

# **Embedded LF1 Notes:**

## **Understanding the DRAM Cell**

### **1. Structure of a DRAM Cell**

A DRAM cell is the fundamental storage unit of Dynamic Random Access Memory, consisting of:

- **Storage Capacitor:** Holds data in the form of electrical charge.
  - **Logic 1:** The capacitor is charged.
  - **Logic 0:** The capacitor is discharged.
  - **Key Property:** Capacitors are imperfect and gradually leak charge, necessitating **refresh operations** to preserve data integrity.
- **Access Transistor:** Acts as a gate, controlling the flow of charge between the capacitor and the connected bitline.

### **2. Hierarchical Organization**

To manage and access billions of DRAM cells efficiently:

- **Subarrays:** A 2D grid of DRAM cells.
  - Each cell connects to a **wordline** (horizontal) and a **bitline** (vertical).
  - Bitlines connect to **sense amplifiers** for data read and write operations.
- **Banks, Chips, and Modules:** Subarrays are grouped into banks, banks into chips, and chips into DRAM modules for scalability.

---

## **DRAM Operations: The Cycle of Read, Write, and Refresh**

### **Step 1: Precharge Phase (Preparing Bitlines)**

- **Goal:** Initialize bitlines to a neutral reference voltage ( $V_{DD}/2$ ).
  - **Why:** Precharging ensures a baseline state for detecting minute voltage changes during a read operation.
  - **Outcome:** Bitlines are equilibrated, and the DRAM is ready for the next access.
- 

### **Step 2: Reading Data from a DRAM Cell**

1. **Activation (ACT Command):**

- The wordline of the target row is activated.
  - This turns on the access transistors in the row, connecting the storage capacitors of the selected cells to their respective bitlines.
2. **Charge Sharing:**
- Charge from the capacitor flows to the bitline, or vice versa, depending on the stored data.
  - This causes a small voltage deviation on the bitline:
    - Charged capacitor (Logic 1): Bitline voltage rises slightly ( $V_{DD}/2 + \Delta$ ).
    - Discharged capacitor (Logic 0): Bitline voltage drops slightly ( $V_{DD}/2 - \Delta$ ).
3. **Sense Amplification:**
- The sense amplifier detects the tiny voltage change ( $\Delta$ ).
  - It amplifies the voltage:
    - Logic 1: Bitline voltage is driven to  $V_{DD}$ .
    - Logic 0: Bitline voltage is driven to 0V.
4. **Outcome:**
- The data is latched into the sense amplifier and row buffer.
5. **Data Output:**
- The memory controller issues a read command to retrieve data from the row buffer.
6. **Recharge (Precharge Command):**
- After the read operation, the bitlines are reset to the neutral state ( $V_{DD}/2$ ), preparing the DRAM for the next access.
- 

### Step 3: Writing Data to a DRAM Cell

1. **Activation:**
    - The wordline is activated, connecting the capacitor to the bitline.
  2. **Data Transfer to Bitline:**
    - The memory controller sets the desired logic level on the bitline:
      - Logic 1: Bitline is driven to  $V_{DD}$ , charging the capacitor.
      - Logic 0: Bitline is driven to 0V, discharging the capacitor.
  3. **Charge or Discharge Capacitor:**
    - The bitline's voltage determines whether the capacitor is charged (logic 1) or discharged (logic 0), storing the desired data.
  4. **Precharge:**
    - The bitline is restored to the reference voltage ( $V_{DD}/2$ ), completing the write operation.
- 

### Periodic Refresh Operations

- **Why Needed:** Capacitors leak charge over time, risking data loss.

- **How It Works:** The DRAM controller periodically reactivates each row, reading and rewriting the data to restore charge levels.
- 

## Key Concepts to Understand DRAM Behavior

1. **Equilibrium State:**
  - Bitlines are held at  $V_{DD}/2$  to prepare for sensing small voltage differences during read operations.
2. **Stable States:**
  - $V_{DD}$  (logic 1) and 0V (logic 0) are the stable voltage levels representing definitive data.
3. **Unstable (Transitional) States:**
  - Intermediate voltages ( $V_{DD}/2 + \Delta$  or  $V_{DD}/2 - \Delta$ ) occur during charge sharing but are quickly resolved by sense amplifiers.

## Row Cloning in DRAM:

---

### What is Row Cloning?

Row cloning is a technique used within DRAM (Dynamic Random Access Memory) to copy data from one row (source row) to another row (destination row) **without moving the data outside the memory chip**. It allows data duplication to happen entirely inside the DRAM array, which saves time and energy.

---

### Why is Row Cloning Done?

1. **Avoiding Data Movement:**
  - Normally, copying data involves reading it out of DRAM into the processor, and then writing it back into DRAM. This back-and-forth movement is **slow and energy-intensive**.
  - Row cloning eliminates this movement by performing the copying directly inside DRAM.
2. **Improved Efficiency for Algorithms:**
  - Certain algorithms (e.g., parallel processing) require multiple copies of data. Row cloning allows these copies to be made faster, optimizing performance.
3. **Building Block for Compute-in-Memory (CIM):**

- Row cloning is a key step in advanced techniques like CIM, where computations happen directly within memory instead of relying heavily on the processor.
- 

## How Does Row Cloning Work?

Row cloning takes advantage of how DRAM cells behave during standard operations like **activation** and **precharge**. Here's the simplified step-by-step process:

1. **Activate the Source Row:**
    - The source row (containing the data to be copied) is activated.
    - This connects its cells (capacitors) to the bitlines, and the data is sensed and amplified by the sense amplifiers.
    - The amplified data is temporarily held in the **row buffer**.
  2. **Activate the Destination Row:**
    - The destination row (where the data is to be copied) is activated immediately after the source row.
    - This connects the destination row's cells to the same bitlines that are now holding stable logic levels based on the source row's data.
  3. **Data Overwrite:**
    - The stable data from the bitlines overwrites the contents of the destination row. Essentially, the data from the source row is copied directly into the destination row.
  4. **Precharge:**
    - Finally, the bitlines are reset to their neutral state to prepare for the next operation.
- 

## Key Points to Remember:

- **No Processor Involvement:** The processor is not used in this process, making it faster and more energy-efficient.
- **Shared Bitlines:** For row cloning to work, the source and destination rows must share the same bitlines within the DRAM subarray.
- **Use of Row Buffer:** The row buffer acts as a temporary holding space for the data being copied.
- **Foundation for Advanced Operations:** Row cloning is not just a copying tool but also the basis for more complex in-memory computations like AND/OR logic.

## Triple Row Activation in DRAM:

---

### What is Triple Row Activation?

Triple row activation is a technique used in DRAM to compute the **majority function**—a Boolean logic operation—directly within the memory array. It allows DRAM to determine the dominant logic value (1 or 0) among three memory cells connected to the same bitline. This process is a foundational step for performing computations directly in memory, a concept known as **in-DRAM processing** or **compute-in-memory (CIM)**.

