

Content Addressable Memories and Associative Processors

Asif Ali Khan

Fall Semester 2024

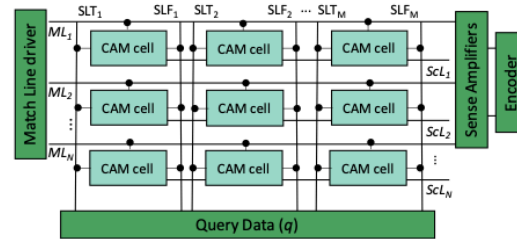
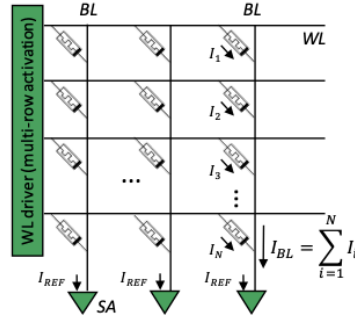
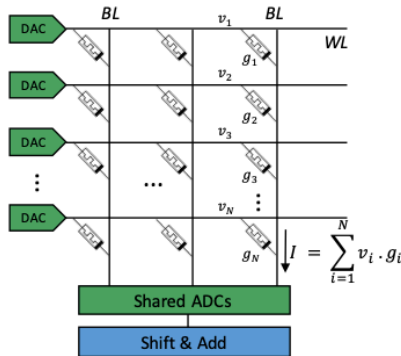
Department of Computer Systems Engineering

UET Peshawar, Pakistan

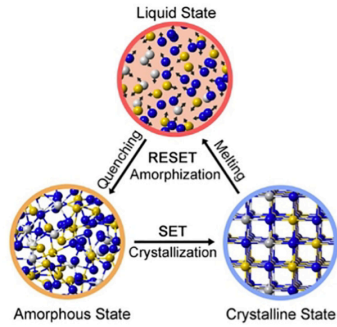
Dec 5, 2024

Recap: Compute-in-memory (CIM)

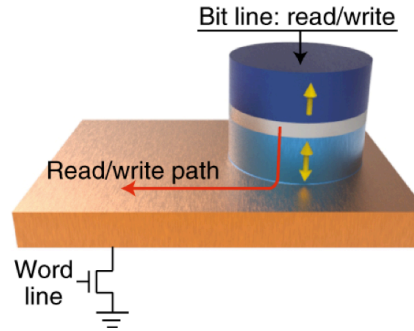
- ❑ The CIM paradigm aims to completely eliminate the data movement
- ❑ The fundamental idea is to exploit the physical properties of the memory devices to perform computations
- ❑ Not every computation can be performed with every technology



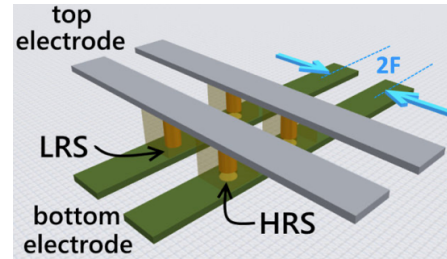
Recap: Emerging nonvolatile memories (NVMs)



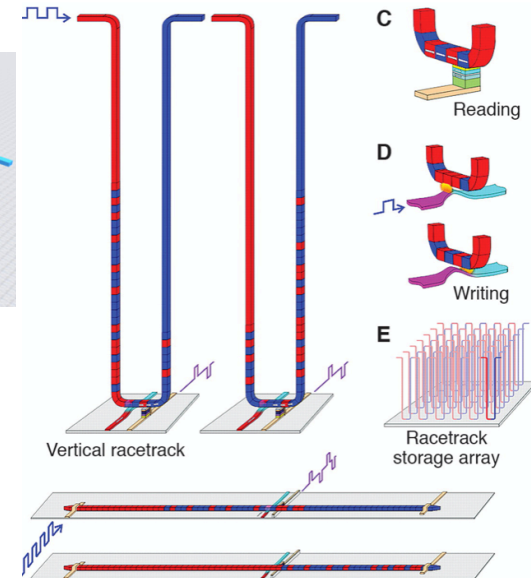
Zhang et al., 2020



G. Yu, 2020



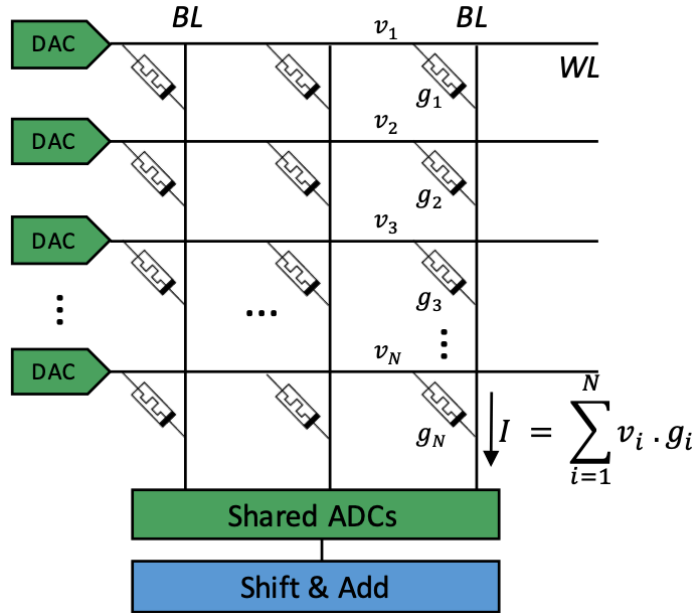
Lim et al, 2015



Parkin et al, 2008

- ❑ Other options: FeFET, FeRAM etc
- ❑ Each technology has its strengths and challenges
- ❑ PCM and MRAM receive a lot of traction in industry

