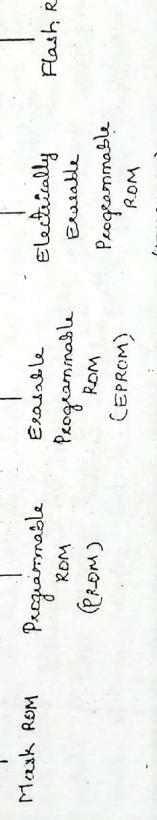
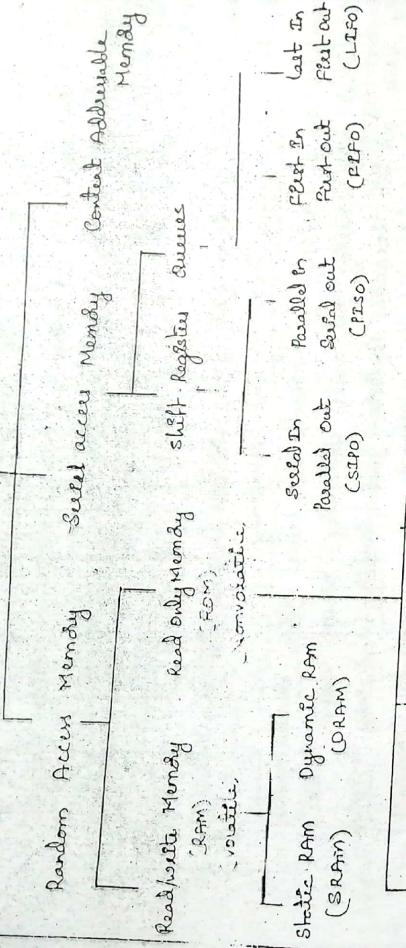


# UNIT-IV

## Memory Subsystem

### Introduction:-

Memory arrays often account for the majority of transistors in a micro system-on-chip.



Random access memory is accessed sequentially with an address and has a latency independent of the address.  
Serial access memories are accessed sequentially so no address is necessary. Content addressable memories determine which address(es) contain data that match a specified key.

Random access memory is commonly classified as read-only memory (ROM) and read/write memory (RAM). A more useful classification is volatile and nonvolatile memory. Volatile memory retains its data as long as power is applied, while nonvolatile memory will hold data indefinitely. RAM is synonymous with volatile memory while ROM is synonymous with non-volatile memory.

The memory cells used in volatile memories can further be divided into static structures and dynamic structures.

Static cells use some form of feedback to maintain their state, while dynamic cells use charge stored on a floating capacitor through an access transistor. Charge will leak away through the access transistor even while the transistor is off, so dynamic cells must be periodically read and rewrites to refresh their state. Static RAMs are faster and less troublesome, but require more area per bit than dynamic counterparts (DRAMs).

Some nonvolatile memories are indeed read-only. The content of a mask ROM are hardened during fabrication and cannot be changed.

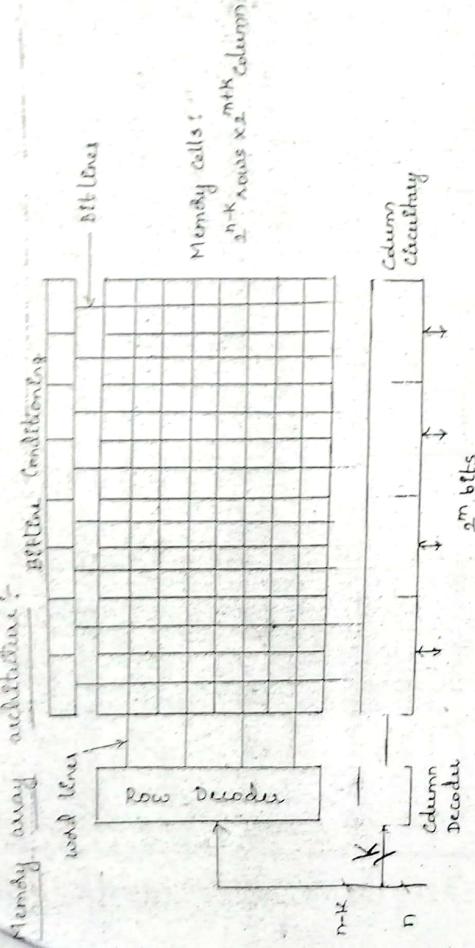
A programmable ROM (PROM) can be programmed once after fabrication by blowing on-chip fuses with a special high programming voltage.

An erasable programmable ROM (EPROM) is programmed by storing charge on a floating gate. It can be erased by exposure to ultraviolet (UV) light for several minutes to knock the charge off the gate. Then the EPROM can be programmed.

Finally, erasable programmable ROMs (EEPROMs) are similar, but can be erased in microseconds without ultraviolet necessary.

Flash memories are a variant of EEPROM that erase entire blocks rather than individual bits. Sharing the erase circuitry across larger blocks reduces the area per bit. Because of their good density and easy in-system reprogrammability, flash memories have replaced other nonvolatile memories in most modern CMOS systems.

Memory cells can have one or more ports for access. On a read/write memory, each port can be read-only, write-only, or capable of both read and write.



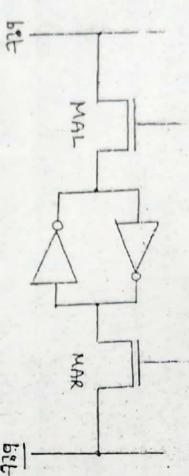
Above figure shows typical small memory array architecture. central to the design is a memory array consisting of  $2^n$  words of  $2^m$  bits each. In the simplest design, the array is organized with one row per word and one column per bit in each word. Often there are far more words in the memory than bits in each word, which would lead to a very tall, skinny memory that is hard to fit in the chip footprint and slow because of the long vertical wires. Therefore, the array is often folded into fewer rows of more columns. After folding, each row of the memory contains  $2^k$  words, so the array is physically organized as  $2^{n-k}$  rows of  $2^m$  columns of bits. The row decoder activates one of the rows by asserting one of the word lines. During a read operation, the cells on this word line drive the bit lines, which may have been conditioned to a known value in advance of the memory access. The column decoder controls or multiplexes the column circuitry to select  $2^m$  bits from the row as the data to access.

Larger memories are generally built from multiple smaller subarrays so that the word lines and bit lines remain reasonably short, fast, and low in power dissipation.

### Static RAM:

SRAM cells use a simple bistable circuit to hold a data bit. A static RAM cell can hold the stored data bit so long as the power is applied to the circuit. SRAMs have three operational modes. When the cell is in a hold state, the value of the bit is stored in the cell for future usage.

During a write operation, a logic OR is fed to the cell for storage. During a read operation, the value of the stored bit is transmitted to the outside world.



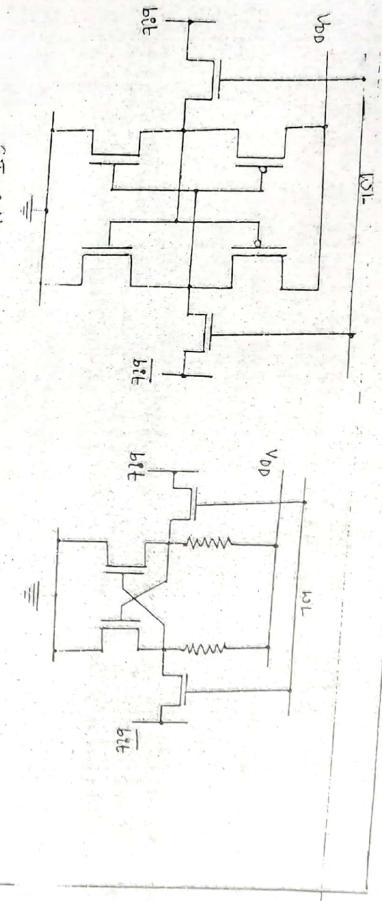
General SRAM cell

Above figure shows basic SRAM cell. A pair of cross-coupled inverters provides the storage, while two access transistors MA1 and MA2 provide read and write operations. The access transistors are controlled by the word line signal WL that defines the operational modes.

When  $WL=0$ , both access FETs are off and the cell is isolated. This defines the hold condition. To perform a read or write operation, the word line is brought up to a value of  $WL=1$ . This turns on the access transistors connecting the dual rail data lines bit and  $\overline{bit}$  to the outside circuitry. These are often called the bit and bit-bar lines, respectively.

A write operation is performed by placing voltage on the bit and  $\overline{bit}$  lines, which then act as inputs. Dual rail logic helps enhance the writing speed.