---

### Why is Triple Row Activation Done?

1. **Efficient Boolean Computation:**
    - Traditional computation requires data to be transferred to the processor for logic operations. Triple row activation enables basic computations like the majority function directly in DRAM, reducing the need for data movement.
  2. **Energy and Time Savings:**
    - Since data is not transferred outside DRAM, this approach saves energy and speeds up operations.
  3. **Foundation for Complex Logic:**
    - The majority function is a building block for implementing other Boolean logic operations (AND, OR, etc.) directly in memory.
- 

### How Does Triple Row Activation Work?

The process leverages DRAM's architecture, specifically the shared bitlines and sense amplifiers. Here's a step-by-step breakdown:

---

#### Step 1: Bitline Precharge

- Before any operation, the bitlines are precharged to a neutral voltage ( $V_{DD}/2$ )
  - This ensures the bitlines are ready to sense and amplify any charge changes during the activation.
- 

#### Step 2: Triple Activation Command

- The **wordlines** of three specific rows (let's call them A, B, and C) are activated simultaneously.
  - This connects the storage capacitors of these rows (representing their stored data) to their respective bitlines.
- 

### Step 3: Charge Sharing and Resultant Voltage

- The charge from the capacitors in rows A, B, and C interacts with the bitlines. The voltage on the bitlines changes depending on the number of capacitors holding a charge (logic 1) versus those without a charge (logic 0).
- 

### Step 4: Majority Function Determination

- The **sense amplifier** detects the resulting bitline voltage and amplifies it:
    - **Case 1:** If two or three cells hold logic 1, the bitline voltage rises above  $V_{DD}/2$ . The sense amplifier amplifies it to  $V_{DD}$ , representing a majority of 1.
    - **Case 2:** If two or three cells hold logic 0, the bitline voltage drops below  $V_{DD}/2$ . The sense amplifier drives it to 0V, representing a majority of 0.
- 

### Step 5: Write-Back

- The majority value determined by the sense amplifier is written back to the storage capacitors of all three activated rows (A, B, and C).
  - This overwrites their original data with the computed majority value.
- 

### Step 6: Precharge

- A **precharge command** resets the bitlines to  $V_{DD}/2$ , preparing the DRAM for the next operation.

# Explaining AND and OR Operations Using Triple Row Activation (TRA) in DRAM

Let's break down how Triple Row Activation (TRA) can implement **AND** and **OR** logic gates within DRAM by combining insights from both sources.

---

## Understanding the Setup

1. **Three Rows (A, B, and C):**
    - **A and B:** These rows store the **data inputs** for the operation (the values you want to compute on).
    - **C (Control Row):** This row determines whether the operation will perform **AND** or **OR** logic.
  2. **Bitline and Charge Sharing:**
    - All three rows are connected to the same **bitline** for each column of memory.
    - When TRA is performed, the charge from these three rows is combined on the bitline, which produces the **final logic result**.
  3. **Equation:**

The final result of TRA is determined by:

$$\text{Bitline Output} = C \cdot (A+B) + C' \cdot (A \cdot B)$$
    - This equation allows us to implement both **AND** and **OR** gates by simply controlling the value of C.
- 

## How the AND Gate Works

1. **Set C = 0**
  - This ensures the first term  $C \cdot (A + B)$  becomes 0.
  - Reason: Since  $C = 0$ , multiplying by 0 nullifies the term.
2. **Simplify the Equation**
  - The equation reduces to:
$$\text{Bitline Output} = C' \cdot (A \cdot B)$$
  - With  $C' = 1$ , the equation simplifies further to:
$$\text{Bitline Output} = A \cdot B$$
3. **Perform TRA (Transfer Row Activate)**
  - Activate rows A, B, and C together.
  - The bitline stabilizes at logic 1 only if both A and B contain logic 1 in their respective cells.
4. **Result**
  - The final state of the bitline reflects the behavior of an AND gate:
$$\text{Output} = 1 \text{ only if both } A = 1 \text{ and } B = 1$$

---

## How the OR Gate Works

1. **Set  $C = 1$** 
    - This ensures the second term  $C' \cdot (A \cdot B)$  becomes 0.
    - Reason: Since  $C' = 0$ , multiplying by 0 nullifies the term.
  2. **Simplify the Equation**
    - The equation reduces to:  
Bitline Output =  $C \cdot (A + B)$
    - With  $C = 1$ , the equation simplifies further to:  
Bitline Output =  $A + B$
  3. **Perform TRA (Transfer Row Activate)**
    - Activate rows A, B, and C together.
    - The bitline stabilizes at logic 1 if either A, B, or both contain logic 1 in their respective cells.
  4. **Result**
    - The final state of the bitline reflects the behavior of an OR gate:  
Output = 1 if  $A = 1$  or  $B = 1$
- 

## What is Ambit?

Ambit is a technology that allows a DRAM chip (the main memory in your computer) to perform certain logic operations (like AND, OR, and NOT) directly inside the memory. This reduces the need to transfer data between the memory and the processor, saving time and energy.

---

## How is Ambit Organized?

Ambit divides the DRAM memory into three main groups:

- **B-group (Computation Rows):**

These rows are specifically used for performing logic operations.  
Example: 8 rows are grouped together for this purpose.
  - **C-group (Constant Rows):**

Two rows store fixed values:

    - C0: All 0s.
    - C1: All 1s.

These rows are essential for implementing the AND, OR, and NOT operations.
  - **D-group (Data Storage Rows):**

The rest of the rows are for regular data storage, like in a conventional DRAM chip.
-



## How Does Ambit Work?

- **Triple Row Activation (TRA):**

Ambit can activate three rows simultaneously.

When three rows share the same bitline, the logic of the three rows is combined based on a majority function:

- If at least 2 out of 3 rows store 1, the result is 1.
  - Otherwise, the result is 0.
- 

## How is Ambit Different from Regular DRAM?

- **Row Decoders:**

Ambit has special row decoders that can activate three rows at once, unlike traditional DRAM, which activates one row at a time.

- **Bitwise Group:**

Ambit reserves certain rows (the B-group) for computations, ensuring efficient logic operations.

- **Minimal Modifications:**

These changes add only a small overhead (less than 1%) to the DRAM chip design, making it practical to implement.

---

## Advantages of Ambit

1. **Faster Logic Operations:**

Logic is performed directly in memory, avoiding the need to send data back and forth between memory and the processor.

2. **Energy Efficiency:**

Reduces energy consumption by performing computations within the memory.

3. **Parallel Processing:**

Ambit can perform multiple logic operations simultaneously across different columns, increasing efficiency.

## Limitations of Ambit

1. **Limited Operations:**

It can only perform basic Boolean logic (AND, OR, NOT, XOR). Complex computations aren't supported.

2. **Error Handling:**

Activating multiple rows at once increases the risk of errors. Error correction mechanisms are necessary.

---

## How Does Ambit Perform AND Logic Using DRAM?

### 1. Prepare the Input Data

- Imagine you have two rows of data in memory:
  - Row A (e.g., 1010).
  - Row B (e.g., 1100).These rows hold the binary values you want to compute the AND operation on.

### 2. Set Up a Constant Row

- Ambit has a special row of memory called C0, which is filled with all 0s (e.g., 0000).

### 3. Copy Data to Temporary Rows

- Ambit uses temporary rows (T0, T1, etc.) for computation:
  - Copy Row A to T0.
  - Copy Row B to T1.
  - Copy C0 (all 0s) to T2.At this point:

- T0 = Row A
- T1 = Row B
- T2 = C0 (all 0s)

### 4. Activate the Rows Simultaneously

- Ambit performs a Triple Row Activation (TRA):
  - T0, T1, and T2 are activated at the same time.