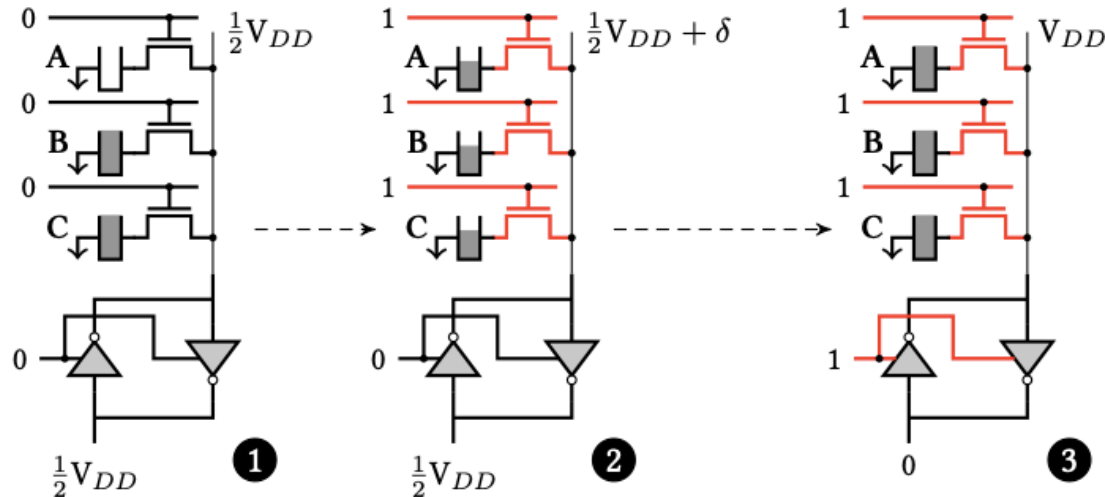
Recap: CIM crossbar



- ❑ Program one operand into memristors devices (conductance)
- ❑ Enable all wordlines simultaneously and apply another operand as input
- ❑ The accumulated current at the bitlines using Kirchoff's law produces the outcome of dot product
- ❑ Analog domain computation – results are approximate

Recap: CIM-logic using DRAM

- ❑ Simultaneous activation of three rows in a DRAM array results in bit-wise majority
 - ❑ At least two cells have to be one, for the output to be one
 - ❑ The operation is called *triple row activation (TRA)*



Recap: CIM-logic using DRAM

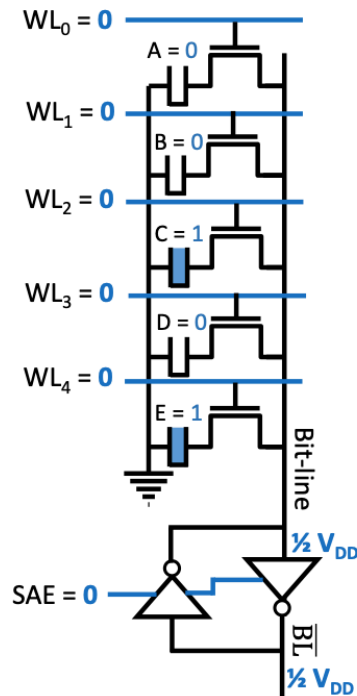
- ❑ Lets say A , B , C represent the state of the three cells
- ❑ The final state of the bitline is: $AB + AC + BC$
- ❑ This can be rewritten as: $C \cdot (A+B) + C' \cdot (AB)$
- ❑ By controlling C , we can implement both AND and OR operations
 - ❑ $C=1$, to implement OR and $C=0$, to implement AND operation
- ❑ Important: TRA destroys contents of the involved cells
 - ❑ Contents need to be copied to a different place first, if important/needed

Recap: In-DRAM addition

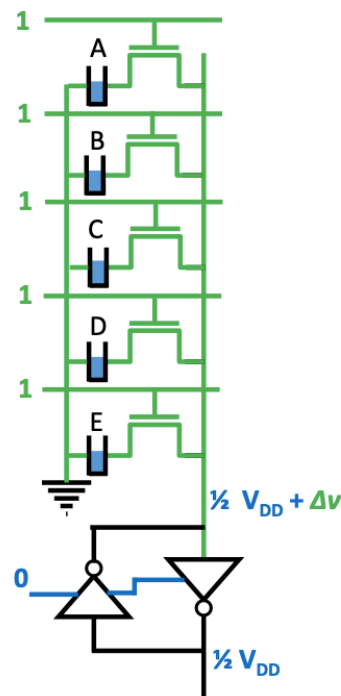
$$C_{out} = Majority(A, B, C_{in})$$

$$Sum = Majority(A, B, C_{in}, \overline{C_{out}}, \overline{C_{out}})$$

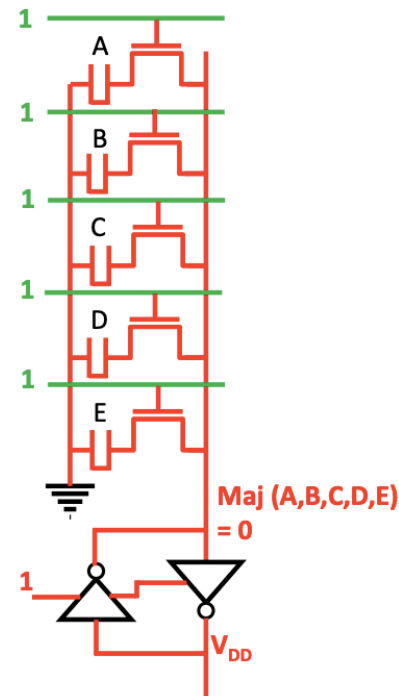
- Note that the operations are bit-serial
- ... but word-parallel
- Particularly useful for bulk additions on low-precision numbers



(1) Initial state



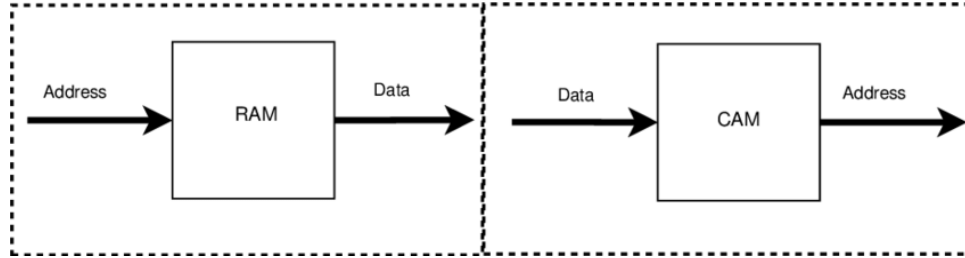
(2) Enable WLS



(3) Enable Sense amp

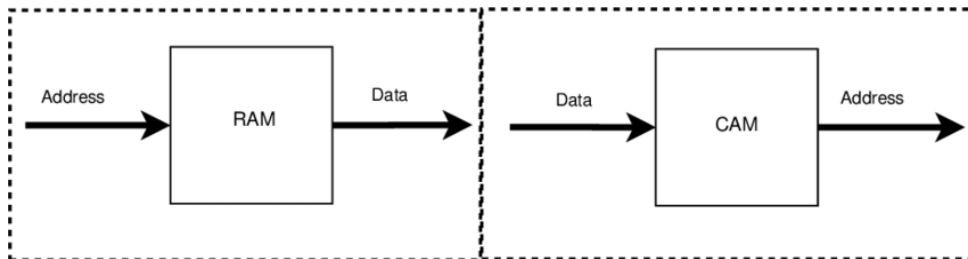
Content-addressable memory (CAM)

- ❑ Computer memory that allows data to be accessed based on content rather than specific addresses



Content-addressable memory (CAM)