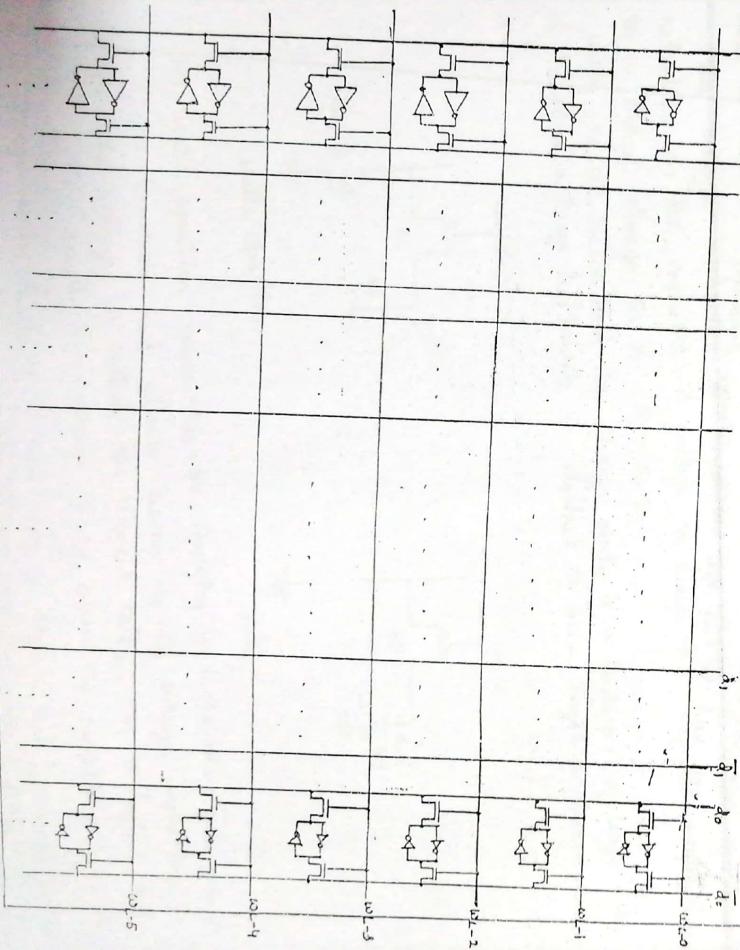
For a read operation, the bit and bit-bar lines act as outputs and are fed into a sense amplifier that determines the stored value or stored state. The distinction between read and write operations is obtained by circuitry outside the cell array.



Two types of CMOS cells are dominant in practice. 6T cell uses standard CMOS inversions. 4T cell uses resistors as load devices in an NMOS circuit.

### SRAM Arrays:

Static SRAM arrays are created by replicating the basic storage cell and adding the necessary peripheral circuitry.



word lines are oriented to run in a vertical direction.

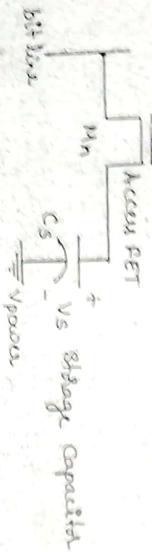
bit-bar lines are patterned vertically. The outputs of a row decoder provide the word line signals to the storage cells. The address word specifies a particular row, which is then driven high. The access transistors of the selected row cells are turned on, permitting the read/write operations to take place.

The row decoder outputs are fed into row driver circuits that are used to drive the word lines of the array. Drivers are needed because of the large capacitive load presented by the long interconnects and the access transistors connected to each word line.

The input/output bit and bit-bar data lines of the array are connected to the memory matrix. The data flow is thus visualized to be vertical for both read and write operations. Once a word line is driven high by the row decoder, every cell in the row is accessible. To choose a particular k-bit word in the row, we must add the group of column decoder circuits that select a particular set of k columns in the matrix.

A read operation requires an output from the cells, so the column decoder circuit acts as multiplexer. For a write operation, the sense mode must be used to store a data word into the proper columns. Column drivers are used to feed the memory column decoders.

### word line (WL)



LT. DRAm cell

**Dynamic RAM (DRAM)** use substantially smaller than SRAM cells. DRAM could lead to higher density storage arrays. The reduced cost per bit makes them attractive for applications requiring large read/write memory sizes. DRAMs are slower than SRAMs, and require more peripheral circuitry at the circuit level, they are simple in structure.

Above figure shows 1T DRAM cell. It consists of a single access M<sub>N</sub> and a storage capacitor C<sub>S</sub>. The cell is controlled by the word line signal WL and a single bit line provides the I/O path to the cell. The bottom of the capacitor is connected to one of the power supply rails, and is denoted as V<sub>DD</sub> or V<sub>SS</sub>, either V<sub>DD</sub> or V<sub>SS</sub> may be used.

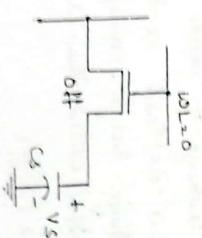
The storage mechanism is based on the concept of temporary charge retention on the capacitor. A voltage V<sub>S</sub> across the capacitor corresponds to a stored charge Q<sub>S</sub> if Q<sub>S</sub> = C<sub>S</sub>V<sub>S</sub>. With V<sub>S</sub>=0V, Q<sub>S</sub>=0 and the charge state is a logic 0. A large value of V<sub>S</sub> gives a large Q<sub>S</sub>, which is defined to be a logic 1 charge state.

$$WL=1$$



Write Operation

$$WL=0$$



Hold

In write operation, V<sub>POWER</sub> = V<sub>DD</sub> = 0V. Applying V<sub>DD</sub> to the input gate turns on the access transistor and allows access to the storage capacitor. The output data voltage V<sub>O</sub> controls the current to/from C<sub>S</sub>. A logic 0 data voltage V<sub>O</sub> = 0V results in a voltage V<sub>S</sub> = 0V across the capacitor, corresponding to a charge state of Q<sub>S</sub> = 0. If we apply a logic 1 data voltage V<sub>O</sub> = V<sub>DD</sub>

equal to the power supply, the voltage on the gate reduces the transmitted signal by an nFET threshold voltage. The target voltage that can be passed to the capacitor is

$$V_S = V_{MAX} = V_{DD} - V_{TH}$$

which gives a maximum charge of  $Q_{MAX} = C_S(V_{DD} - V_{TH})$ . The hold state is achieved by turning off the access transistor with the word line signal of  $WL=0$ .

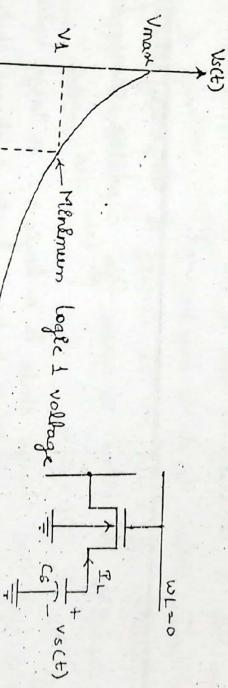
The dynamic aspect of the cell arises during a data hold time. A MOSFET that is biased off with  $V_S < V_T$  still admits small leakage currents.

A logic 1 voltage  $V_S > V_{MAX}$  on the storage capacitor provides the electromotive force for the leakage current to flow away from  $C_S$ . This can be described by

$$I_L = - \left( \frac{dQ_S}{dt} \right)$$

which shows that the current removes charge from the capacitor.

$$\therefore T_L = -C_S \left( \frac{dV_S}{dt} \right) \text{ so that } V_S \text{ also drops.}$$



charge leakage in DRAM cell.

The hold time  $t_h$  is defined as the longest period of time that the cell can maintain a voltage large enough to be interpreted as a logic 1, the hold time is also called the retention time.

$$I_L = -C_S \left( \frac{\Delta V_S}{\Delta t} \right)$$

where  $\Delta V_S$  and  $\Delta t$  represent changes in the variables.

$$t_h = | \Delta t | = \left( \frac{C_S}{I_L} \right) \Delta V_S$$

- Above expression shows that the hold time may be increased by using a large capacitance and minimizing the leakage current.

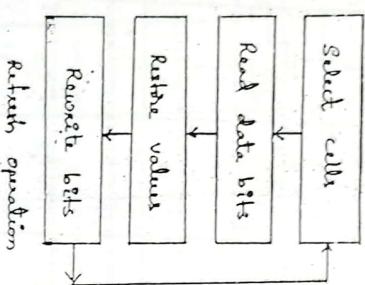
### Refresh Operation

Memory cells must be able to hold data so long as the power is applied to overcome the charge leakage problem. DRAM arrays employ a refresh operation where the data is periodically read from every cell, amplified, and then rewritten.

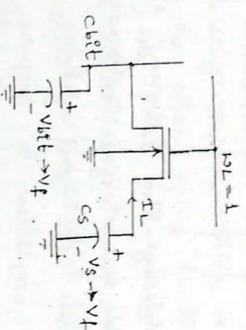
The cycle must be performed on every cell in the array with a minimum refresh frequency of about

$$\text{refresh} = \frac{1}{2t_h}$$

refresh circuitry is included to insure success logic that guarantees the cell array. The refresh cycle is designed to operate in the background and is therefore transparent to the user.



Read operation in a DRAM cell:



The voltage  $V_S$  on the capacitor at the read time provides the voltage to move charge from  $C_S$  to the bit line capacitance  $C_{BL}$ , which sets up a charge sharing situation.  $C_{BL}$  includes the bit line capacitance and other parasitic contributions such as the input capacitance of the sense amplifier. The initial charge on the capacitor is

$$Q_S = C_S V_S$$

where  $V_S > 0$  for a logic 0, and  $V_S < 0$  for a logic 1.