### 5. Compute the Result

- When the rows are activated:
  - The memory cells share their charges on the bitlines.
  - The Sense Amplifier calculates the "majority" value:
    - If at least 2 out of the 3 rows have 1, the output is 1.
    - Otherwise, the output is 0.
- Because T2 is all 0s, the result behaves like an AND operation:
  - If both T0 (Row A) and T1 (Row B) have 1, the output is 1.
  - Otherwise, the output is 0.

### 6. Store the Result

- The final result of the AND operation is stored in one of the temporary rows, such as T0, or copied back to a storage row.

---

## How is the OR Operation Different?

- For an OR operation:
  - Instead of C0 (all 0s), Ambit uses C1 (all 1s).
  - The rest of the steps are the same.
  - The result behaves like the OR logic table:
    - If either Row A or Row B has 1, the output is 1.

---

## Key Benefits of Ambit

- **Speed:** Logic operations happen directly inside memory without needing to transfer data to the processor.
- **Energy Efficiency:** Saves energy by reducing data movement.
- **Parallelism:** Many rows in memory can compute logic operations at the same time.

## Key Limitation

- The result overwrites the data in the rows used for computation (like T0, T1, and T2). If you need the original data, you must copy it elsewhere first.
- 

## In-DRAM Addition:

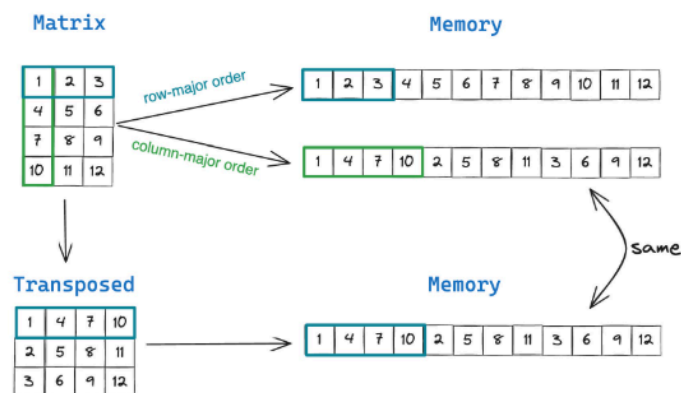
In-DRAM addition enables arithmetic operations, such as addition, to occur directly within the memory (DRAM).

### Key Concepts:

- **Bit-Serial:** The addition is done one bit at a time, starting from the least significant bit.
- **Word-Parallel:** Multiple additions are performed simultaneously across columns in the DRAM array.

### Steps for In-DRAM Addition

1. **Initialization:**
  - **Operands Storage:** Binary numbers (operands) A and B are stored in a column-major layout.



- **Carry Row Initialization:** The carry row  $C$  is set to 0 initially.
  - 2. **Bit Addition (Using the Majority Function):**
    - **Activating Rows:** Rows containing bits of  $A$ ,  $B$ , and  $C$  are activated simultaneously.
    - **Sum Calculation:** The bitline outputs the sum bit ( $S$ ) based on the majority logic of  $A$ ,  $B$ , and  $C$ .
    - **Carry Calculation:** Another majority operation determines the carry-out ( $C_{out}$ ) for the current bit position.
    - **Storing Results:**
      - The sum  $S$  is written to a result row.
      - The carry  $C_{out}$  is stored in the carry row for the next step.
  - 3. **Carry Propagation:**
    - The carry-out becomes the carry-in for the next bit position.
    - Steps in the bit addition process are repeated for all bit positions.
- 

## Example: 2-Bit Addition

Let's add  $A=10$  and  $B=01$  (in binary):

### Column-Major Layout

**Example: 2-Bit Addition**

Let's add  $A = 10$  and  $B = 01$  (in binary):

**Column-Major Layout**

Column	$A_0$	$B_0$	$C_0$ (initial)	$A_1$	$B_1$	$C_1$
Data	0	1	0	1	0	0

### Bit 0 Addition

- Activate rows for  $A_0 = 0$ ,  $B_0 = 1$ , and  $C_0 = 0$ .
- **Sum ( $S_0$ ) Calculation:** Majority of 0, 1, 0 = 1.
- **Carry ( $C_1$ ) Calculation:** Majority for carry = 0.
- **Result:**  $S_0 = 1$ ,  $C_1 = 0$ .

### Bit 1 Addition

- Activate rows for  $A_1 = 1$ ,  $B_1 = 0$ , and  $C_1 = 0$ .
- **Sum ( $S_1$ ) Calculation:** Majority of 1, 0, 0 = 1.
- **Carry ( $C_2$ ) Calculation:** Majority for carry = 0.
- **Result:**  $S_1 = 1$ ,  $C_2 = 0$ .

### Final Output

The result is stored as  $S = 11$  (binary for decimal 3).

---

### Advantages

1. **In-Memory Computation:** Eliminates the need to transfer data between memory and the processor, saving time and energy.
2. **High Parallelism:** Supports simultaneous processing of multiple bit additions across columns.

### Challenges

1. **Bit-Serial Nature:** Carry propagation makes the process sequential, which can slow down operations for large numbers.
2. **Error Probability:** Activating multiple rows simultaneously increases the chances of DRAM errors, requiring robust error correction.

### Purpose of Johnson Counters in DRAM

Johnson counters are used to **implement masked counting operations** within DRAM arrays. These counters operate directly within the memory, making arithmetic operations more efficient by avoiding unnecessary data transfers.

### Key Applications

1. **Masked Counting:** Counting specific patterns or values in memory.

2. **Efficient In-Memory Operations:** Performing operations like shifts and bit manipulations directly in DRAM.

### Mechanism of Johnson Counters

#### Basic Shifting Operation

- Each bit  $b_i$  in the counter is updated based on the values of:
  - The **mask bit  $m$** : Determines whether the operation is active or skipped.
  - The previous bit  $b_{i-1}$ : Used for the shifting mechanism.

#### Logic for Masked Counting

The update logic for each bit  $b_i$  is given by:  $b_i = (b_i \wedge \neg m) \vee (b_{i-1} \wedge m)$  Where:

- $\wedge$  is the logical AND operation.
- $\vee$  is the logical OR operation.
- $\neg m$  is the NOT of the mask bit  $m$ .

## Steps for Masked Counting Using Johnson Counters

### 1. Initialization:

- The mask value  $m$  and the initial bit pattern  $b_i$  are loaded into the DRAM array.
- Rows corresponding to the bits and mask are activated.

### 2. Bit Update:

- For each bit position  $i$ , the new value is calculated based on the logical expression.
- If  $m = 0$ : The bit  $b_i$  remains unchanged.
- If  $m = 1$ : The new value of  $b_i$  is copied from  $b_{i-1}$ .

### 3. Circular Shifting:

- After updating all bits, the last bit  $b_n$  wraps around to the first position  $b_1$ , creating a circular shift.

### 4. Repeat for Multiple Steps:

- The shifting and bit manipulation continue for the desired number of cycles, depending on the counting requirements.



## Example of Masked Counting with a 3-Bit Johnson Counter

Let's consider a 3-bit Johnson counter ( $b_1, b_2, b_3$ ) with an initial pattern of 001 and mask  $m = 1$ .

