Over 4x lower energy/bit/search

18.8x reduction in Area (@16nm)

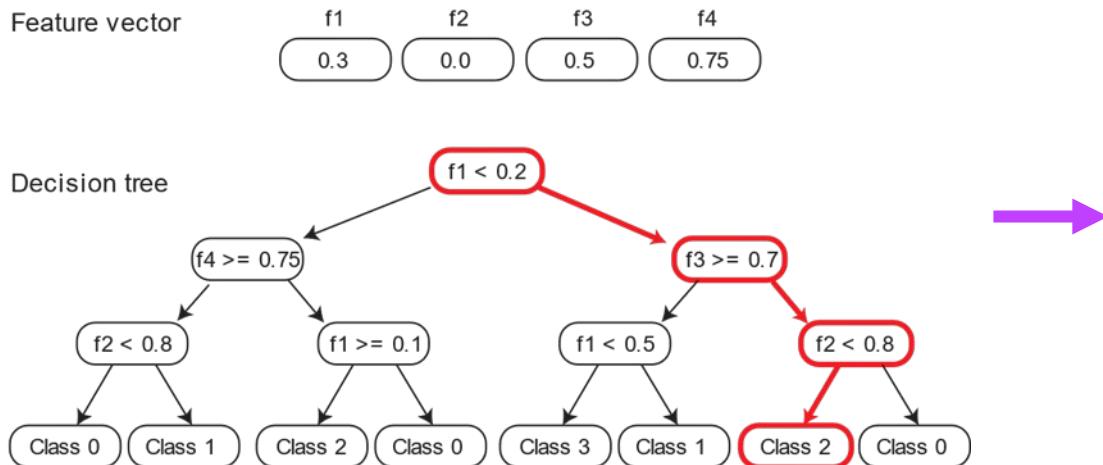
and a-CAM is non-volatile

Supply	Energy/search	Energy/search/ equiv. TCAM bit
ML precharging	102.9 fJ	0.007 fJ
SLhi driver	298.5 fJ	0.021 fJ
Others	86.4 fJ	0.006 fJ
DAC	52.1 fJ	0.004 fJ
<b>Total</b>	<b>539.9 fJ</b>	<b>0.037 fJ</b>

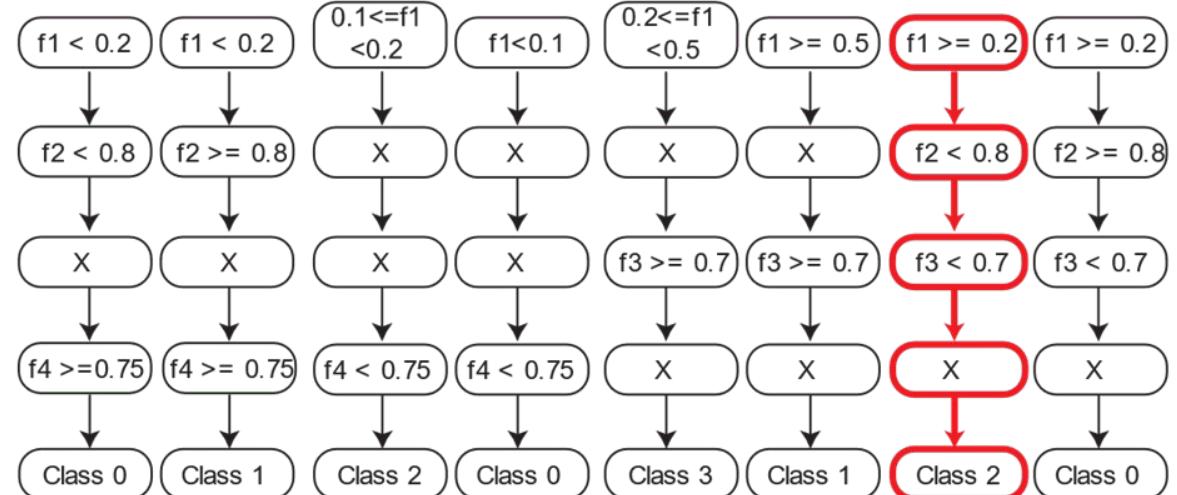
vs 0.165 fJ with highly optimized SRAM-TCAM  
(SRAM has >10x power saving techniques not implemented in our a-CAM )