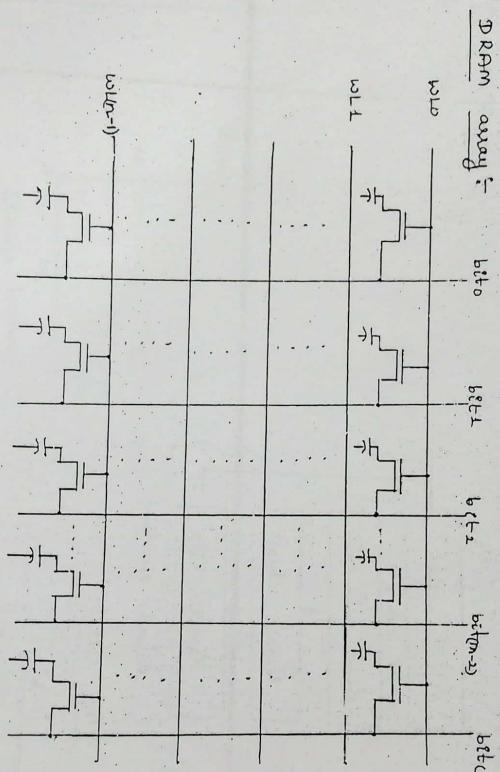
current flows from  $C_s$  to  $C_{bit}$  continues until the voltages are equal to the final voltage  $V_f = V_{bit} + V_s$ . Thus change in node establishing according to  $C_s = C_S V_f + C_{bit} V_f$

The initial and final values of  $C_S$  must be equal by charge conservation, so

$$V_f = \frac{C_S}{C_S + C_{bit}}$$

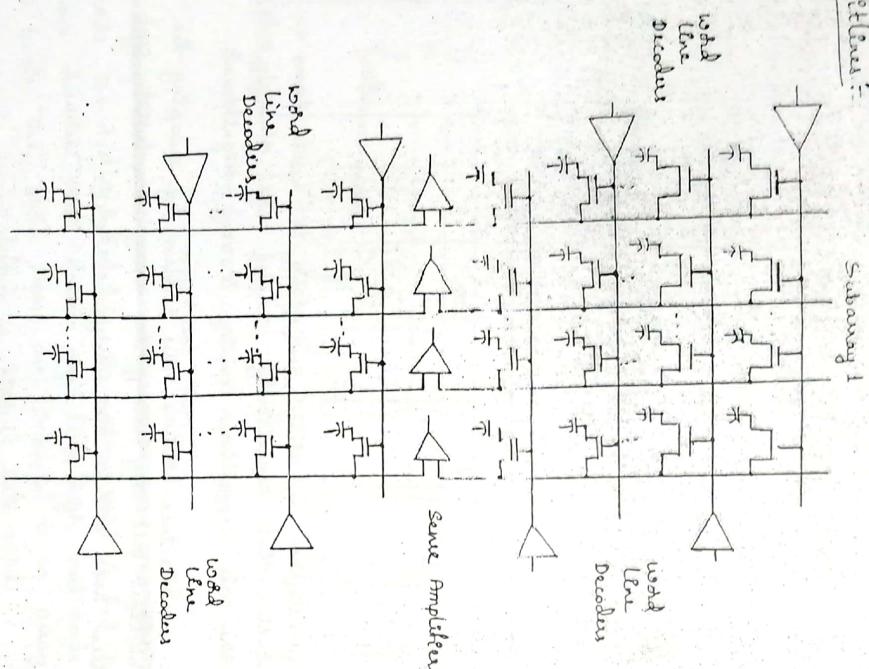
$$\therefore V_f = \left[ \frac{C_S}{C_S + C_{bit}} \right] V_S$$

In steady state  $V_f = V_S$  i.e., steady state loss is reduced to a few tenths of a volt, so that the number of bits per array becomes a critical factor.



large DRAMs are divided into multiple subarrays. The subarray size represents a tradeoff between density and performance. Larger subarrays are slow and have small bitline swings because of high wordline and bitline capacitance. The array uses a sense amplifier to compare the bitline voltage to that of an idle bit line (procharged to  $V_{DD}$ ). The sense amplifier must also be very compact to fit the tight pitch of the bitline. The low-swing bitlines are very sensitive to noise. These bitline architectures, open, folded and twisted, offer different compromises between noise and area.

### Open-bitline:

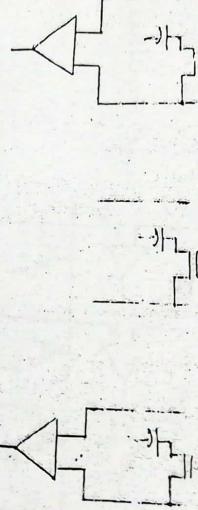
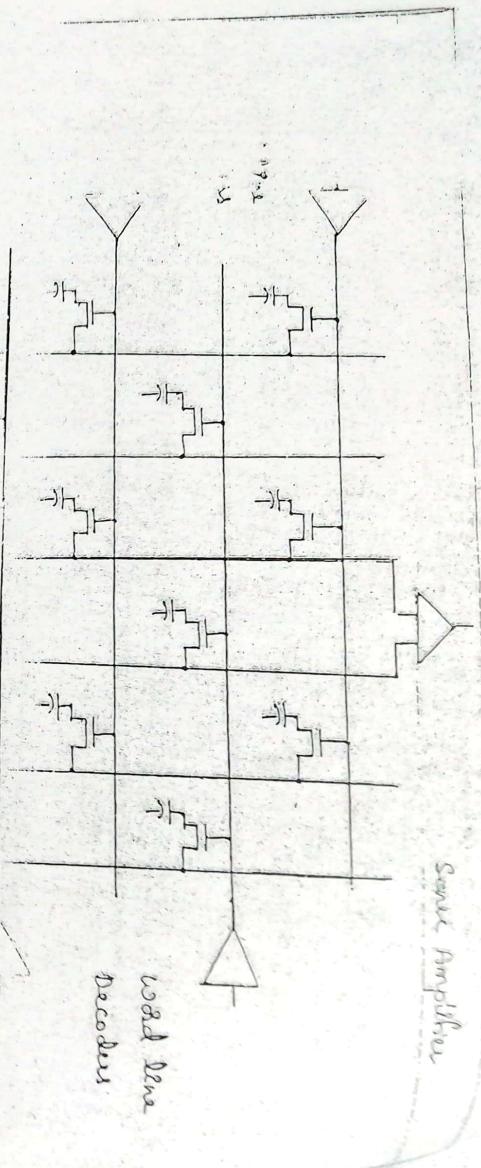


Subarray<sup>2</sup>

In this architecture, the sense amplifier receives one bitline from each of two subarrays. The wordline is only asserted in one array, leaving the bitlines in the other array floating at the reference voltage. The arrays are very dense. And noise that affects one array more than the other will appear as differential noise at the sense amplifiers. Thus, open bit lines have unacceptably low signal-to-noise ratios for high-density DRAM.

### Fused bitline:

In this architecture, each bitline connects to only half as many cells. Adjacent bitlines are organized in pairs as inputs to the sense amplifiers. When a wordline is asserted, one bitline will switch while its neighbor serves as the quiet reference. Many noise sources will couple equally onto the two adjacent bitlines so they tend to appear as common mode noise that is rejected by the sense amplifier. The noise advantage comes at the expense of greater layout area.



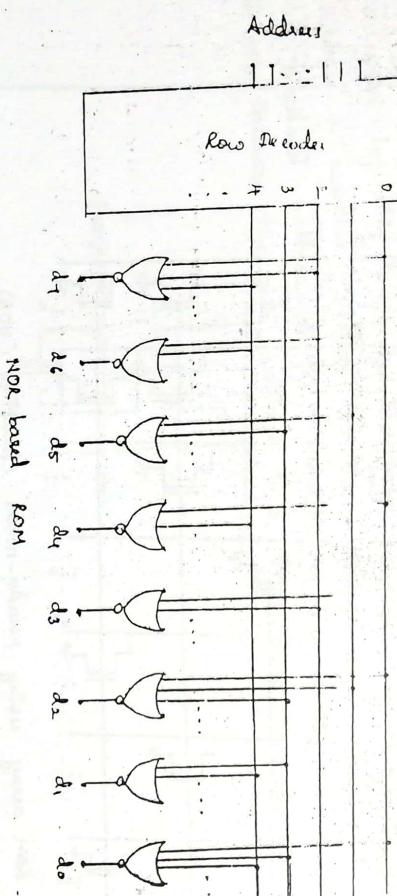
Sense amplifier

Traditional bitline architecture is still susceptible to noise from a switching bitline that capacitively couples more strongly onto bitlines in the pair. Capacitive coupling is very significant procedure.

Interdigitated bitline architecture solves this problem by sweeping the total bitline part way along the array. This reduces the amount of extra area within the array.

### ROM arrays:

Read only memories are used for permanent bit storage. The structure of a ROM array is similar to that used for RAMs, but the individual bit cells are much simpler. The data stored in a basic read only memory is created by the selective placement of PEs. Since this is accomplished in the physical design, the data cannot be altered once the chip is fabricated.

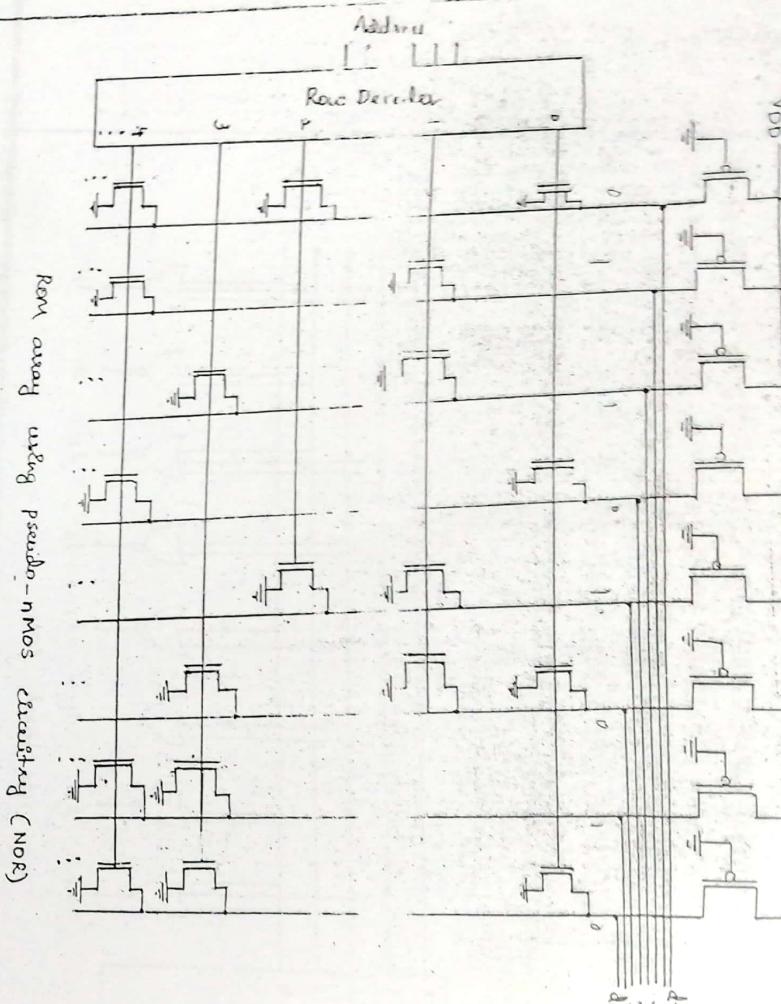


Above figure shows a ROM array that uses NOR gates to store 8-bit data words. D = address word. An address word is fed into an active high row decoder that drives one line high and keeps the other at logic 0 levels. The word lines are connected to an array of NOR gates such that each row defines a distinct data word.