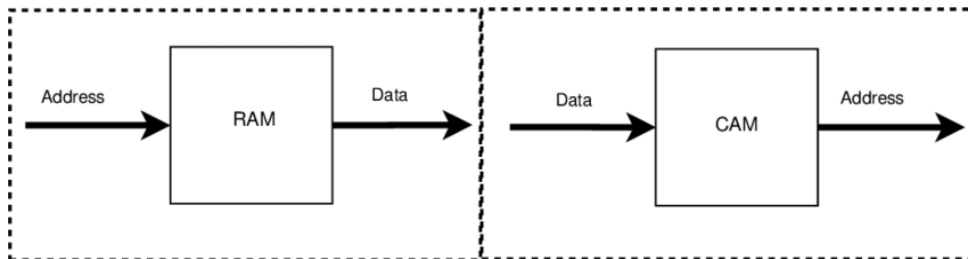
- ❑ Computer memory that allows data to be accessed based on content rather than specific addresses



- ❑ Search entire memory for matching content word
 - ❑ Outputs locations of matching content

Content-addressable memory (CAM) or Associative Memory

- ❑ Computer memory that allows data to be accessed based on content rather than specific addresses



- ❑ Search entire memory for matching content word
 - ❑ Outputs locations of matching content
- ❑ Typical use-cases: Database searches, packet forwarding in networking, ML, etc.

CAM types

- ❑ CAM array can be:
 - ❑ Binary (BCAM)
 - ❑ Ternary (TCAM)

CAM types

- ❑ CAM array can be:
 - ❑ Binary (BCAM)
 - ❑ Ternary (TCAM)
- ❑ BCAMs only allow for exact match

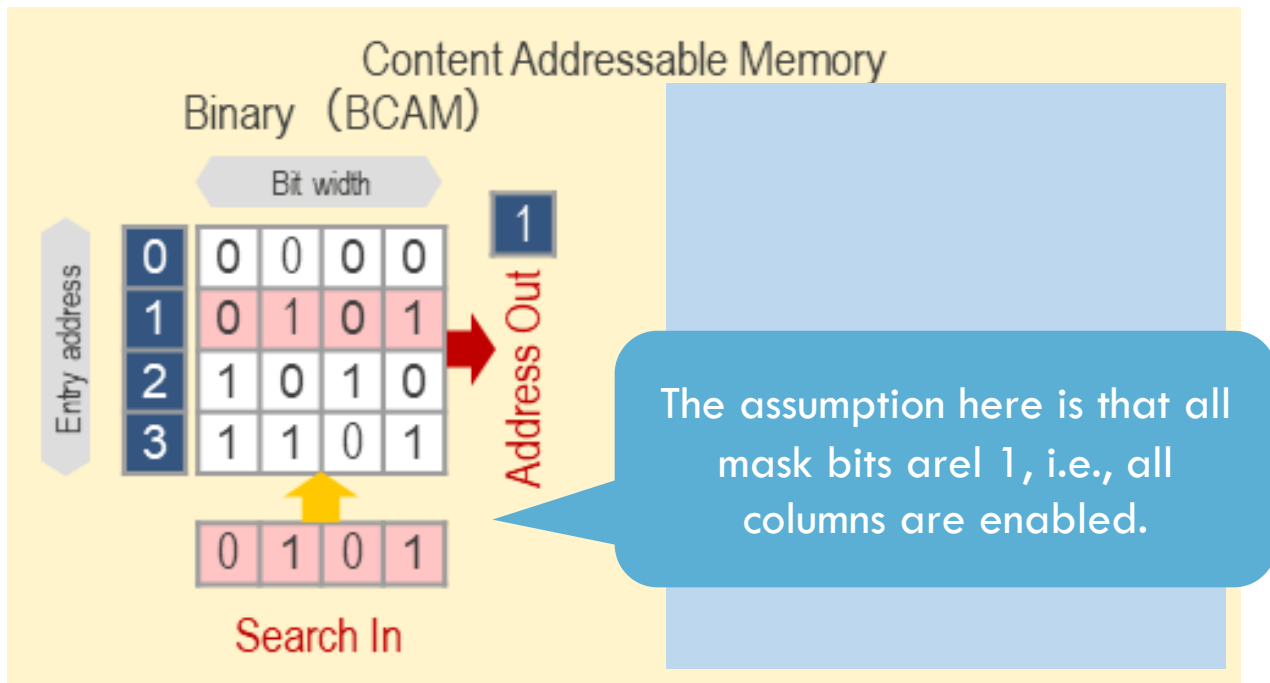
CAM types

- ❑ CAM array can be:
 - ❑ Binary (BCAM)
 - ❑ Ternary (TCAM)
- ❑ BCAMs only allow for exact match
- ❑ In TCAMs, a bit can be in `set`, `reset`, `don't care (x)` state

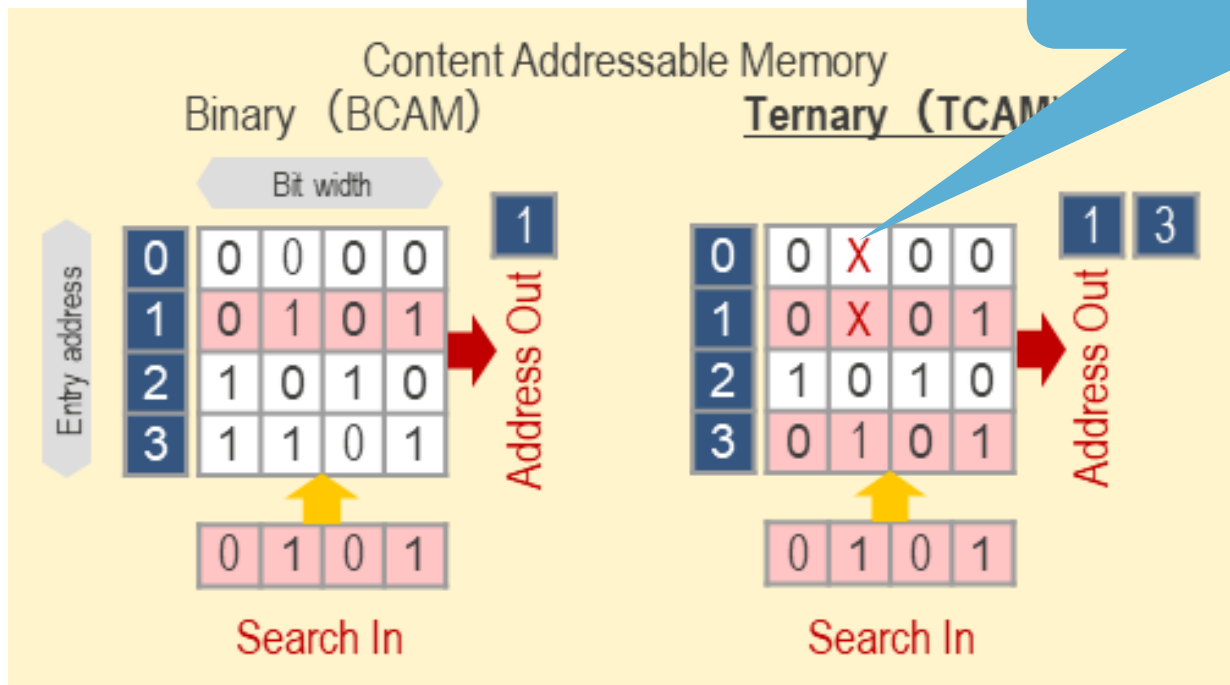
CAM types

- ❑ CAM array can be:
 - ❑ Binary (BCAM)
 - ❑ Ternary (TCAM)
- ❑ BCAMs only allow for exact match
- ❑ In TCAMs, a bit can be in `set`, `reset`, `don't care (x)` state
- ❑ This opens the possibility to do approx. and range based searches in TCAMs