# NEW APPLICATIONS ENABLED: DECISION TREES WITH A-CAM

**Decision Trees & Random Forests:** Easy to train; support smaller training sets;  
Highly interpretable and verifiable



Represent each root to leaf path as a chain with a series of nodes to be evaluated on and combine nodes associated with same feature and add 'Don't Care' or 'X' nodes for features not evaluated in chains

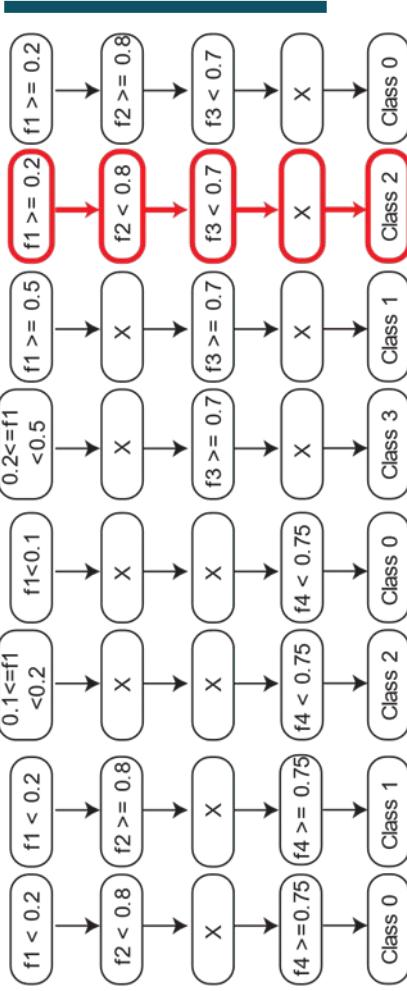


# NEW APPLICATIONS ENABLED: DECISION TREES WITH A-CAM

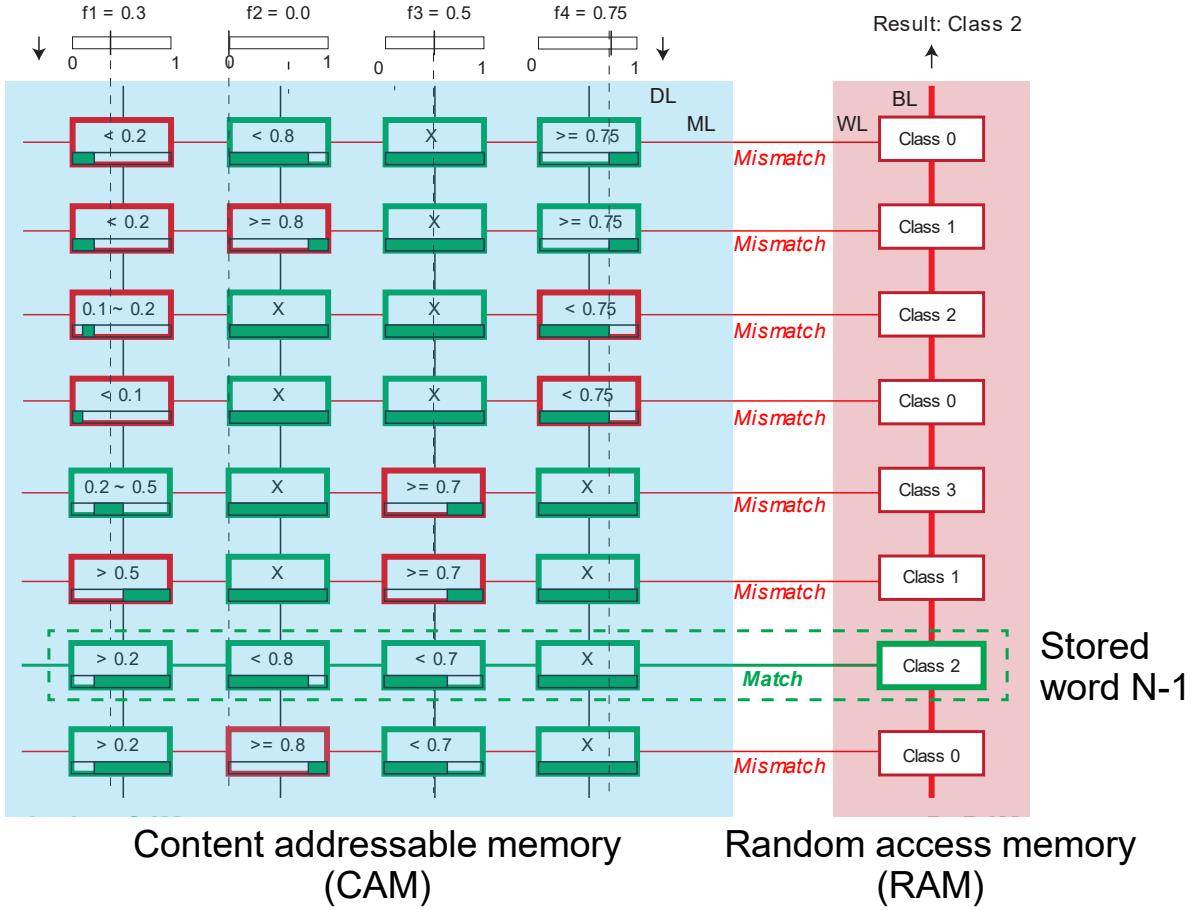


Search feature vector on cols

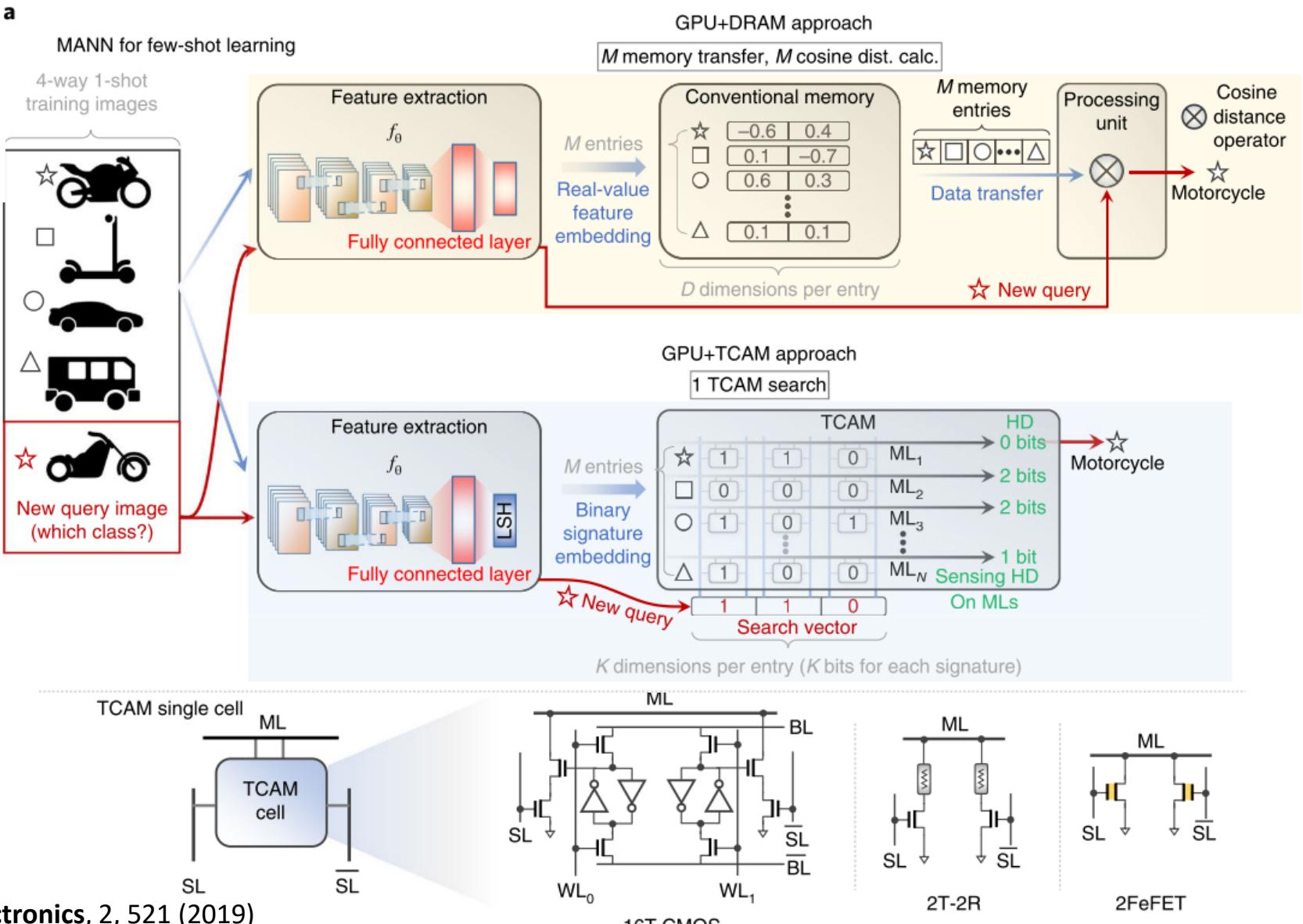
Represent each root to leaf path as a chain with a series of nodes to be evaluated on and combine nodes associated with same feature and add 'Don't Care' or 'X' nodes for features not evaluated in chains