Ex :- Row  
0 01101010  
1 10010011  
2 01110111  
3 11011000  
4 00101100

### ROM array using pseudo-nMOS:

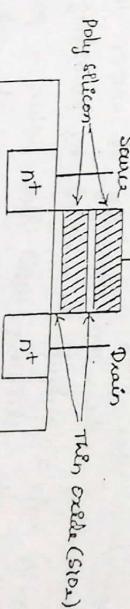
Pseudo-nMOS uses only one pull-up PFET for each NOR gate, the task of programming centers on placement & programming controls on placement of PFETs that act as pull down devices. A logic 0 output is obtained by providing a PFET with its gate connected to the driving word line. When a pull-down transistor turns on, it provides a good connection to ground and pulls the output low.



### User programmable ROM :-

Electrically programmable ROMS (EPROM) allow the user to store data as required by the application. Special voltage settings are used to write to the cells. Read operations are performed with normal voltage levels, so that the data remains unchanged. Many ROM devices provide for erasing and reprogramming the contents of the array. Optical erasure using UV light was used in early EPROM designs, but these have been replaced by electrically alterable devices.

A reprogrammable ROM array is built using special PETS that use a pair of stacked poly gates and has the circuit symbol as shown in figure.

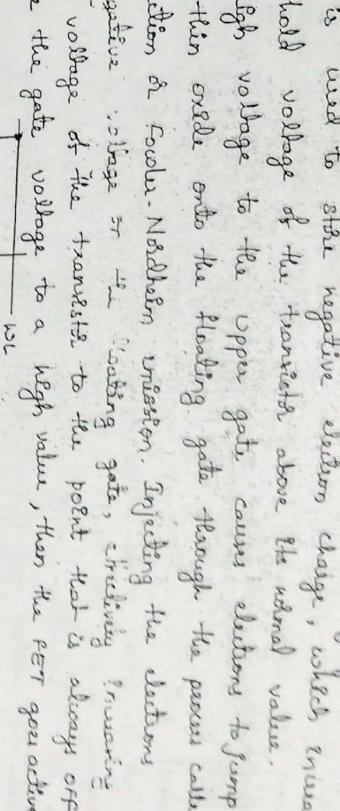


the topmost

gate constitutes the usual gate terminal of the transistor.

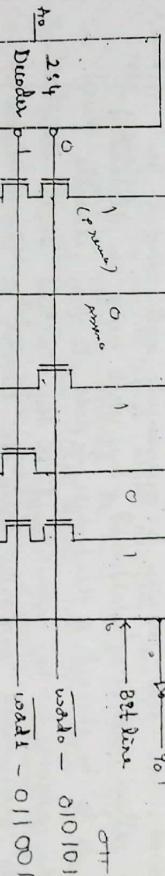
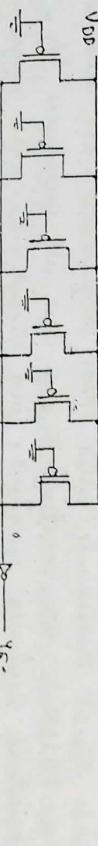
Another polygate layer is sandwiched in between the top poly and the silicon substrate. It is not electrically connected to any part of the transistor, and is therefore called an electrically floating gate. The floating gate is used to store negative electron charge, which increases the threshold voltage of the transistor above its normal value.

Applying a high voltage to the upper gate causes electrons to jump through the thin oxide onto the floating gate through the process called avalanche injection or Fowler-Nordheim injection. Injecting the electrons in excess to negative charge in the floating gate, effectively increasing the threshold voltage of the transistor to the point that is always off. If we increase the gate voltage to a high value, then the PNP goes active.



EEPROM cell with write line

Pseudo-nMOS NAND ROM



NAND ROM uses active low methods. Transistors are placed in series associated with the selected rows are on. If the no transistor is present, the bit line will remain high.

In NOR ROM the size of the cell is limited by the ground line.

In NAND ROM delay grows gradually with the number of selected transistors decreasing the lifetime.

#### Types of ROMS:

The required paths in a ROM may be programmed in two different ways

— Mask is called mask programming and is done by the semiconductor company during the last fabrication process of the chip. It needs a ROM service that the customer tells out the corresponding mask for the paths to produce the 1's and 0's according to the customers truth table. This procedure is costly because the vendor charges this customer a special fee for custom masking the particular ROM. For some ROM configuration is to be altered.

called programmable read-only memory or PROM. When ordered units contain in the ROM are blown by application of a high-voltage pulse to the device through a special pin. A blown fuse defines a binary 0 state and an intact fuse gives a binary 1 state. This allows the user to program the PROM in their laboratory to achieve the desired relationship between input address and stored words. Special instruments called PROM programmers called PROM programmers are available commercially to facilitate this procedure. In any case, all procedures for programming ROMs are hardware procedures even though the word program-

— The hardware procedure for programming ROMs or PROMs is irreversible and, once programmed, the fixed pattern is permanent and cannot be altered. Once a bit pattern has been established, the unit must be discarded if the bit pattern is to be changed. A third type of ROM is the erasable PROM or EPROM. The EPROM can be restructured to the initial state even though it has been programmed previously. When the EPROM is placed under a special UV light for a given period of time, the short wave radiation decomposes the internal floating gates that serve as the programmed connection. After erasure, the EPROM returns to its initial state and can be reprogrammed to a new set of values.

The fourth type of ROM is the electrically erasable PROM (EEPROM). It is like the EEPROM except that the previously programmed connections can be erased with an electrical signal instead of UV light. The advantage is that the device can be erased without removing it from its socket.

### Serial Access Memories:

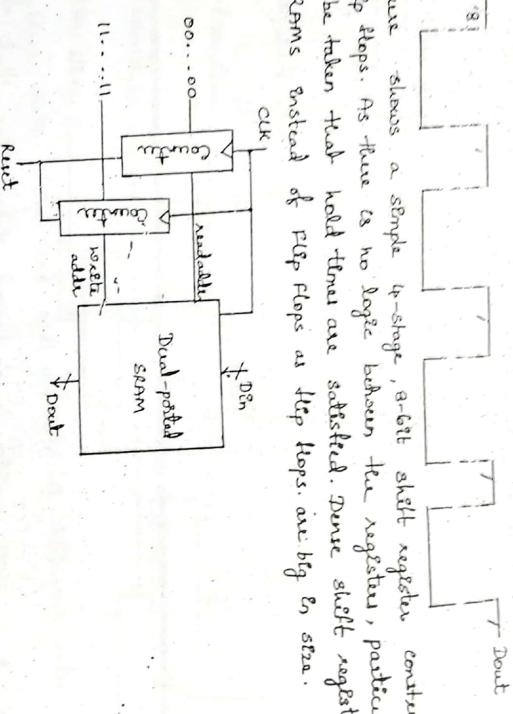
Using the basic SRAM cell and its registers, we can construct a variety of serial access memories including shift registers and queues. These memories avoid the need for external logic to back addresses for reading or writing.

### Shift Registers:

A shift register is commonly used in signal processing applications to store and delay data.

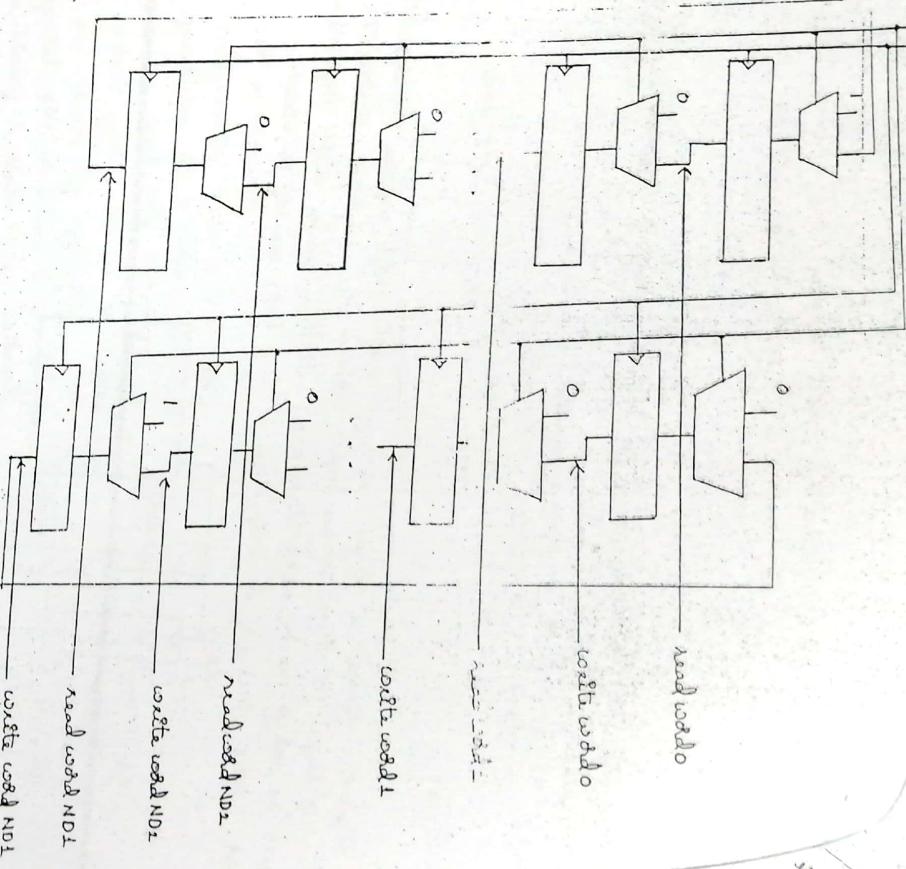


Above figure shows a simple 4-stage, 8-bit shift register constructed from 32 flip flops. As there is no logic between the registers, particular care must be taken that hold times are satisfied. These shift registers uses dual port RAMs instead of flip flops as flip flops are big & slow.