Background: Search with BCAM



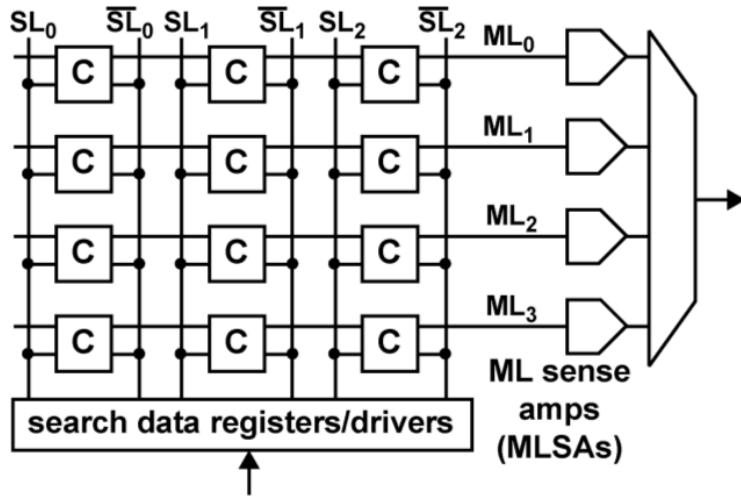
Background: Search with TCAM



The don't care bit results in a match

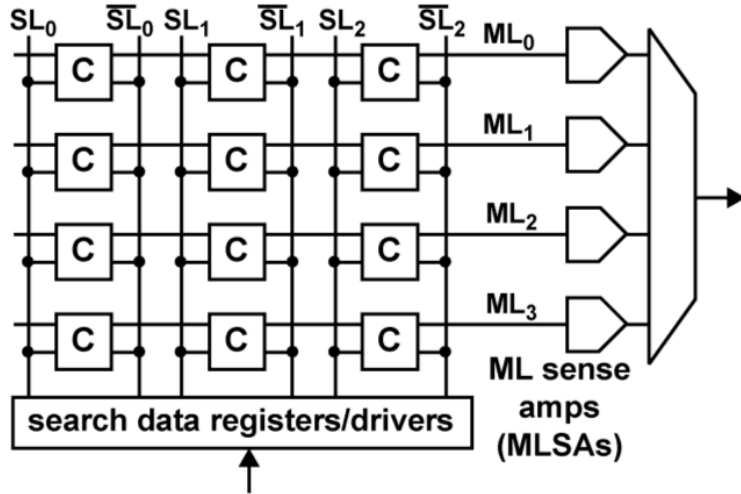
CAM architecture and functionality

- Initially the data to be searched is stored into the search data reg.

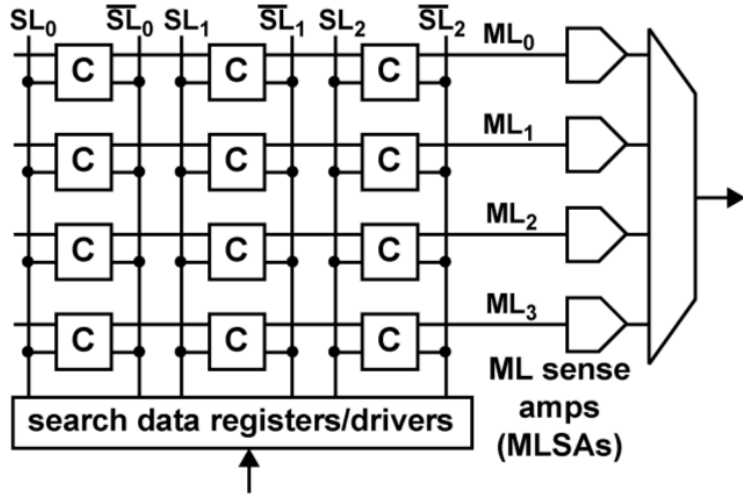


CAM architecture and functionality

- Initially the data to be searched is stored into the search data reg.
- All matchlines (MLs) are precharged to the match state (high)

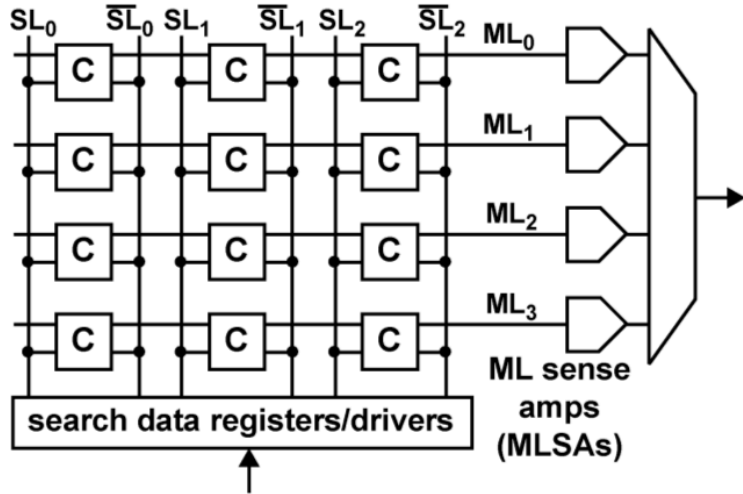


CAM architecture and functionality



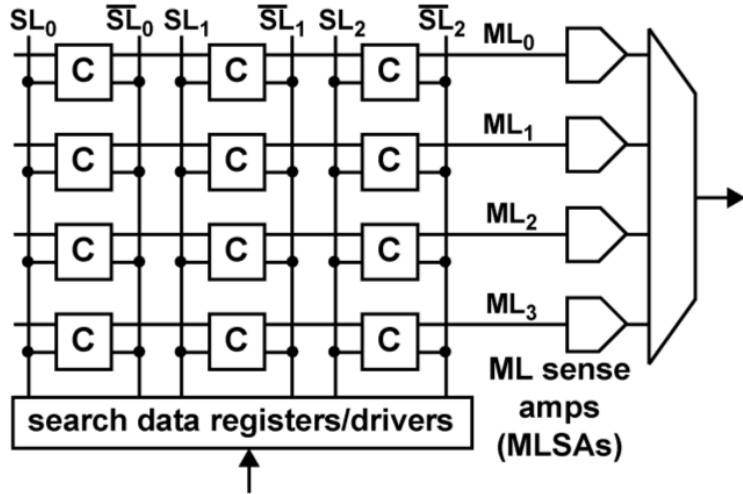
- Initially the data to be searched is stored into the search data reg.
- All matchlines (MLs) are precharged to the match state (high)
- The searchline (SL) drivers broadcast the search word onto the differential SLs