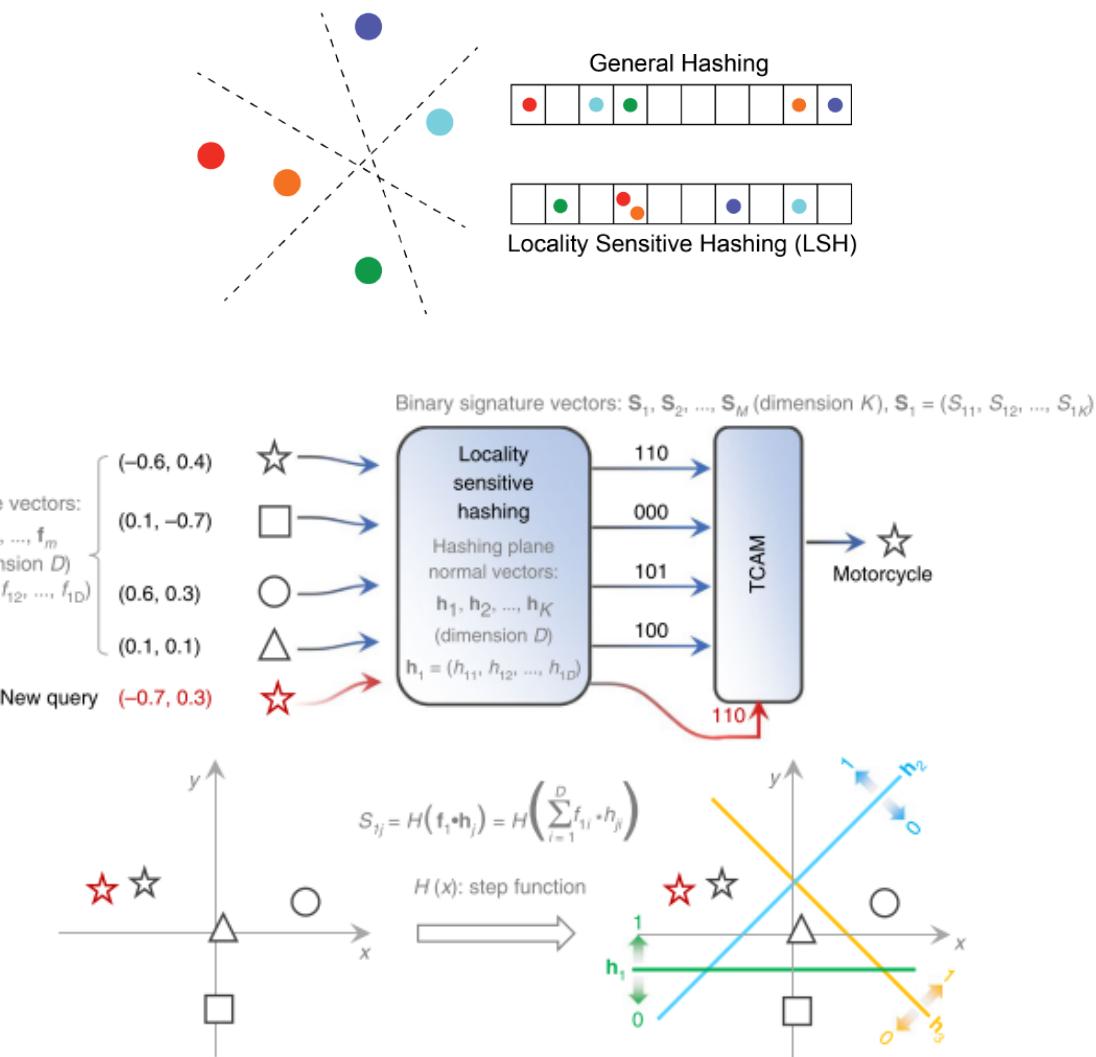
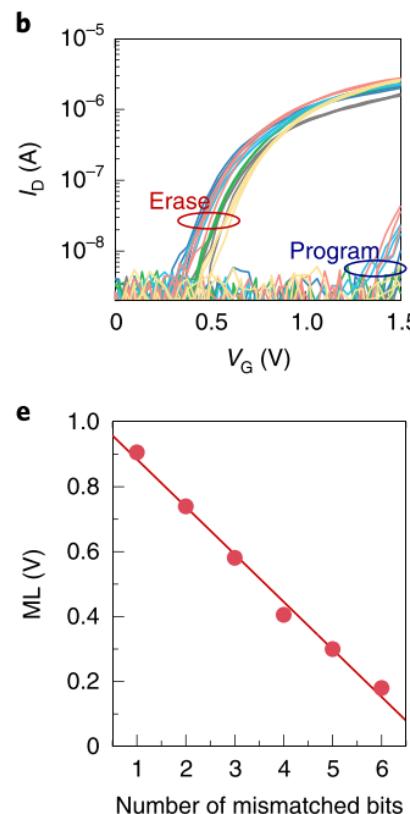