The RAM is configured as a circular buffer with a pair of counters specifying where the data is read and written. The read counter is initialized to the first entry and the write counter to the last entry in the array. The counters in an N stage shift register can use two lot of N bit registers to track which entries should be read and written. Again one is initialized to point to the first entry and the other to the last entry. These registers can drive the wordlines directly without the need for a separate decoder.

N hot register is a register which contains 'w' as loaded (in one flip flop) values Ex: for 4 bit register 1000 or 0001 is loaded.



- o
- o number of stages of delay.

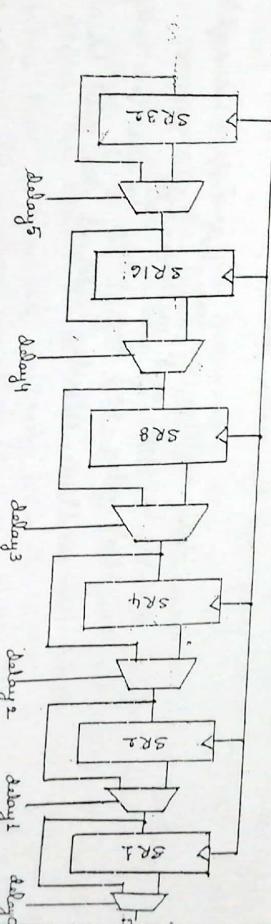
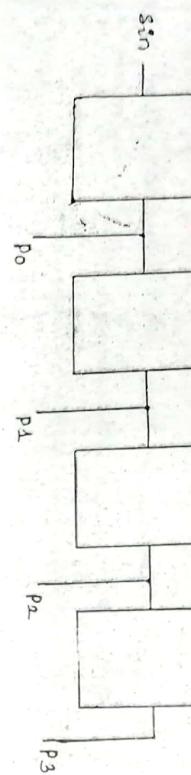
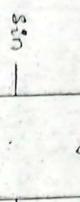


Figure shows a 4 stage tapped delay line that could be used in VLSI processing system. Delay blocks are built from 32, 16, 8, 4, 2 and 1 stage shift registers. Multiplexers control pass-around of the delay blocks to provide the minimum total delay.

Another variant is a serial/parallel memory. These are useful in signal processing and communications systems.

CLK



Above diagram shows 4-stage serial in parallel out (SISO) memory. It contains 4 stages.

Queues :

Queues allow data to be read and written at different rates.



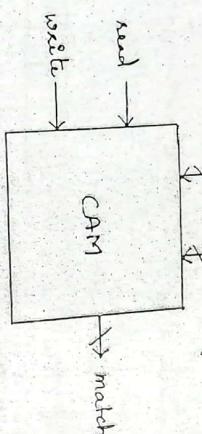
Above figure shows an interface to a queue. The read and write operations each are controlled by their own clocks that may be asynchronous. The queue emits the FULL flag when there is no room remaining to write data and the EMPTY flag when there is no data to read. Because of other system delays, some queues also provide ALMOST-FULL and ALMOST-EMPTY flags to communicate the impending state and half write or read requests. The queue internally maintains state and half write or read pointers indicating which data should be accessed next.

The pointers can be contained in 1-of-N hot registers.

First In First Out (FIFO) queues are commonly used to exchange data between two asynchronous streams. The FIFO is organized as a circular buffer. On next, the read and write pointers are both initialized to the first element and the FIFO is EMPTY. On a write, the write pointer advances to the next element. If it is about to catch read pointer, the FIFO is FULL. On a read, the read pointer advances to the next element. If it catches the write pointer, the FIFO is empty again.

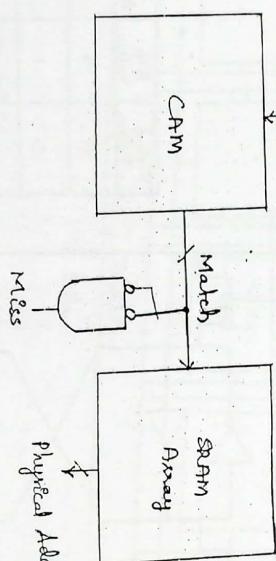
Last In First Out (LIFO) queue also known as stacks are used in applications such as subroutines or element stacks in microcontroller. The LIFO uses a single pointer for both read and write. On next, the pointer is initialized to the first element and the LIFO is EMPTY. On a write, the pointer is incremented. If it reaches the last element, the LIFO is FULL. On a read, the pointer is decremented. If it reaches the first element, the LIFO is empty again.

Addressable Memory :-



Above figure shows symbol for a content-addressable memory (CAM). The CAM acts as an ordinary SRAM that can read or written given address and data but also performs matching operations. Matching asserts a match line output to indicate if the item that receives a specified

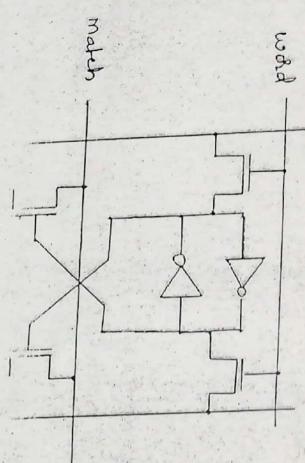
A common application of CAM is translation lookaside buffers (TLB). A common application of CAM is translation lookaside buffers (TLB).  
 In microprocessor supporting virtual memory, the virtual address is given as the key to the TLB CAM. If this address is in the CAM, the corresponding matchline is asserted. This matchline can serve as the enables to access a RAM containing the associated physical address. A NOR gate processing all of the matchlines generates a miss signal for the RAM.



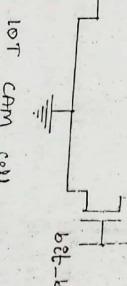
Translation: lookaside Buffer (TLB) using CAM

CAM cell:  
 10T and 9T CAM cells consist of a normal SRAM cell with additional transistors to put in the match. Multiple CAM cells in the same word are tied to the same matchline. The matchline is either precharged or pulled high as a distributed pseudo-MOS gate. The key is placed on the bottom. If the key and the value stored in the cell differ, the matchline will be pulled down. Only if all of the key bits match all of the bits stored in the

word of memory than the match line for that word remain high. The key can contain a "don't care" by setting both bit and bit-b-low.

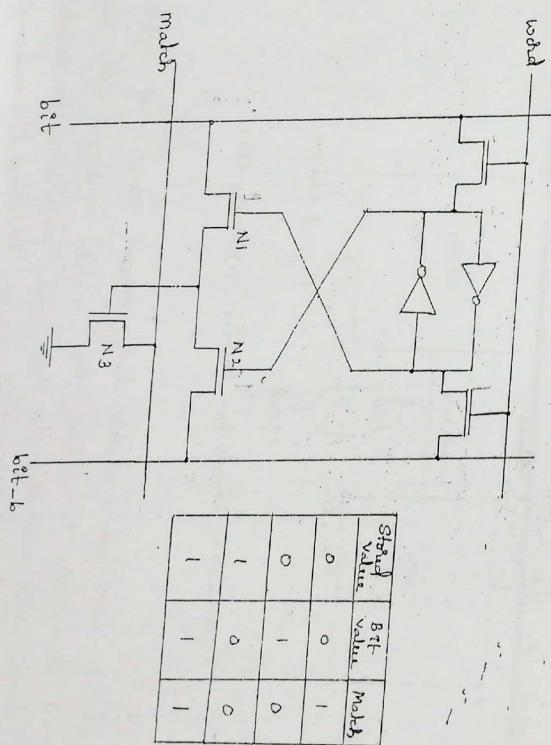


Stored value	Bit value	Match
0	0	1
0	1	0
1	0	0



1-bit CAM cell

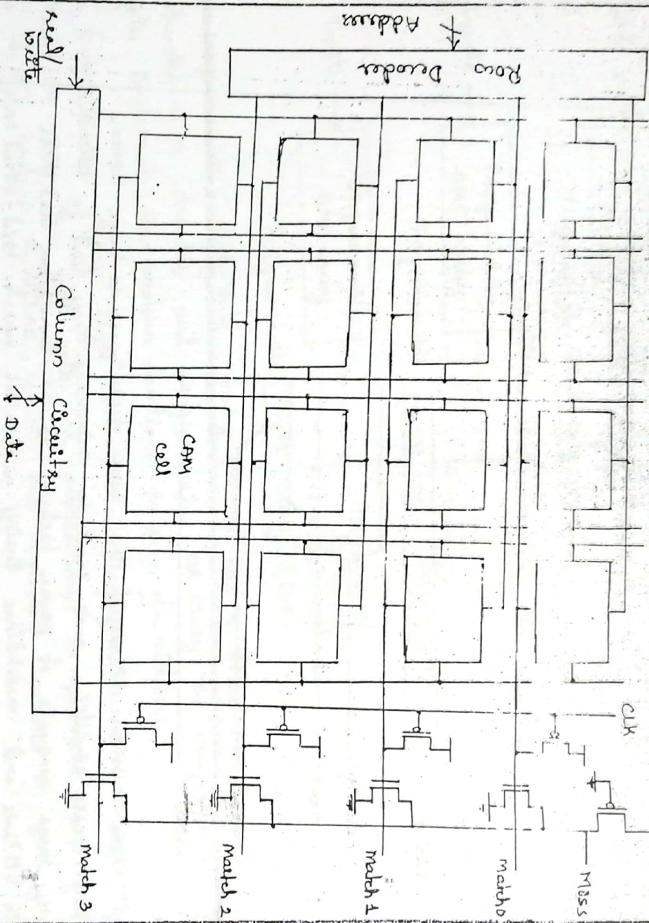
CAMs generally have about twice the area of SRAM cells. A CAM cell is designed with 9 transistors. N1 and N2 perform an XOR of the key and cell data. If the values disagree, N3 is turned on to pull down the wordline.