### 1. Initial State:

- $b_1 = 0, b_2 = 0, b_3 = 1$

### 2. Step 1 (Mask = 1):

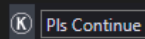
- $b_1 = b_3 = 1$
- $b_2 = b_1 = 0$
- $b_3 = b_2 = 0$
- New pattern: 100

### 3. Step 2 (Mask = 1):

- $b_1 = b_3 = 0$
- $b_2 = b_1 = 1$
- $b_3 = b_2 = 0$
- New pattern: 010

### 4. Step 3 (Mask = 1):

- $b_1 = b_3 = 0$
- $b_2 = b_1 = 0$
- $b_3 = b_2 = 1$
- New pattern: 001



This sequence continues cyclically.



## Advantages of Using Johnson Counters in DRAM

1. **In-Memory Processing:** Enables efficient counting and bit manipulation directly in DRAM, reducing data movement and latency.
  2. **High Parallelism:** Multiple counters can operate simultaneously across columns.
  3. **Low Hardware Complexity:** Johnson counters rely on simple shift and logical operations.
- 

## Challenges and Considerations

1. **Error Propagation:** Inverted feedback can amplify errors, requiring robust error correction mechanisms.
2. **Complex Implementation:** Efficiently managing mask values and counter initialization requires precise control logic.

## CIM-Logic Using Memristors

CIM-logic allows computation directly within the memory, reducing the need to transfer data between memory and processors. Memristors, which are resistive memory devices, are ideal for this due to their unique properties.

---

### Key Advantages of Memristors for CIM-Logic

1. **Inherent Statefulness:**
    - Memristors can store their resistance state even when power is turned off.
    - This makes them capable of storing data and performing computation in the same device.
  2. **Variable Resistance:**
    - The resistance of a memristor can be changed by applying voltage pulses.
    - This ability allows them to represent different logic states.
  3. **Threshold-Based Switching:**
    - Memristors have a specific voltage threshold that triggers their transition between high resistance (logic 0) and low resistance (logic 1).
    - This behavior is crucial for performing logic operations.
-



## How CIM-Logic is Implemented Using Memristors

1. **Data Representation:**
    - Memristors represent data as resistance:
      - High resistance = Logic 0
      - Low resistance = Logic 1
  2. **Logic Operation Steps:**
    - a. **Row Activation:**
      - Input data (operands) is stored in rows. To perform an operation (e.g., AND), rows containing the operands are activated by applying voltage.
    - b. **Parallel Evaluation:**
      - Logic operations are evaluated simultaneously for all columns.
      - For each column, the combined current or voltage is compared against a reference threshold to determine the result.
    - c. **Threshold-Based Comparison:**
      - A reference voltage/current is predefined for each logic operation (AND, OR, XOR). The output is decided by comparing the resulting bitline voltage/current with these reference values.
  5. **Result Storage:**
    - The outcome of the logic operation is stored in the resistance state of specific memristors.
- 

### Example: AND Operation Using Memristors

- **Operands:** Two rows store the input bits AAA and BBB.
  - **Process:**
    - Both rows are activated by applying voltage.
    - For each column:
      - If both memristors (representing AAA and BBB) are in a low-resistance state (logic 1), the resulting bitline current exceeds the AND gate threshold, giving a logic 1 output.
      - If either memristor has high resistance (logic 0), the current does not reach the threshold, resulting in logic 0.
- 

### Majority Function Using Memristors

- Memristors can perform a majority function by activating multiple rows simultaneously.
- The resulting bitline voltage reflects the majority value (most common logic state) among the activated memristors.
- This enables more complex logic operations without requiring additional sensing mechanisms.

---

## Comparison with DRAM-Based CIM-Logic

- Unlike DRAM-based CIM-logic, memristors:
    - Do not require dedicated sensing circuits.
    - Do not rely on destructive reads (data remains intact after computation).
  - Memristor-based CIM-logic is simpler and more efficient due to the intrinsic threshold-switching behavior.
- 

## Memristor Technology (RRAM) - Challenges and Limitations

1. **Write Operation Speed:**
    - Memristors, including RRAM, face slow write speeds due to the "write-verify" process. This involves writing data and then reading it back to verify its accuracy, which requires multiple repetitions to stabilize.
    - In comparison, DRAM and SRAM offer significantly faster write speeds.
  2. **Energy Consumption During Write:**
    - The write operation in NVM technologies like RRAM and PCM (Phase-Change Memory) is energy-intensive, similar to flash memory.
    - This can be a disadvantage for applications that require frequent write operations, as it consumes more energy compared to faster memory technologies.
  3. **Limited Write Endurance:**
    - Memristors, like flash memory, have limited write endurance. They can only withstand a finite number of write cycles before performance degrades.
    - The typical write endurance is about  $10^8$  cycles, which is considerably lower compared to DRAM and SRAM, which have practically unlimited write endurance.
  4. **Density Reduction in CAM Implementations:**
    - Using memristors (e.g., RRAM) to build Content Addressable Memories (CAMs) reduces memory density due to the additional circuitry required for comparison and match line operations.
    - This reduction in density limits the storage capacity of memristor-based CAMs.
- 

## Current Use and Potential:

- Despite the challenges, **RRAM** is commercially available and being utilized in some applications.
  - Notably, **crossbar arrays** are used in machine learning applications, where RRAM is employed for specific computational tasks.

- **Slow Write Speed** and **Limited Write Endurance** remain significant obstacles for the wider adoption of memristor technology.
- 

### Conclusion:

- Memristor-based technologies like RRAM offer promising applications, particularly in machine learning and certain computational tasks, but their widespread use is hindered by the slow write speed, energy consumption during write operations, and limited write endurance.
- Although these limitations present challenges, RRAM remains a commercially viable option, with ongoing developments to address these issues for future applications.

### Conclusion

Memristors, with their ability to combine data storage and computation, are a promising technology for efficient in-memory computing. Their threshold-switching, statefulness, and variable resistance allow logic operations to be performed directly in memory arrays, offering significant performance improvements over traditional approaches.

## Student Questions and Answers

### Q1: Why do we need RowClone?

- **A:** It enables efficient data copying within the array, such as swapping two variables in memory.
- 

### Q2: What if rows are in different arrays?

- **A:** Row cloning becomes complex; additional logic is required.
- 

### Q3: What is the maximum number of rows that can be activated simultaneously?

- **A:** Depends on technology. In DRAM, the maximum reported is **five**.

---

**Q4: Are there reliability issues with multi-row activation?**

- **A:** Yes, bit flips occur, and the probability increases with more activated rows.

---

**Q5: Does RowClone increase energy consumption?**

- **A:** No. Precharge operations are standard in memory reads and do not add significant overhead.

---

**Q6: How is addition in DRAM performed?**

- **A:** Using the majority operation and carry propagation in a **bit-serial, word-parallel** manner.

---

**Q7: How are operands stored for addition?**

- **A:** In **column-major order**, enabling shared bitline access.

---

**Q8: Can you give an ML example of bulk addition?**

- **A:** Multiply-accumulate (MAC) operations in ML often use bulk addition, especially for low-precision numbers.

---

**Q9: What are the overheads of replacing multiplication with addition?**

- **A:** Includes bit-slicing matrices and final accumulation steps. Fixed-weight matrices reduce this overhead.

# **Embedded LF2 Notes:**

## **Non-Volatile Memories**

### **Phase-Change Memory (PCM)**

- A mature technology previously used in Intel's Optane DC products but later discontinued.
- Positioned between DRAM and hard disks, faced challenges in its intended placement.
- Utilizes changes in material resistance for data representation.

### **MRAM (Magnetoresistive RAM)**

- Stores data based on the magnetic orientation of materials.
- Though not stored in resistance states, reading and writing leverage resistance states.
- A mature technology used in smartwatches like Huawei's and for matrix-matrix multiplication in crossbar accelerators.