CAM architecture and functionality



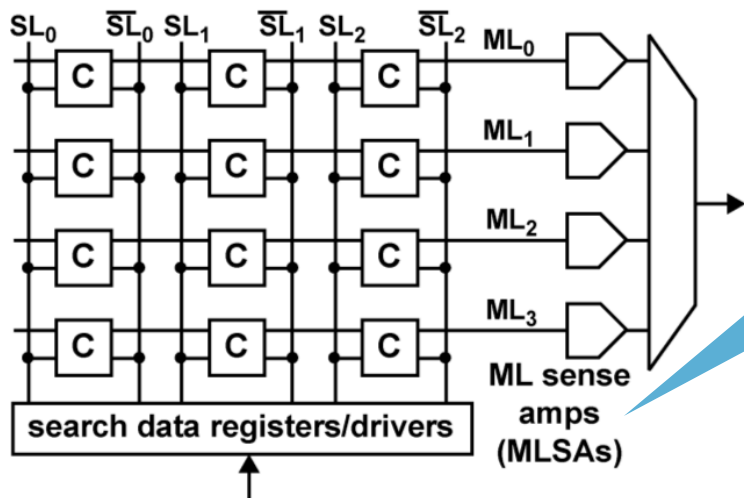
- Initially the data to be searched is stored into the search data reg.
- All matchlines (MLs) are precharged to the match state (high)
- The searchline (SL) drivers broadcast the search word onto the differential SLs
- Each CAM cell compares its stored value to the value on the SLs

CAM architecture and functionality



- Initially the data to be searched is stored into the search data reg.
- All matchlines (MLs) are precharged to the match state (high)
- The searchline (SL) drivers broadcast the search word onto the differential SLs
- Each CAM cell compares its stored value to the value on the SLs
- Only MLs where all bits match remain in high state, all others discharge to ground

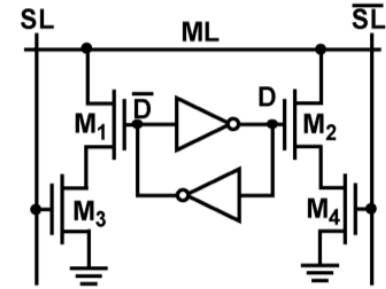
CAM architecture and functionality



- Initially the data to be searched is stored into the search data registers
- The matchline sense amplifier (MLSA) finally sense the match/mismatch state
- The matchline (SL) drivers broadcast the search word onto the differential SLs
- Each CAM cell compares its stored value to the value on the SLs
- Only MLs where all bits match remain in high state, all others discharge to ground

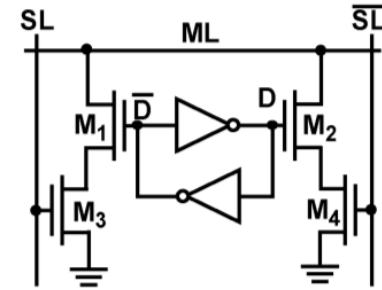
CAM architecture and functionality

- ❑ A CAM cells both stores data and performs the XOR-like comparison



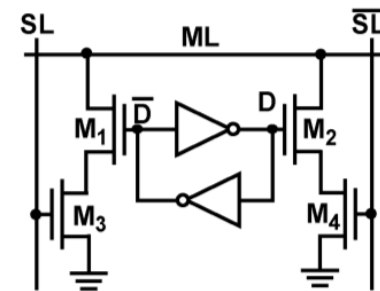
CAM architecture and functionality

- ❑ A CAM cells both stores data and performs the XOR-like comparison
- ❑ Consider the SRAM-based CAM cell in the figure
 - ❑ Does not show the bitlines and wordlines for simplicity
- ❑ The four comparison transistors, M_1 through M_4 , performs the matching and implements the pull-down path



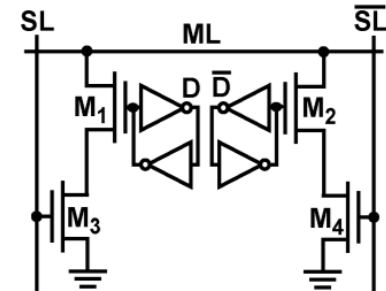
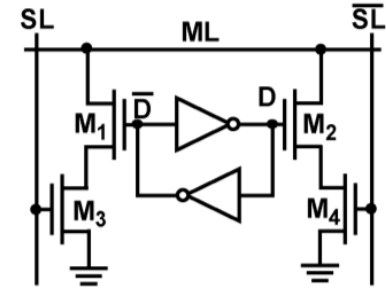
CAM architecture and functionality

- ❑ A CAM cell both stores data and performs the XOR-like comparison
- ❑ Consider the SRAM-based CAM cell in the figure
 - ❑ Does not show the bitlines and wordlines for simplicity
- ❑ The four comparison transistors, M_1 through M_4 , performs the matching and implements the pull-down path
- ❑ A match of SL and D disables both pull down path

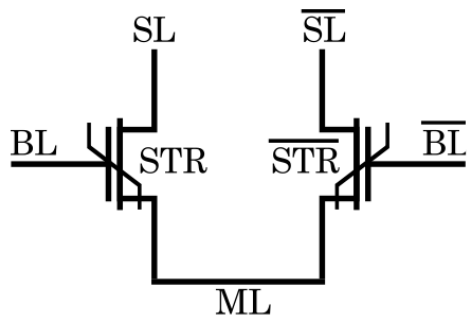


CAM architecture and functionality

- ❑ A CAM cell both stores data and performs the XOR-like comparison
- ❑ Consider the SRAM-based CAM cell in the figure
 - ❑ Does not show the bitlines and wordlines for simplicity
- ❑ The four comparison transistors, M_1 through M_4 , performs the matching and implements the pull-down path
- ❑ A match of SL and D disables both pull down path
- ❑ For ternary cells, we need to store two bits (3 states in total)