9T CAM cell

### CAM array

Like an SRAM, CAM array consists of an array of cells, a decoder, and column circuitry. Each row produces a dynamic matchline. The matchlines are packaged with the clocked pmos transistors. The match signal is produced with a distributed pseudo-nmos NOR.



When the matchlines are used to access a RAM, the monotonicity problem must be considered. Initially, all the matchlines are high. During CAM operation, the lines pull down, leaving at most one line asserted to indicate which row contains the key. The RAM requires a monotonically rising address. Strobed AND gates are used to derive the wordline as early as possible after the matchline has settled. The strobe can be timed with an inverter chain or replica delay line in the same way that the sense amplifier clock for a SRAM was generated. Clocked sense amplifiers consume power only when activated, but require a timing chain to activate at the proper time.

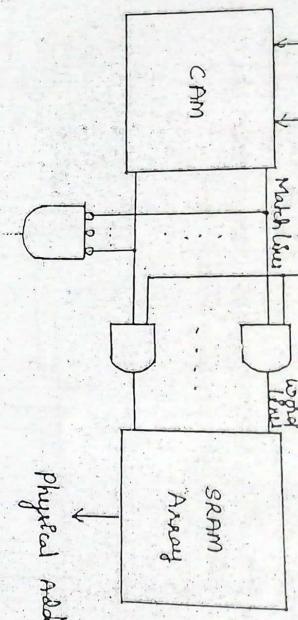
Vertical Address CLK

Matchline

Word line

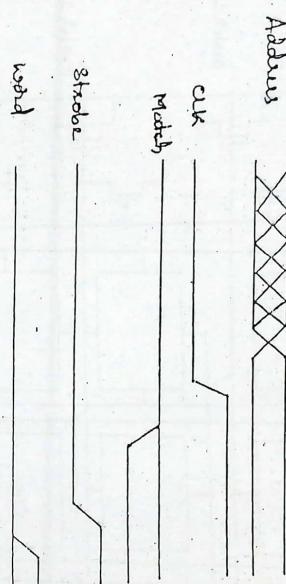
SRAM

Array



Resulting paths with mnemonic wordlines

miss

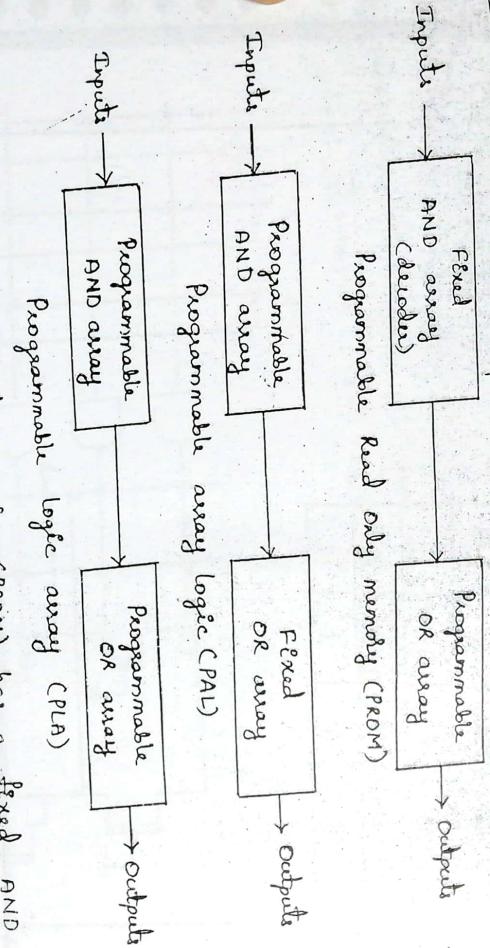


Large CAMs can use many of the same techniques as large RAMs, including sense amplifiers and multiple subarrays. They tend to consume relatively large amounts of power because of the transitions that take place in buffers and matchlines during a parallel search. Mid-sized CAM using a pmos match-line driver reduces the swing of matchlines and saves the power.

### (a) Programmable Logic Devices

Combinational programmable logic device (PLD) is an integrated circuit with combinational logic gates divided into an AND array and an OR array to produce an AND-OR sum of product implementation. There are three major types of combinational PLDs and they differ in the placement of the programmable connections in the AND-OR array.

Basic configuration of these PLDs:

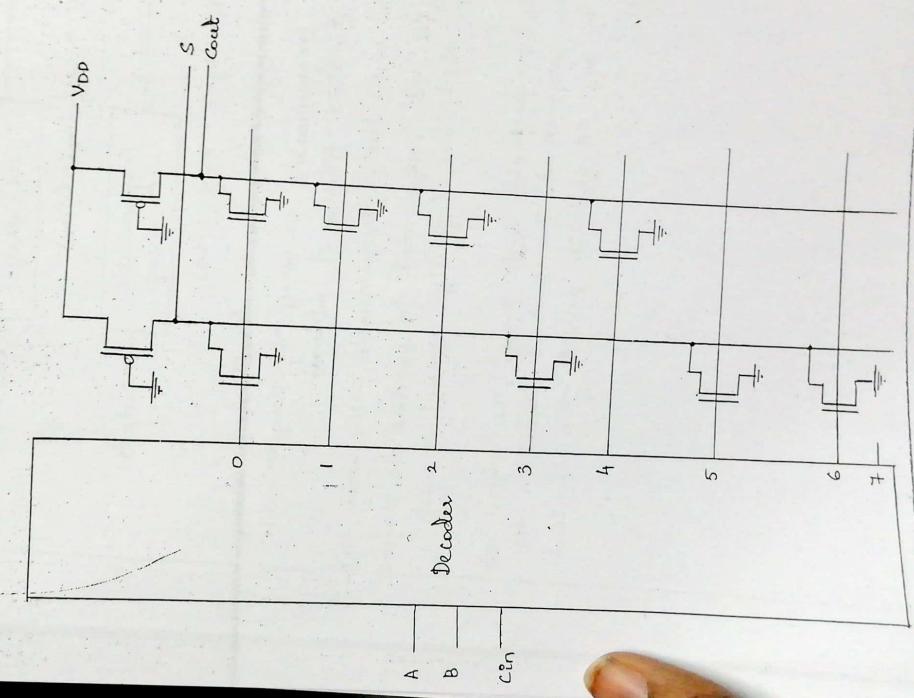


The programmable read-only memory (PROM) has a fixed AND array constructed as a decoder and programmable OR array. The programmable OR gates implement the boolean functions in sum of minterms. The programmable array logic (PAL) has a programmable AND array and a fixed OR array. The AND gates are programmed to provide the product terms for the boolean functions, which are logically summed in each OR gate.

The most flexible PLD is the programmable logic array (PLA) where both the AND and OR arrays can be programmed. The product terms in the AND array may be shared by any OR gate to provide the required sum of products implementation.

Implementation of Full adder using PROM:

A	B	C <sub>in</sub>	S	C <sub>out</sub>
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



and the following two boolean functions with a pin:

$$f_{\alpha}(a, b, c) = \min(\rho, s, t, r)$$

$\delta C$	$\delta C$	$\delta C$	$\delta C$
$\bar{D}_1$	$D_1$	$D_1$	$D_1$
$\bar{D}_2$	$D_2$	$D_2$	$D_2$
$\bar{D}_3$	$D_3$	$D_3$	$D_3$

$$F_1 = \overline{A}\overline{B} + \overline{AB} + \overline{BC}$$

$$\frac{1}{R_1} = \frac{\alpha c + \alpha s + bc}{\alpha c + \alpha b + bc}$$

卷之二

$$F_2 = \pi \delta \bar{c} + nc + \rho b$$

$$F_2 = \sqrt{c} + \sqrt{a}b + a\sqrt{b}c$$

$$\partial_\tau F_x = \overline{\partial c} + \overline{\partial} B + \partial \overline{B} \overline{c}$$