### **Resistive RAM (RRAM)**

- Uses resistive materials like Hafnium Oxide.
- A mature technology frequently used in machine learning applications and sold by various companies.

### **Racetrack Memory**

- Introduced in 2008 and remains experimental.
- Developed by IBM and followed by Samsung, but no prototype exists yet.

---

## **Challenges with Non-Volatile Memories**

1. **Slow Write Operations:**
  - Writing involves a write-verify mechanism similar to flash drives.
  - Data is written, read back, compared, and rewritten if necessary. This process may repeat up to 51 times.
2. **Energy Consumption During Write:**
  - Write operations are energy-intensive in PCM and RRAM.
3. **Limited Write Endurance:**

- These memories have a finite number of write cycles ( $\sim 10^8$ ), much lower than DRAM or SRAM ( $\sim 10^{15}$ ).
- 

## Crossbar Architecture for In-Memory Computing

- External voltage or current is applied to program a vector or matrix into memory cells.
  - Pulse width or amplitude of the voltage signal represents input values.
  - Matrix-vector multiplication is performed by applying a second input vector across the memory array, generating the product in each cell.
- 

## Questions from the Lecture

1. **How is a mask used in CAMs?**
    - Masks enable or disable specific columns during the search.
  2. **What is the use of the address returned by the CAM?**
    - It retrieves data from a memory connected to the CAM, such as caches or lookup tables.
  3. **How do NVMs store multiple levels compared to SRAM or DRAM?**
    - NVMs use resistance values to represent multiple levels, while SRAM and DRAM are limited to binary states.
  4. **Are there other technologies used in databases besides CAMs?**
    - CPUs and traditional memories, where the CPU performs comparisons to search data.
  5. **Can CAMs be implemented in CNM architectures?**
    - Yes, though CNM systems are relatively new.
  6. **What is the error margin using decision trees for CAM computations?**
    - Accuracy depends on the training of the decision tree, independent of the CAM.
-

## What Are Content Addressable Memories (CAMs)?

Content:

CAMs are a special type of memory where **data is searched based on its content** rather than its address. Unlike traditional memories (RAM or ROM), where you provide an address to get data, in CAMs:

1. You provide the **data** to search for.
  2. CAMs return the **address** or row where that data is stored.
- 

## How CAMs Work (General Process)

CAMs perform searches quickly and in parallel. Here's the general process of how they work:

1. **Data Input and Search:**  
The data you want to search for is loaded into a special register called the **search data register**.
2. **Match Line Precharge:**  
Each row in the CAM has a **match line** (a wire running horizontally). All match lines are precharged to a **high state** (e.g., high voltage), assuming every row is a potential match.
3. **Search Line Activation:**  
The search data is applied to **search lines** (vertical wires). These lines carry the data to each cell in the memory array.
4. **Data Comparison:**  
Each memory cell compares its stored data with the search data on the search line:
  - If they **match**, the cell does nothing, and the match line stays high.
  - If they **mismatch**, the cell discharges the match line to ground, making it low.
5. **Result Detection:**  
Special circuits called **match line sense amplifiers** check the state of the match lines:
  - If a match line is still high → That row contains the matching data.
  - If a match line is low → The data in that row does not match.

The CAM outputs the **address** of the matching row(s) or a special code representing the match result.

---

## Types of CAMs

CAMs are categorized into two main types based on how they perform matches: **Binary CAMs (BCAMs)** and **Ternary CAMs (TCAMs)**.

---

# 1. Binary CAMs (BCAMs)

## Purpose:

BCAMs are used for **exact matches**. Every bit of the search data must match the stored data perfectly for a row to be flagged as a match.

## Key Features:

### 1. Data Representation:

- BCAMs use simple **binary values**: 0 and 1.
- Each memory cell stores **one bit** of data.

### 2. Exact Search:

- Every bit of the search data is compared with the stored data.
- If even one bit differs, the row is flagged as a mismatch.

### 3. Example:

- Suppose a BCAM stores 4-bit words.
- If the search data is 0101, it will match only rows containing the exact sequence 0101.
- Rows like 0111 or 1101 will not match.

## How BCAM Search Works:

### 1. Load Search Data:

- Input the search data (e.g., 0101) into the **search data register**.

### 2. Precharge Match Lines:

- All match lines are precharged to a **high state**, assuming every row is a match initially.

### 3. Activate Search Lines:

- Search lines apply the search data to the memory array for comparison.

### 4. Cell Comparison:

- Each cell compares its stored data with the corresponding search data bit:
  - If the bits match → The match line stays high.
  - If the bits mismatch → The match line discharges to ground (low state).

### 5. Sense Amplifiers Detect Results:

- Amplifiers check the match line state:
  - High match line → Row matches the search data.
  - Low match line → Row does not match.

## Applications of BCAMs:

- **Network Routing**: Exact matching of IP addresses.
- **Cache Memory**: Matching tags for memory blocks.

## Limitations:



- BCAMs are not suitable for searches involving patterns or approximate matches.
- 

## 2. Ternary CAMs (TCAMs)

### Purpose:

TCAMs allow for approximate searches using an additional state called "Don't Care" (X), making them more flexible than BCAMs. They also enable searches for exact binary numbers while still considering rows with "Don't Care" bits.

### Key Features:

- **Three States:**  
TCAM cells store 0, 1, or X ("Don't Care").
  - The "X" state allows specific bits to be ignored during searches.
- **Exact and Approximate Matching:**  
TCAMs can match exact binary numbers or approximate patterns. For example:
  - A search for 0101 matches rows like:
    - 0101 (exact match).
    - 0X01 (where "X" acts as "Don't Care").

### How TCAM Search Works:

1. **Load Search Data:**  
Input the search data (e.g., 0101) into the search data register.
2. **Precharge Match Lines:**  
As with BCAMs, all match lines are precharged to a high state.
3. **Activate Search Lines:**  
Search lines apply the search data to the memory array.
4. **Cell Comparison:**  
Each cell compares its stored data with the search data:
  - 0 matches 0, and 1 matches 1.
  - X always matches, regardless of the stored bit.
5. **Sense Amplifiers Detect Results:**  
Amplifiers check the match line state:

- High match line → Row matches the search criteria (considering "Don't Care" bits).
- Low match line → Row does not match.

For this example:

- Row storing 0101 → Match detected.
- Row storing 0X01 → Match detected (as "X" in the second position allows flexibility).

This flexibility makes TCAMs ideal for tasks like routing table lookups, where both exact and approximate matches are critical.

### Advantages of TCAMs:

- **Flexibility:** TCAMs support searches involving patterns or ranges.
- **Efficient for Complex Queries:** They handle multiple criteria in a single search.

### Applications of TCAMs:

- **Networking:**
  - Packet classification, filtering, and IP routing with wildcards (e.g., 192.168.X.X).
- **Databases:**
  - Searching with partial matches or ranges.
- **Pattern Recognition:**
  - Identifying patterns in image processing or bioinformatics.

### Limitations:

- **Complexity:** The additional "X" state increases circuit complexity.
- **Lower Density:** TCAMs store fewer bits in the same physical space compared to BCAMs.