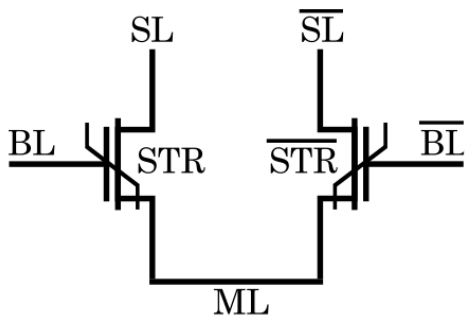
NVM-based CAMs (e.g., FeFET based)



SL (V_D)	STR (V_{PGM})	ML (I_D)
0 V	-5 V	10.9 nA
0 V	+5 V	3.21 pA
50 mV	-5 V	2.79 pA
50 mV	+5 V	13.1 nA

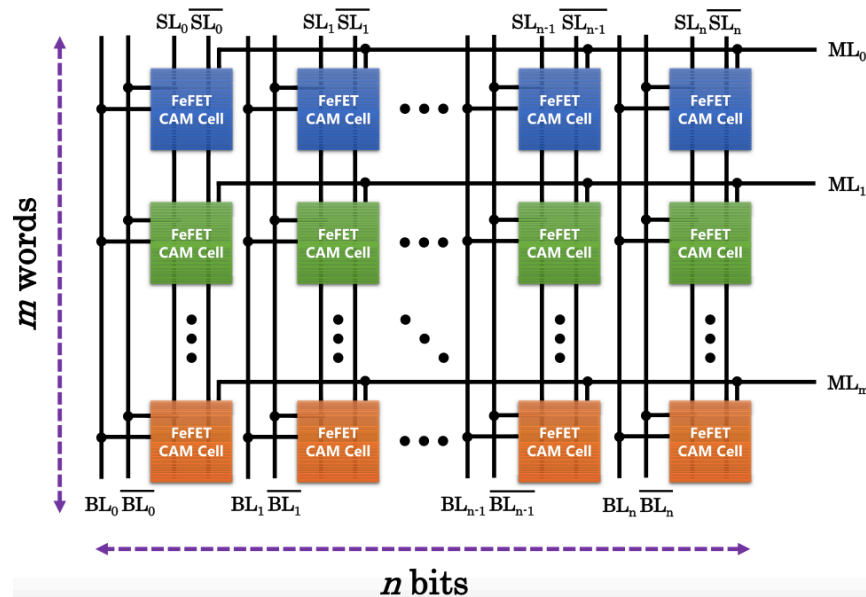
$$ML = XNOR(SL, STR)$$

NVM-based CAMs (e.g., FeFET based)