the sum number of product terms is

$$\text{and } F_2 = AB + AC + \overline{AB}C$$

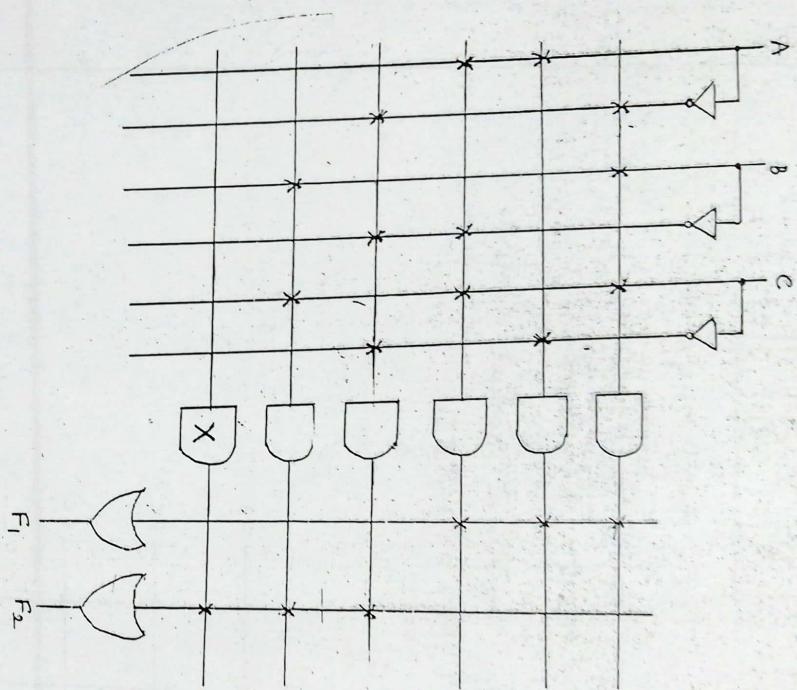
The circuit diagram consists of three logic gates. The first is an AND gate with inputs labeled  $\bar{A}$  and  $\bar{B}$ . Its output is connected to one input of a second AND gate, which has inputs labeled  $C$  and  $D$ . The output of this second AND gate is connected to one input of an OR gate. The other input of the OR gate is labeled  $E$ . The output of the OR gate is labeled  $F_1$ .

Scanned by CamScanner

1) Implement following functions with PAL:

$$F_1 = \bar{A}BC + A\bar{C} + A\bar{B}\bar{C}$$

$$F_2 = \bar{A}\bar{B}\bar{C} + BC$$

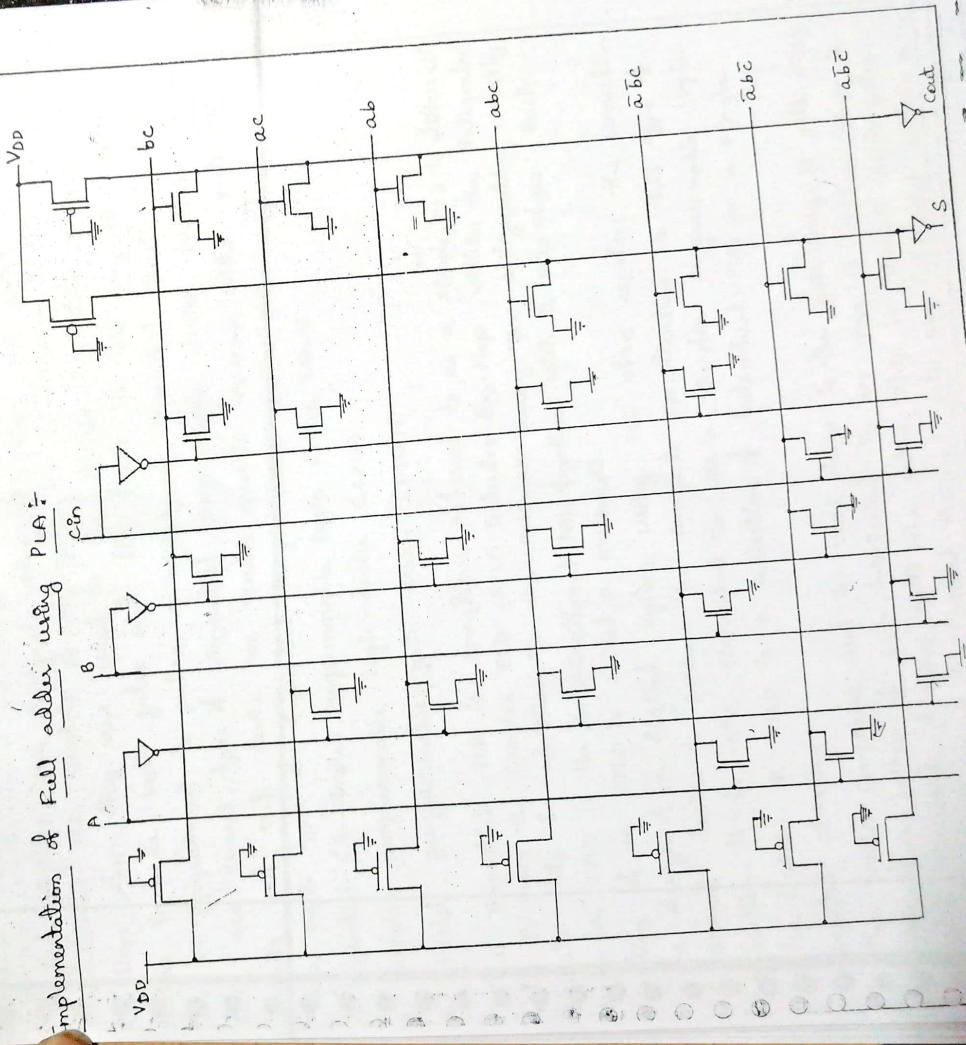


### Programmable logic array:

Programmable logic array (PLA) is similar to the PROM in concept. But the PLA does not provide full decoding of the variables. It does not generate all the minterms. The decoder is replaced by an array of AND gates that can be programmed to generate any product term of the input variables. The product terms are then connected to sum of products for the required Boolean function.

The most straightforward PLA uses a pseudo-NMOS NOR gate. Advantages of this PLA include simplicity and small size. Disadvantages include the static power dissipation of the NOR gates and slow pull-up response.

### Implementation of Full adder using PLA:



when implementing a combinational circuit with a PLA, careful simplification must be undertaken in order to reduce the number of distinct product terms, since a PLA has a finite number of AND gates. This can be done by simplifying each boolean function to a minimum number of terms. The number of literals in a term is not important since all the input variables are available anyway. Both the true and complement of each function should be simplified to see which one can be expressed with fewer product terms and which one provides product terms that are common to other functions.

## Programmable Array Logic:

programmable array logic (PAL) is a programmable logic device with AND OR array and a programmable AND array. Because only the AND array is programmable, the PAL is easier to program, but is not as flexible as the PLA.

When designing with a PAL, the boolean functions must be simplified to fit into each section. Unlike the PLA, a product term cannot be shared among two or more OR gates. Therefore, each function can be simplified by itself without regard to common product terms. The number of product terms in each section is fixed.

## Sequential Programmable devices:

Digital systems are designed using flip-flops and gates. Since the combinational PLD consists of only gates, it is necessary to include external flip flops when they are used in the design. Sequential programmable devices include both gates and flip flops. In this way, the device can be programmed to perform a variety of sequential circuit functions. There are several types of sequential programmable devices available commercially and each device has vendor specific variant within each type. Three major types of sequential programmable devices are

- (i) Sequential (or simple) programmable logic device (SPLD)
- (ii) Complex programmable logic device (CPLD)
- (iii) Field programmable gate array (FPGA)

The sequential PLD is sometimes referred to as a simple PLD to differentiate it from the complex PLD. SPLD includes flip-flops within the integrated circuit chip in addition to the AND-OR array. The configuration mostly used for SPLD is the combinational PAL together with D flip flops. Each section of an SPLD is called a macrocell.

The design of a digital system using PLD often requires the connection of several devices to produce the complete specification. For this type of application, it is more economical to use a complex programmable logic device (CPLD). A CPLD is a collection of individual PLDs on a single integrated circuit.

The basic component used in VLSI design is the gate array. A gate array consists of a pattern of gates fabricated in an area of silicon that is repeated thousands of times until the entire chip is covered with gates. Arrays of one thousand to hundred thousand gates are fabricated within a