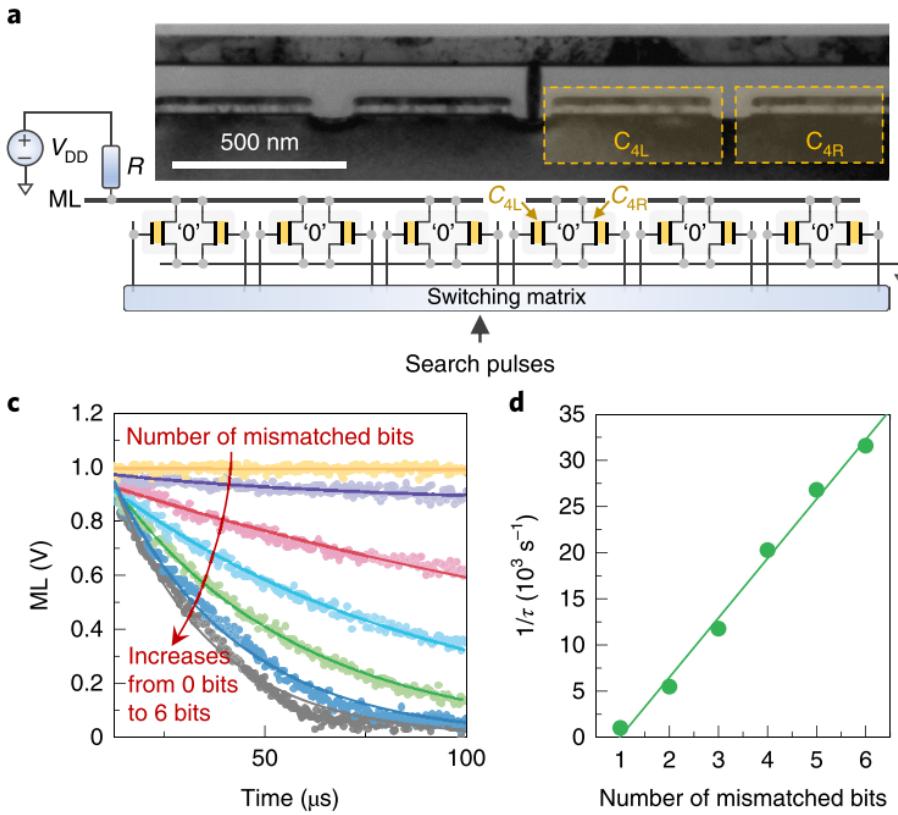
Rotate chains and map onto analog CAM and associated RAM for classification