## Comparison: BCAM vs. TCAM

Feature	BCAM	TCAM
States	0, 1	0, 1, x
Matching	Exact	Approximate (with x)
Applications	Cache, IP exact match	Routing, filtering, queries
Flexibility	Low	High
Complexity	Simple	More complex

## Conclusion:

- **BCAMs** are great for exact searches where all bits must match.
- **TCAMs** shine in scenarios requiring flexibility and pattern matching.

## Architecture of CAMs

The architecture of CAMs consists of several key components that work together to enable parallel search operations.

### 1. Memory Array

- The memory array is the central part of a CAM. It is organized as a grid of cells arranged in rows and columns.
- **Binary CAMs (BCAMs):** Each cell stores either a '0' or a '1'.
- **Ternary CAMs (TCAMs):** Each cell can store '0', '1', or 'X' (Don't Care state).
- This array serves both as storage for data and as the comparison unit for search operations.

### 2. Search Lines

- Search lines run vertically across the memory array and connect to each cell in a column.
- These lines carry the search data and distribute it to all cells in the array simultaneously, enabling parallel comparisons.

### **3. Match Lines**

- Each row of the memory array has a corresponding match line(horizontal).
- Match lines are initially set to a high voltage level, representing a potential match.
- The state of the match line changes based on the comparison result between the stored data and the applied search data.

### **4. Search Data Register**

- This register temporarily holds the data pattern that you want to search for within the CAM.
- The search data is then applied to the memory array through the search lines.

### **5. Match Line Sense Amplifiers**

- Sense amplifiers monitor the state of each match line to determine whether the stored data in a particular row matches the search data.
  - These amplifiers output the addresses of rows where matches are found or provide an encoded representation of the matching results.
- 

## **Functionality of CAMs**

The functionality of CAMs revolves around their ability to compare the search data against stored data across all rows simultaneously. This process is broken down into several steps:

### **1. Load Search Data**

- The data pattern to be searched is loaded into the search data register.

### **2. Precharge Match Lines**

- All match lines are precharged to a high voltage level, signifying an initial "match" state for every row.

### **3. Enable Search Lines**

- The search data is applied to all cells in the array simultaneously through the search lines.

### **4. Cell Comparison**

- Each cell in the memory array compares its stored data bit with the corresponding bit on the search line.
- The comparison logic depends on the type of CAM:
  - **Binary CAM (BCAM):**
    - **Match:** If the stored bit matches the search bit, the cell does not alter its match line.
    - **Mismatch:** If there's a mismatch, the cell discharges the match line to ground, signaling a "mismatch" for that row.
  - **Ternary CAM (TCAM):**
    - In addition to handling '0' and '1', TCAMs include the "Don't Care" ('X') state.
    - If the search bit is "X," the cell treats it as a match, regardless of its stored value.
    - This allows TCAMs to perform approximate searches and match patterns more flexibly.

## 5. Sense Amplifiers Output

- The sense amplifiers monitor the voltage of each match line to determine whether the stored data matches the search data.
  - Rows with match lines that remain in the high state are identified as matches.
  - The output is typically the addresses of matching rows or an encoded result.
- 

## Key Points about CAM Operation

- **Parallel Search:** CAMs compare search data with all rows at the same time, making searches extremely fast.
- **Match Line Logic:** A high match line state indicates a match, while a low state indicates a mismatch.
- **Complexity and Density:** CAM cells, especially in TCAMs (due to their three-state logic), are more complex and occupy more space compared to traditional memory cells.

## Slide (27)

### 1. SRAM Foundation:

- A CAM cell is built upon a basic SRAM (Static Random Access Memory) cell.
- The SRAM cell consists of two cross-coupled inverters, enabling the storage of either a logic 0 or logic 1.

### 2. Additional Components:

- The CAM integrates four transistors (labeled as M1, M2, M3, M4 in the circuit diagram) for comparison functionality.
  - Key components include:
    - **Search Lines (SL and SL\_bar):** Used to input search data (original and complementary values).
    - **Match Line (ML):** Used to indicate whether the stored data matches the search data.
- 

## Transistor Roles in Comparison

### ● Comparison Transistors (M1 to M4):

- These transistors control the comparison process by connecting the SRAM cell's storage nodes to the search and match lines.
- Their roles include:
  1. **M3 and M4:** Connect the storage nodes of the SRAM cell to the search lines (SL and SL\_bar).
  2. **M1 and M2:** Connect the match line (ML) to the search lines.

### ● Pull-Down Path:

- If the stored data and search data do not match, the transistor arrangement creates a path to discharge the match line (ML) to ground, indicating a mismatch.
- 

## Comparison Process

### 1. Precharge:

- Initially, the match line (ML) is precharged to a high voltage level. This high state assumes a match until proven otherwise.

### 2. Search Data Application:

- The search data is applied to the search lines (SL and SL\_bar):
  - SL carries the data (e.g., logic 1).
  - SL\_bar carries the complementary value (e.g., logic 0).

### 3. Transistor Switching:

- Based on the stored data and the search data, the transistors (M1–M4) switch ON (conducting) or OFF (non-conducting):
  - If the stored data matches the search data, the match line remains high.

- If the stored data differs from the search data, a discharge path is created, and the match line drops to a low voltage.
4. **Match Determination:**
- **Match (ML High):** When stored data equals search data, no discharge path exists for ML, and it stays high.
  - **Mismatch (ML Low):** When stored data differs, ML discharges to ground, indicating a mismatch.
- 

### Simplified Example

- **Stored Data (Logic 1):**
    - If the search data is logic 1:
      - Transistors connected to the matching storage node and search line remain OFF, blocking the discharge path.
      - ML remains high, indicating a match.
    - If the search data is logic 0:
      - A discharge path is created through one of the transistors.
      - ML discharges to ground, indicating a mismatch.
- 

### Ternary CAMs (TCAMs)

- **Three States:**
    - CAM cells are extended to ternary CAMs to handle three states (logic 0, logic 1, and don't care X).
    - This requires storing two bits per cell to represent these states.
- 

### Key Insights from the Circuit

- **Matching Process:** The four transistors (M1 to M4) perform the matching operation and implement the pull-down path.
- **Efficiency:** The CAM architecture allows fast, parallel searches by comparing the search data across all stored rows simultaneously.
- **Scalability Challenges:** CAMs require careful design to handle larger datasets and more complex operations (e.g., multiple matches).

## What is an NVM-based CAM?

- **Traditional vs NVM-based CAMs:**

Traditional CAMs are typically based on volatile memories like SRAM. However, using Non-Volatile Memories (NVMs), such as FeFET (Ferroelectric Field-Effect Transistor), allows the CAM to retain its content even when the power is off. This makes NVM-based CAMs suitable for applications requiring persistent storage.

---

## How Does an NVM-based CAM Work?

1. **Data Storage:**

- The NVM cells in the CAM array store the data to be searched.
- Each cell can represent multiple states (e.g., 0, 1, or a "don't care" state for ternary CAMs).

2. **Search Operation:**

- The data to be searched is loaded into a search data register.
- All the match lines, which are horizontal lines across the memory array, are precharged to a known state (e.g., a high voltage).
- The search lines, which are vertical lines across the memory array, are enabled, supplying the search data to each cell in the corresponding column.
- Each cell compares its stored value with the value on the connected search line.

3. **Comparison & Match Line Behavior:**

- If the values match, the cell does not change the state of the connected match line, leaving it in the precharged high state.
- If the values mismatch, the cell provides a path for the match line to discharge to the ground, bringing it to a low state.