SL (V_D)	STR (V_{PGM})	ML (I_D)
0 V	-5 V	10.9 nA
0 V	+5 V	3.21 pA
50 mV	-5 V	2.79 pA
50 mV	+5 V	13.1 nA

$$ML = XNOR(SL, STR)$$

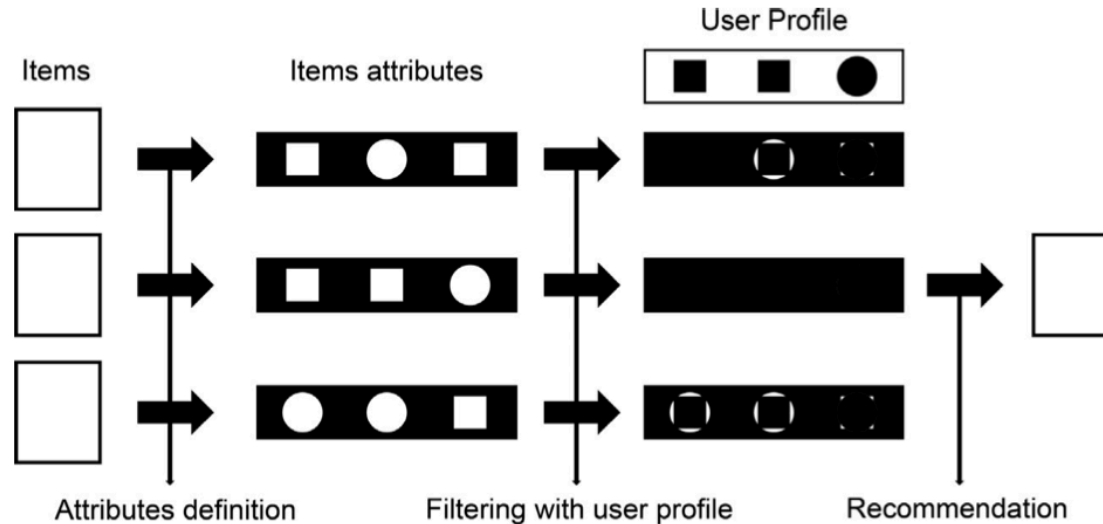


CAM applications

- ❑ CAMs can be used to implement applications from different domains

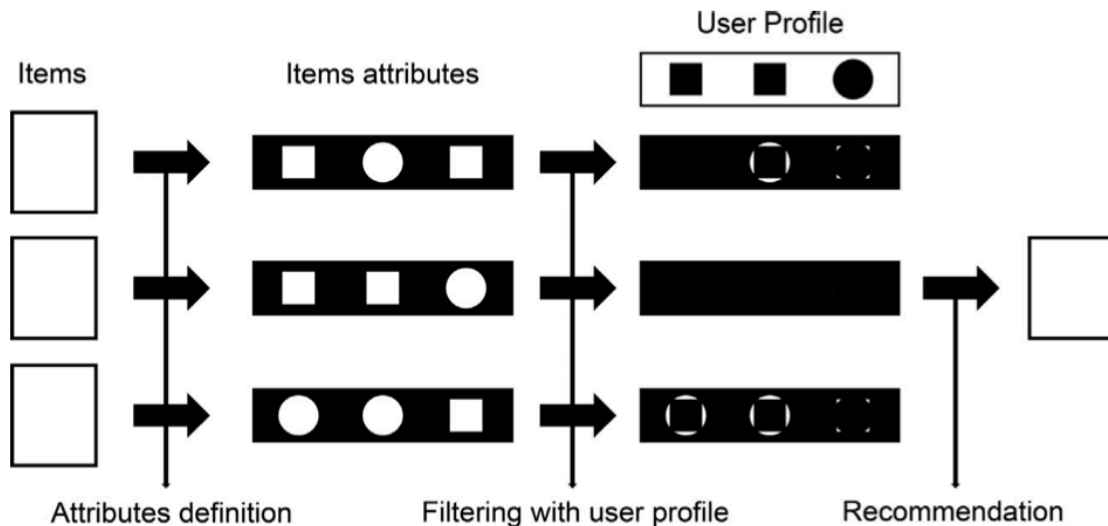
CAM applications

- ❑ CAMs can be used to implement applications from different domains



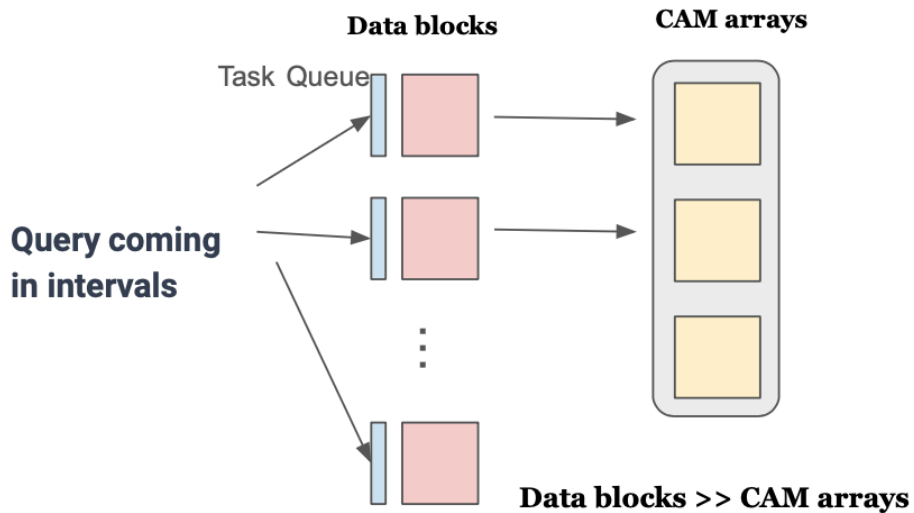
CAM applications

- ❑ CAMs can be used to implement applications from different domains



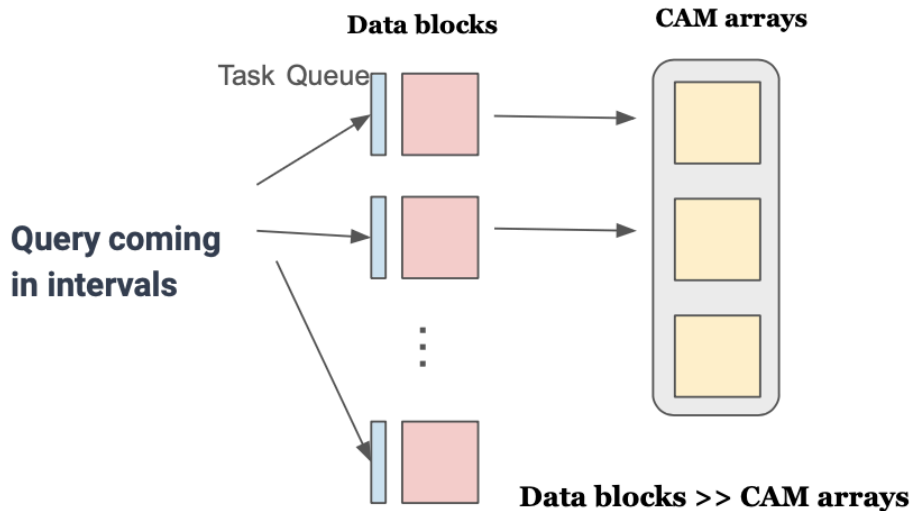
- ❑ Other applications include DNA comparison, K-NNs etc.