# OTHER RECENT CAM COMPUTING WORK: LOCALITY SENSITIVE HASHING WITH FE-CAM



# OTHER RECENT CAM COMPUTING WORK: LOCALITY SENSITIVE HASHING WITH FE-CAM



# OTHER RECENT CAM COMPUTING WORK: ASSOCIATIVE PROCESSING

AP operates concurrently on every word where the data is stored, eliminating the need for data movement.

Word-level parallelism points to applications with high data-level parallelism

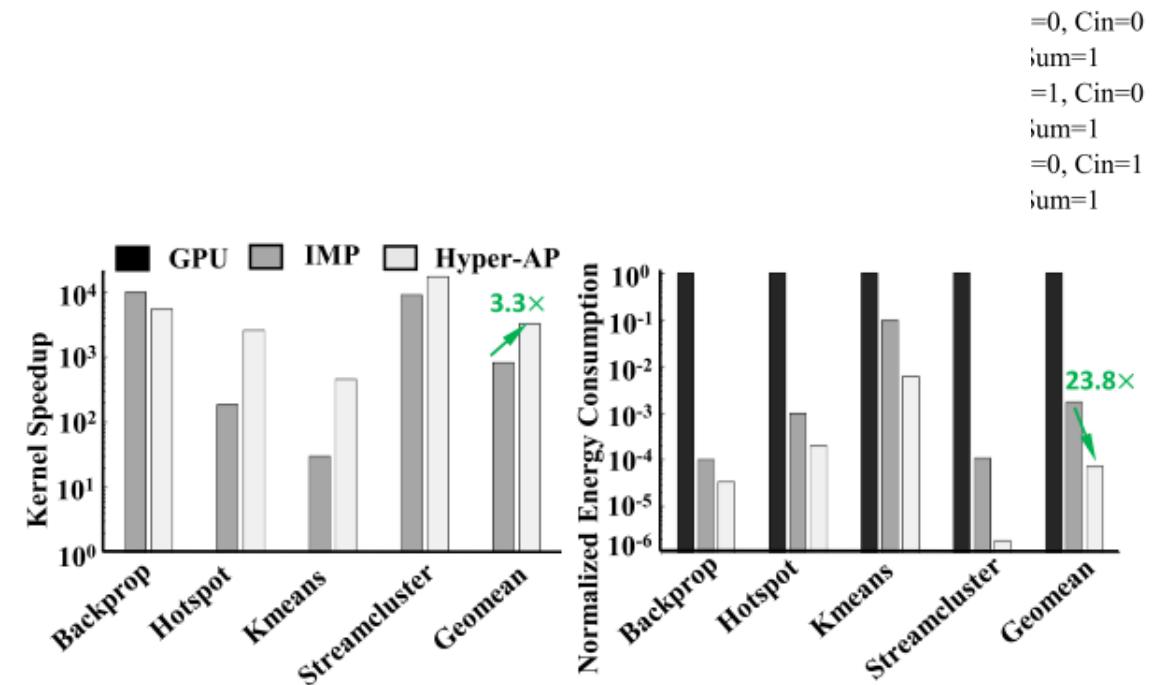
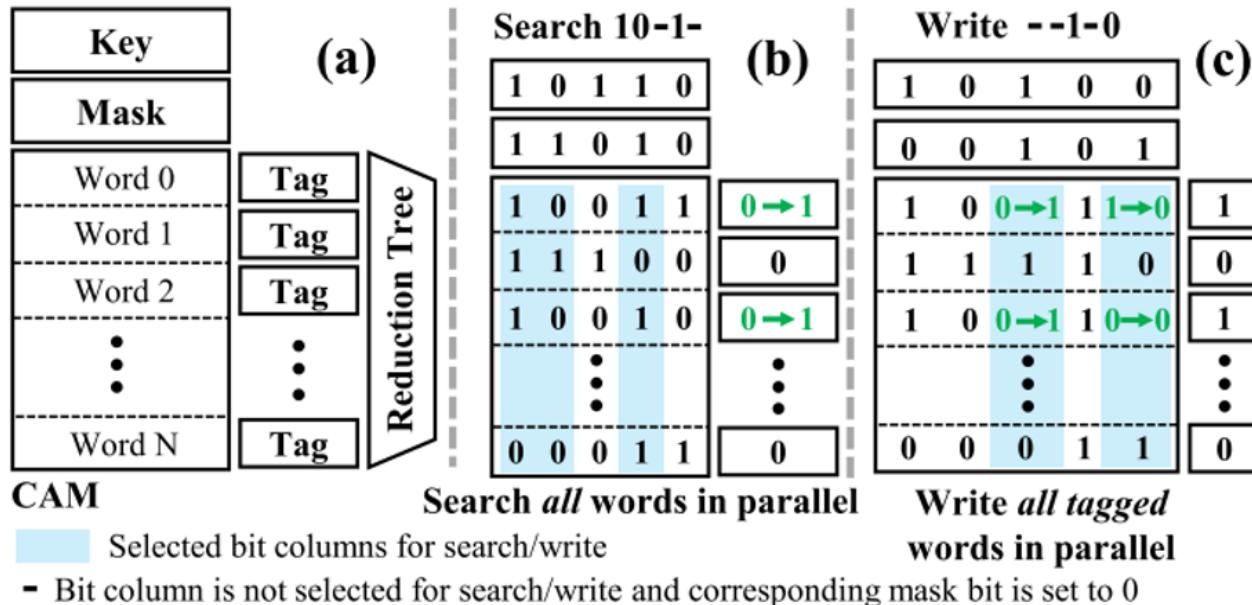