4. **Output Generation:**

- After the comparison, match line sense amplifiers check the state of each match line.
  - The match lines that remain in the high state correspond to the rows where the search data was found. This effectively gives the addresses of the matching data.
- 

## Key Features and Considerations:

- **Parallel Search:**

CAMs search all rows simultaneously in a single step, making them significantly faster for content-based searches than traditional memory architectures.



- **Density Trade-off:**  
NVM-based CAMs require additional circuitry for comparison and match line operations, which reduces the overall memory density. Each comparison typically requires four transistors, consuming valuable chip area.
  - **Applications:**  
NVM-based CAMs are well-suited for the following applications:
    - **Database Searches:** Quickly finding specific records based on content.
    - **Lookup Tables:** Efficiently retrieving values associated with specific keys.
    - **Recommendation Systems:** Matching user profiles with large item databases to suggest relevant items.
    - **Pattern Recognition and Machine Learning:** Implementing algorithms like K-Nearest Neighbors (KNN) and decision trees.
- 

### Example: FeFET-based CAM

- **FeFET (Ferroelectric Field-Effect Transistor):**  
FeFET is an NVM technology that can store multiple states by varying the polarization of the ferroelectric material. This makes it suitable for implementing ternary or even higher-order CAMs.
  - **Working Principle of FeFET-based CAM:**  
The working principle of an FeFET-based CAM is similar to the general CAM description, with the FeFET cells performing the comparison and controlling the discharge paths for the match lines.
- 

### Conclusion:

- While NVMs offer persistent storage for CAMs, the added complexity and density trade-off should be carefully considered. The choice of NVM technology and CAM architecture depends on specific application requirements such as search speed, storage capacity, and power consumption.

## Using Content Addressable Memories (CAMs) in Recommendation Systems

Content Addressable Memories (CAMs) are well-suited for recommendation systems due to their ability to perform fast parallel searches on large datasets. Below is a detailed explanation of their application:

---

## How Recommendation Systems Work

### 1. Item Database:

- Recommendation systems operate on large databases, such as movies on Netflix or products on Amazon.
- Each item in the database has multiple attributes (e.g., movie genre, actors, product category, customer reviews).

### 2. User Profile:

- User interactions (e.g., watching movies, browsing products) are tracked to create a user profile.
- This profile captures the user's preferences and interests based on their activity history.

### 3. Matching and Recommendation:

- The system compares the user profile attributes to the item database to find matches.
  - The aim is to efficiently identify items that align with the user's preferences, especially in large datasets.
- 

## Advantages of CAMs in Recommendation Systems

### 1. Parallel Search Capability:

- Traditional CPU-based searches are sequential and slow when working with large datasets.
- CAMs can search all items simultaneously in a single step, making them highly efficient.

### 2. Content-based Search:

- CAMs search for data based on its content rather than its address, aligning perfectly with the need to match user profile attributes to item attributes.
- 

## Detailed Example: Movie Recommendation Using CAMs

### 1. Store Movie Attributes in CAM:

- Each row in the CAM array represents a movie.

- Each column corresponds to a specific attribute (e.g., genre, actors, director).
  - The cells store attribute values for each movie.
  - 2. **Encode User Profile:**
    - The user profile is encoded into a search data pattern reflecting the user's preferences.
    - For example, if a user enjoys action movies with Tom Cruise, the search pattern represents these preferences in the relevant columns.
  - 3. **Parallel Search:**
    - The encoded user profile is presented to the CAM.
    - All rows in the CAM (movies) are simultaneously compared to the search pattern.
  - 4. **Match Identification:**
    - Rows that match the user profile remain in a high state on their match lines.
    - These match lines indicate the addresses of the matching movies.
  - 5. **Recommendation Generation:**
    - The details of the matching movies are retrieved and presented to the user as recommendations.
- 

## Challenges and Practical Considerations

1. **CAM Array Size:**
    - CAM arrays have limited storage capacity compared to massive databases.
    - Techniques like clustering and hierarchical searches are needed for very large datasets.
  2. **Data Duplication and Reloading:**
    - When the database exceeds the size of the CAM array, it must be loaded and searched sequentially.
    - This can introduce latency, requiring strategies to minimize data reloading for optimal performance.
  3. **Complexity and Density:**
    - CAMs implemented with NVMs (non-volatile memories) require additional circuitry, reducing memory density.
    - This trade-off between functionality and storage capacity must be carefully managed.
- 

## Key Challenges

1. **Mismatch Between CAM Capacity and Dataset Size:**
  - CAM arrays typically have limited storage capacity (e.g., kilobytes or megabytes), while databases often span gigabytes or terabytes.
  - This disparity necessitates dividing datasets into smaller chunks that fit into CAM arrays.

## 2. Data Reloading:

- Large datasets must be incrementally loaded into CAM arrays, increasing latency and reducing performance.
  - Efficient scheduling can mitigate these issues.
- 

## Application Mapping and Scheduling on CAMs

When dealing with large datasets that exceed the capacity of CAM arrays, application mapping and scheduling strategies are crucial to optimizing performance. The diagram illustrates the flow of tasks, highlighting the role of CAM arrays in managing query intervals and data blocks.

---

### Application Mapping

#### 1. Dividing Data into Chunks:

- Large datasets are split into smaller, manageable blocks that fit within the CAM's storage capacity.
- Example: A 1 GB database is divided into 1 MB chunks if the CAM can handle 1 MB at a time.

#### 2. Mapping Chunks to CAM Rows:

- Each chunk is mapped to specific rows in the CAM array.
  - For instance, with a row storing 128 bits, an 8-row CAM array could accommodate 8 chunks simultaneously.
- 

### Scheduling Search Operations

#### 1. Sequential Search with Single Reload:

- The first data chunk is loaded into the CAM, and all queries are processed sequentially against it.
- After completing the queries, the next chunk is loaded, and the process repeats.
  - **Advantage:** Minimizes the number of data reloads.
  - **Drawback:** High latency, as queries must wait for others to complete.

#### 2. Interleaved Search with Potential Parallelism:

- While one query searches a specific chunk, another query can search a different chunk simultaneously.
- This approach improves throughput, provided the CAM architecture supports parallel searches.
  - **Advantage:** Reduces latency for individual queries.
  - **Drawback:** Managing and merging intermediate results becomes more complex.

---

## Goals of Mapping and Scheduling

### 1. Minimizing Reloads:

- Reduce the frequency of reloading data blocks into CAM arrays to minimize latency and improve system efficiency.

### 2. Optimizing Latency and Throughput:

- Strike a balance between reducing data reloads and processing queries quickly.
  - Scheduling strategies must align with the specific requirements of the application.
- 

## Practical Considerations

### 1. Data Duplication for Chunk Boundaries:

- Records spanning multiple chunks need duplication at chunk boundaries to prevent data loss during searches.
- Example: A record overlapping chunks 63 and 67 should be fully included in both chunks.

### 2. Handling Intermediate Results:

- Efficiently merging and interpreting partial results from interleaved searches is a significant challenge.
  - This adds to the overhead and requires robust system design.
- 

## Conclusion

Effective application mapping and scheduling are vital for leveraging CAMs with large datasets. The objectives are to:

- Minimize latency by optimizing data reloads.
- Enhance throughput with parallel query processing, where possible.
- Balance system efficiency while managing intermediate results.

The efficient handling of intermediate results, as highlighted in the diagram, is critical for achieving these goals.