Application mapping and scheduling on CAMs

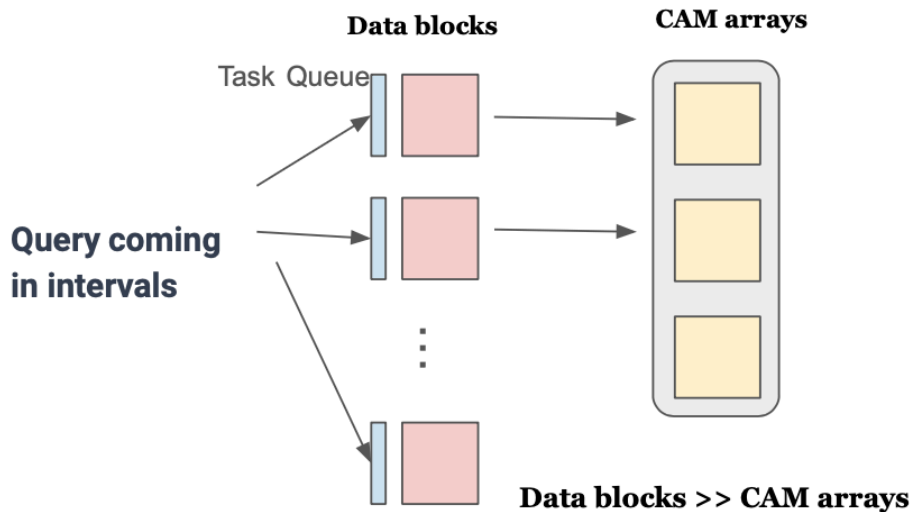


Application mapping and scheduling on CAMs

- The objective usually is to optimize for latency, throughput, and/or energy

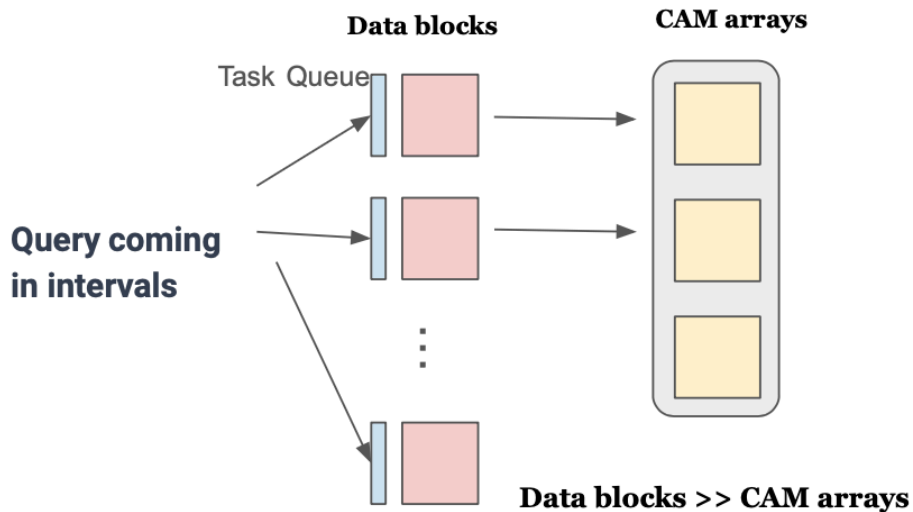


Application mapping and scheduling on CAMs



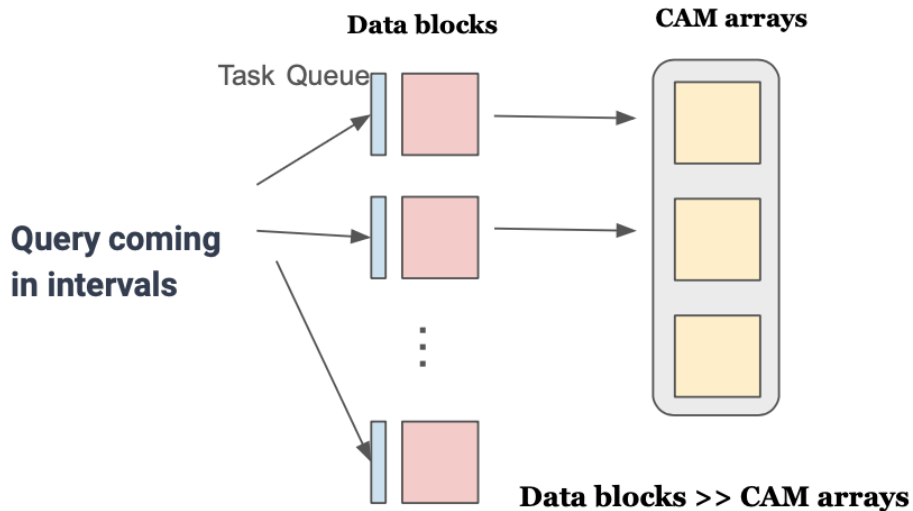
- ❑ The objective usually is to optimize for latency, throughput, and/or energy
- ❑ The best-mapping/scheduling should minimize the number of reloads to the CAM arrays

Application mapping and scheduling on CAMs



- ❑ The objective usually is to optimize for latency, throughput, and/or energy
- ❑ The best-mapping/scheduling should minimize the number of reloads to the CAM arrays
- ❑ In the worst case, every new query starts when the previous one finishes

Application mapping and scheduling on CAMs



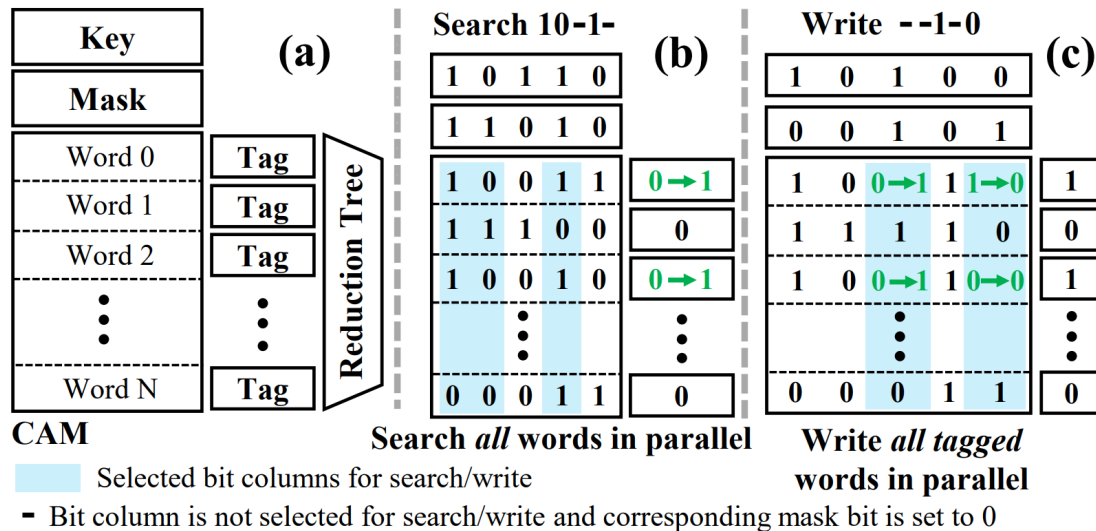
- ❑ The objective usually is to optimize for latency, throughput, and/or energy
- ❑ The best-mapping/scheduling should minimize the number of reloads to the CAM arrays
- ❑ In the worst case, every new query starts when the previous one finishes
- ❑ Efficient handling of intermediate results is the trick here

Associative processors

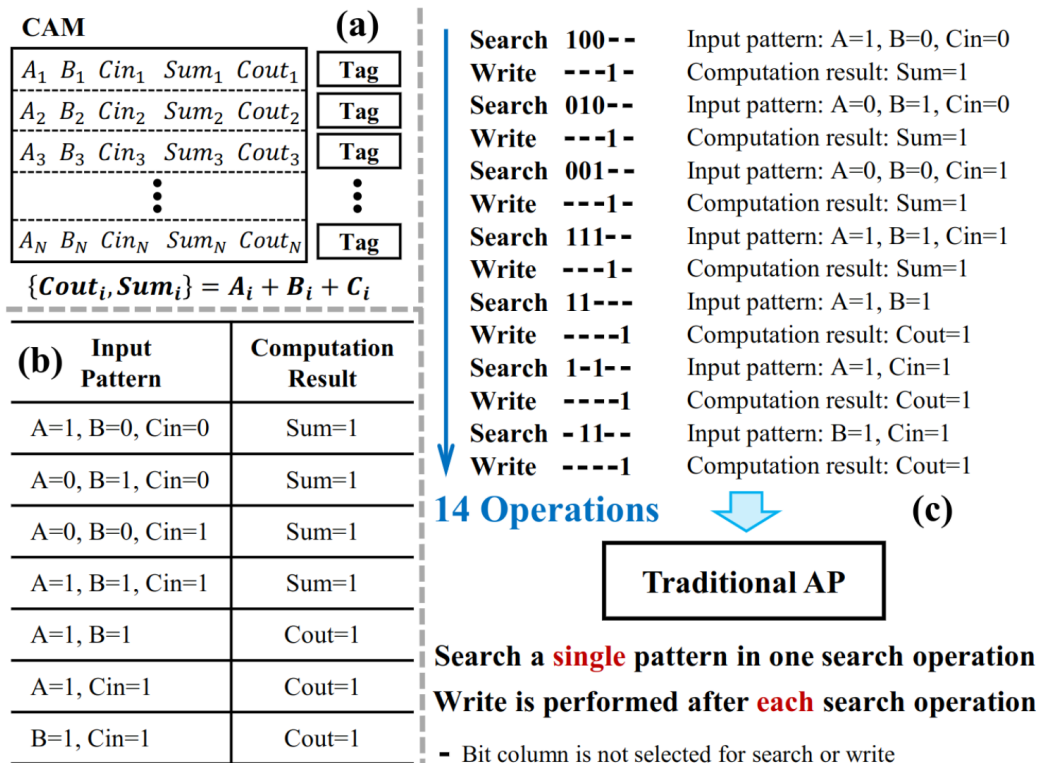
- ❑ Perform computations using CAM

Associative processors

- ❑ Perform computations using CAM
- ❑ The idea is to replace computations using search-and-write operations



Associative processors (Full adder)



Thank you!
asif.ali@uetpeshawar.edu.pk