Fig. 18. The speedup (left) and energy consumption (right) results for the Rodinia kernels. The energy consumption of *Hyper-AP* and *IMP* are normalized to that of *GPU*.

=0, Cin=0  
↓sum=1  
=1, Cin=0  
↓sum=1  
=0, Cin=1  
↓sum=1

# OTHER RECENT CAM COMPUTING WORK: GRAPH PROCESSING

- Scatter Gather Apply operations mapped
- CAM + Crossbar architecture efficiently maps sparse vector matrix multiplication

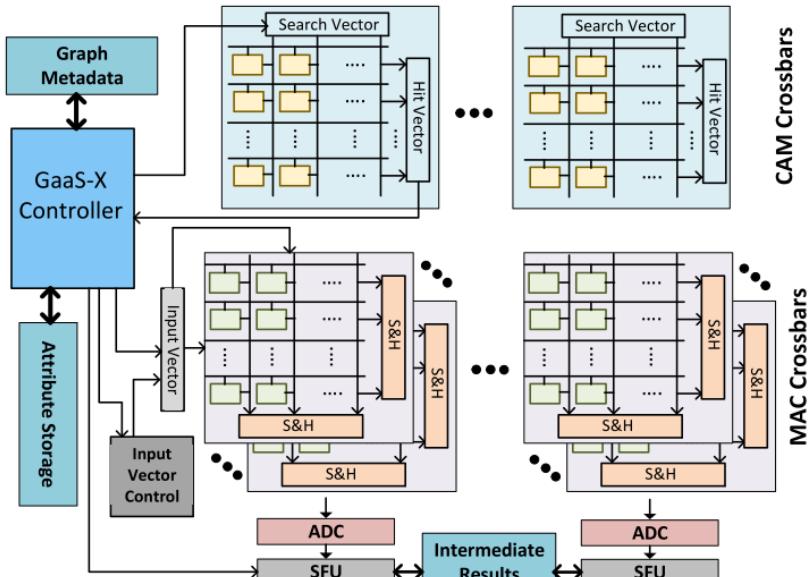


Fig. 6. GaaS-X architecture overview

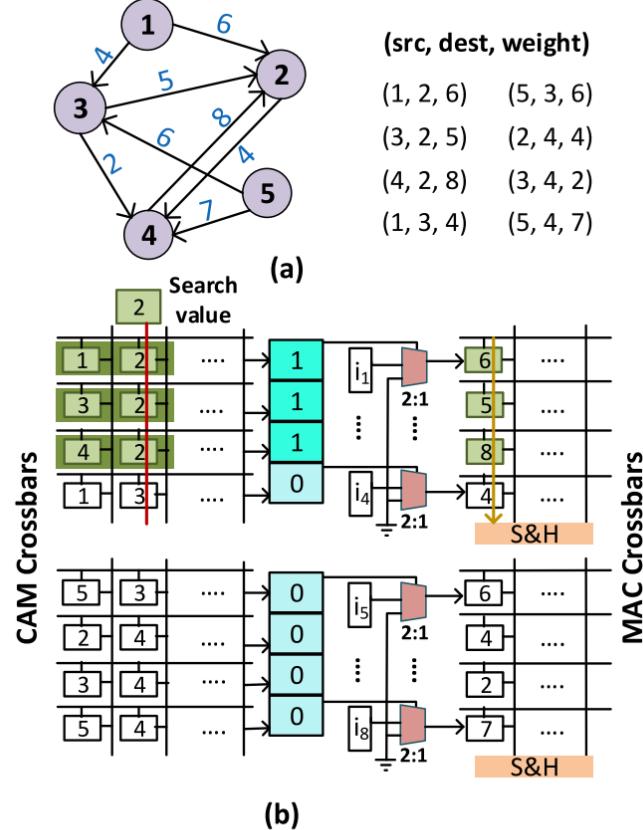


Fig. 7. GaaS-X computing operation

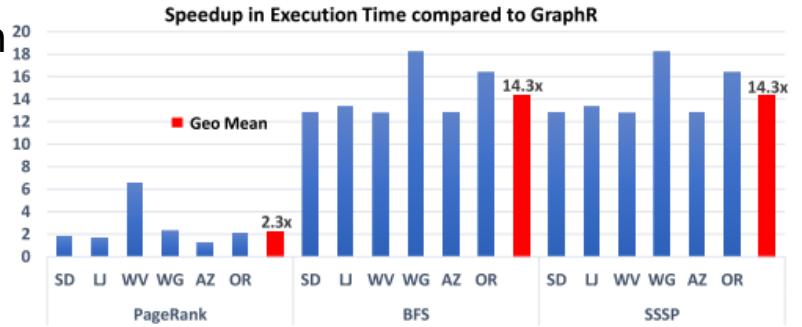


Fig. 11. Speed up in execution time of GaaS-X accelerator for various graph datasets compared to GraphR [33] accelerator.

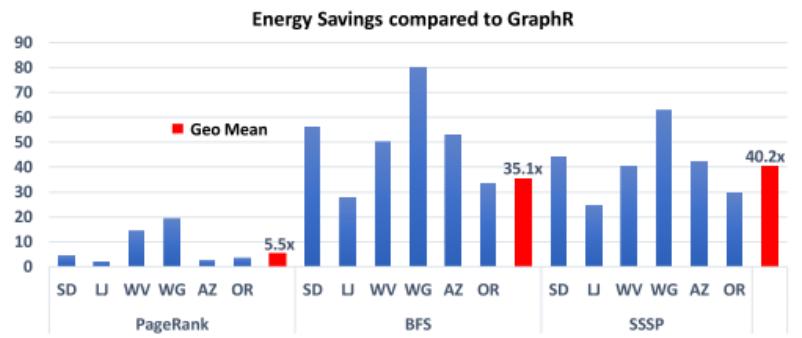


Fig. 12. Energy savings of GaaS-X accelerator for various graph datasets compared to GraphR [33] accelerator.

## SUMMARY

---

- A “new” computing primitive Content Addressable Memory (CAMs): analog and digital CAMs for applications in Security, Genomics, Decision Trees, and more
- CAMs can map a diverse range of computing models
  - Finite State Machines (Reg Ex matching, Edit Distance)
  - Hamming distance calculation, Locality Sensitive Hashing
  - Associative Processing
  - Graph Analytics
- In-memory analog computing platform using non-volatile memristors integrated with CMOS, low current memristors down to 25nm, tuned for CAM properties with application demonstrations
- Analog CAM circuit provides new opportunities
  - ‘Almost matching’ capabilities or range matching



# Thank you

## Hewlett Packard Labs

John Paul Strachan  
Can Li  
Xia Sheng  
Martin Foltin  
Lenny Kiyama  
Xuem Li  
Jim Ignowski  
Rob Wessel  
Jacqui Ingemi  
Cullen Bash

## Program support



Jeffrey Weinschenk (IARPA)  
Richart Slusher  
Karl Roenigk  
John Daly (ARO)  
Chad Meiners (MIT LL)  
Matthew Hardy (LTS)