## Associative Processors: Keys, Masks, and Tags in CAMs

**Content Addressable Memories (CAMs)** are the core components of associative processors, allowing for parallel search and retrieval operations. To understand how associative processors function, it's essential to grasp the role of keys, masks, and tags:

1. **Keys**

Keys represent the data used to query the CAM. For instance, in a database, the key corresponds to the record you are trying to locate. The CAM searches for matches between the key and its stored data.

2. **Masks**

Masks enable selective column activation during searches, allowing focused queries. For example, if you want to search only the first three columns of data, a mask can be applied to disable all other columns. A mask with **11100** would activate the first three bits and ignore the rest.

3. **Tags**

Tags are identifiers for data in cache lines that use CAMs for quick lookups. During a processor's data request, the CAM searches for matching tags. If a match is found, the associated data is fetched, ensuring high-speed access compared to accessing main memory.

---

## Full Adder Implementation Using CAMs

A full adder calculates the sum of three bits: two inputs (A, B) and a carry-in (Cin). Outputs are the sum (S) and carry-out (Cout). Using CAMs, pre-stored truth tables eliminate real-time computation.

### Step 1: Building the CAM

The CAM stores the full adder's truth table with five columns:

<i>A</i>	<i>B</i>	Cin	<i>S</i>	Cout
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

Each row in the CAM represents a single input-output mapping for the full adder operation.

### Step 2: Presenting Inputs

Suppose we want to add the binary numbers **101** and **011**. The addition is performed bit-by-bit, starting with the least significant bit (LSB):

Bits from Numbers	<i>A</i>	<i>B</i>	Cin (Carry-in)
1st Bit	1	1	0
2nd Bit	0	1	1
3rd Bit	1	0	1

### Step 3: Search and Retrieve

For each bit, search the CAM to find matching rows:

1. **1st Bit:**  $A = 1, B = 1, C_{in} = 0 \rightarrow S = 0, C_{out} = 1$ .
2. **2nd Bit:**  $A = 0, B = 1, C_{in} = 1 \rightarrow S = 0, C_{out} = 1$ .
3. **3rd Bit:**  $A = 1, B = 0, C_{in} = 1 \rightarrow S = 0, C_{out} = 1$ .

### Step 4: Result Compilation

Combine the results:

- Sum: 000
- Final Carry-out: 1

The result for **101 + 011** is **1000**.

---

## Advantages

1. Efficiency: Parallel search retrieves results instantly.
2. Simplicity: Precomputed logic simplifies runtime operations.
3. Scalability: Larger truth tables can handle complex tasks.

## Applications

- Cryptographic operations
- Signal processing
- High-speed data lookup in networking hardware

## **Embedded LF3 Notes:**

## **Embedded LF4 Notes:**

### Introduction to Quantum Computing

- Quantum computing uses principles of quantum mechanics, enabling parallel computation through properties like superposition.
- **Superposition:** A qubit can exist in multiple states simultaneously.
- Quantum mechanics governs microscopic systems (e.g., protons, electrons) and is represented using complex numbers.



## Swapping Operation in Quantum Computing (Slide 10)

A black box with two holes swaps the state of the input ball. For example, when a white ball is thrown into the top hole, a black ball emerges from the bottom, and vice versa. This operation is analogous to a **NOT gate** in classical computing, which inverts the input state.

When two such black boxes are used in sequence, the combinations depend on the input states. For instance, if two white balls are input, the output remains two white balls. However, if a white and a black ball are input, they emerge swapped, demonstrating how states are exchanged within the boxes.

It highlighted that while in classical computing, one might not analyze all possible input combinations for a known operation, in quantum mechanics, all combinations are considered due to the nature of physical experiments.

## (Slide 11)

If the control input is "white," the target input remains unchanged. This is because "white" represents the state where the gate is disabled.

If the control input is "black," the target input is inverted. This means "black" enables the gate, leading to the NOT operation being applied.

The key takeaway is that the CNOT gate provides a mechanism for conditional manipulation of quantum states, where the state of one qubit (the control) influences the operation performed on another qubit (the target). This controlled behavior is fundamental to the power of quantum computing.

## (Slide 12)

If the control input is "white," symbolizing the gate's disabled state, the states of the two target inputs remain unchanged.

If the control input is "black," signifying the gate's enabled state, the states of the two target inputs are swapped.

## (Slide 13)

Therefore, stacking two NOT gates results in the output being the same as the input. The first NOT gate inverts the state, and the second NOT gate inverts it back to its original state. The source states that this behavior is similar to classical computing where stacking two NOT gates results in the input being the same as the output

## (Slide 17)

1. **First Gate (CNOT):**
  - The control input is black, so the target input is inverted.
  - Output: White and black.
2. **Second Gate (CSWAP):**
  - The control input is black, so both target inputs are inverted.
  - Output: Black and white.
3. **Third Gate (CNOT):**
  - The control input is white, so the target input remains unchanged.
  - Output: Black and white.

**Final Output:** Black, black, white.

The Pete gate introduces probabilistic behavior, adding an element of randomness to quantum operations. It takes a single input, which can be either a "white" ball or a "black" ball, and produces a random output.

## (Slide 24)

For a black ball, the output is also a superposition, but with the black ball marked by a negative sign, indicating a phase difference. This phase difference is crucial in quantum mechanics, as the phase of each state affects how states interact and evolve. While the analogy simplifies the concept, the negative sign highlights the importance of phase relationships in superposition, which is key to phenomena like interference. To fully understand these concepts, further exploration of the mathematical framework of quantum mechanics is necessary, particularly the role of complex numbers in representing quantum states.

## (Slide 27)

## Quantum Superposition and Interference in Pete Gates

1. **First Pete Gate:**
  - A white ball input creates a **superposition** of a white and black ball.
2. **Second Pete Gate:**
  - The superposition from the first gate is processed, creating a more complex state with four possibilities.
3. **Interference:**
  - A **negative phase** from the black ball components causes them to **cancel**, while white ball components **constructively interfere**.
4. **Deterministic Output:**
  - The final output is a single white ball. Similarly, a black ball input results in a deterministic black ball output.

This highlights how **quantum mechanics** uses superposition and interference for deterministic outcomes

## Entanglement:

Entanglement is a uniquely quantum phenomenon that sets quantum computing apart from classical mechanics. It occurs when two quantum particles are created in opposite states. If the state of one particle changes, the state of the other particle is instantly affected, even if they are separated by vast distances. This phenomenon challenges Einstein's theory of relativity, which asserts that no two objects can communicate faster than the speed of light.

An analogy for entanglement is a pair of shoes ordered online. If one shoe is delivered and found to be right-handed, it immediately reveals that the other shoe, wherever it may be, is left-handed. This analogy illustrates that entangled particles have predefined opposite states, and knowing one reveals the state of the other.

However, unlike physical objects such as shoes, measuring a quantum particle changes its state. When one entangled particle in a superposition is measured, the superposition state of the other particle also changes instantly. This instantaneous change, even over great distances, is often referred to as "spooky action at a distance."

For example, consider two electrons with opposite spins: one spin up and the other spin down. Until measured, the electrons remain in a superposition, meaning there is no definite answer to which electron has which spin. Once one electron is measured, the spin of the other instantly becomes the opposite.

Entanglement is not just a random pairing of particles but involves particles created specifically with opposite states, making it a fundamental aspect of quantum mechanics.

**Embedded LF5 Notes:**