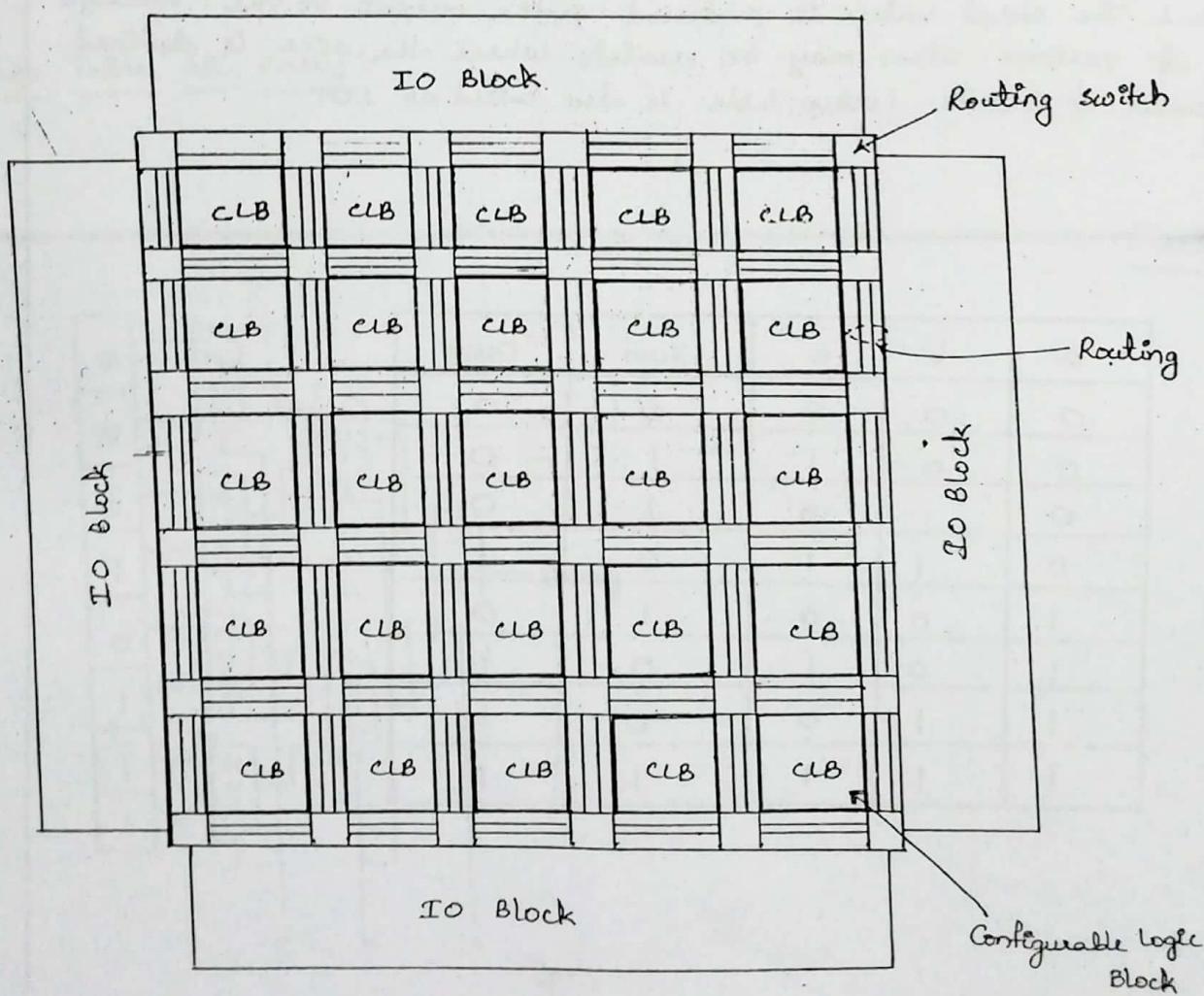
single to chip depending on the technology used. The design with gate array requires that the customer provide the manufacturer the desired interconnection pattern.

A field programmable gate array (FPGA) is a VLSI circuit that can be programmed in the user's location. A typical FPGA consists of an array of hundreds or thousands of logic blocks, surrounded by programmable input and output blocks and connected together via programmable interconnections.

## Programmable Gate Arrays:

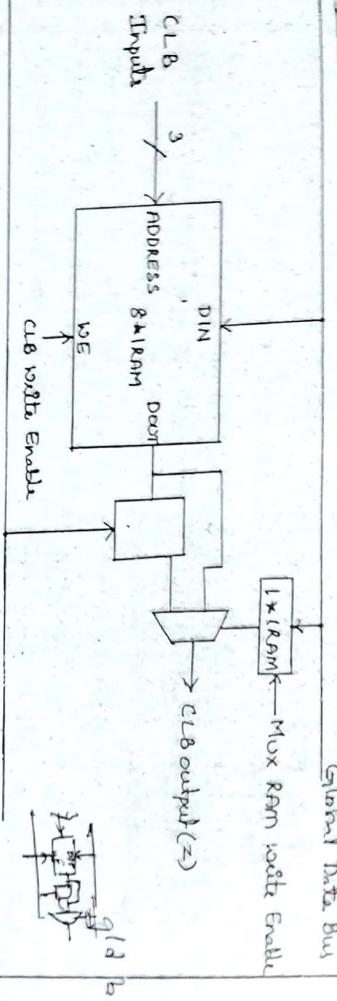
A Field Programmable Gate Array (FPGA) is a programmable logic device that supports implementation of relatively large logic circuits. FPGAs can be used to implement a logic circuit with more than 20,000 gates. The general structure of an FPGA contains three main types of resources. They are logic blocks, I/O blocks for connecting to the pins of the package, and interconnection wires and switches. The logic blocks are arranged in a two-dimensional array, and the interconnection wires are organized as horizontal and vertical routing channels between rows and columns of logic blocks. Programmable connections also exist between the I/O blocks and the interconnection wires. The routing channels contain wires and programmable switches that allow the logic blocks to be connected in many ways.

## General structure of an FPGA:



A typical FPGA logic block consists of look-up tables, multiplexers, gates and flip-flops. The look-up table is a truth table stored in a SRAM and defines the combinational circuit functions for the logic block.

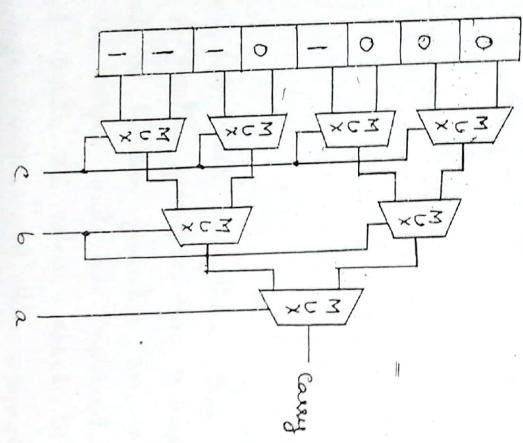
Simple FPGA logic cell :-



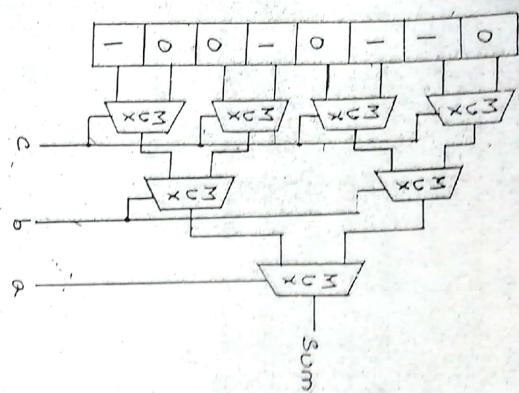
Look-up table contains storage cells that are used to implement a small logic function. Each cell is capable of holding a single logic value either 0 or 1. The stored value is produced as the output of the storage cell. LUTs of various sizes may be created, where the size is denoted by the number of inputs. Look-up table is also called as LUT.

Full adder :-

a	b	c	Sum	Carry
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



Look up Table for carry:



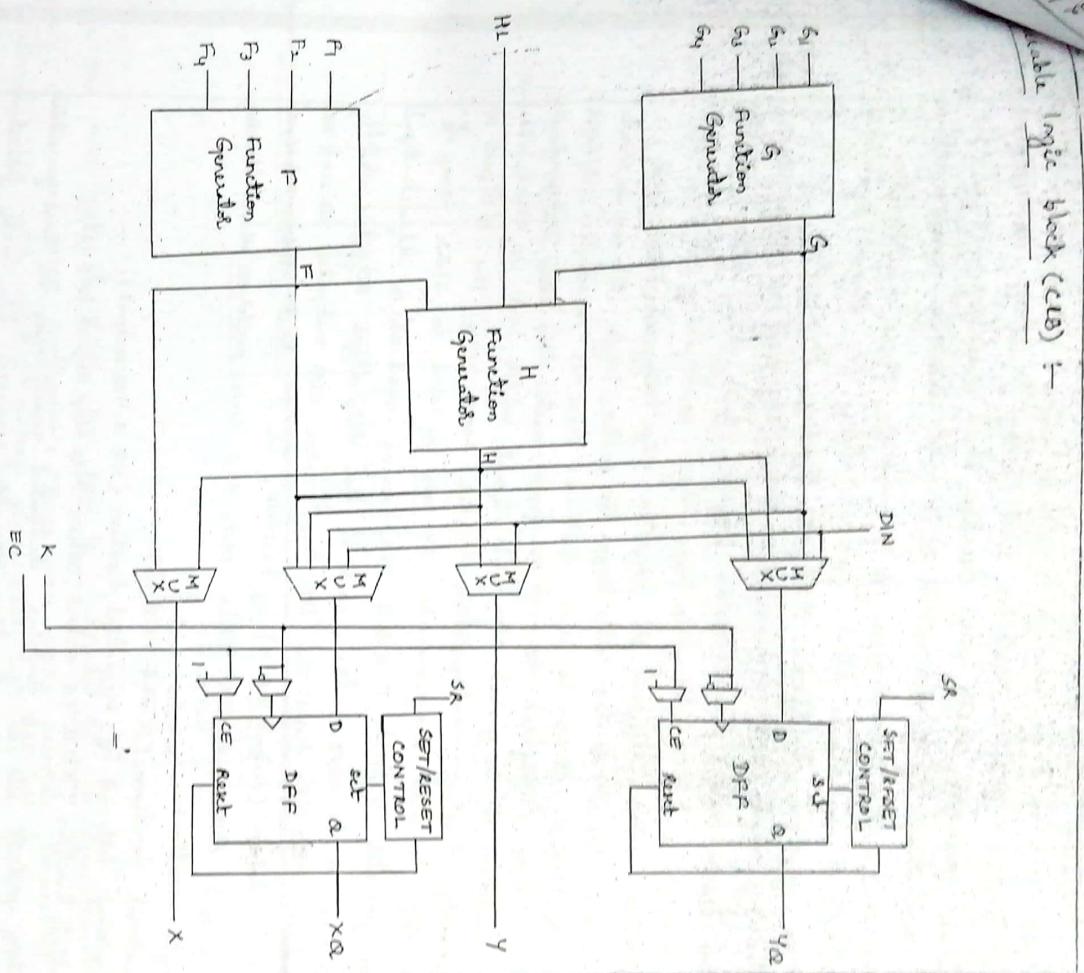
Field programmable. Auto designs use high circuit densities in order to construct ICs. They are completely programmable even after they are placed in the field. Two basic versions of FPGA are SRAM based and EEPROM based. First one uses fuses or antifuses to permanent program interconnect and programmable logic. They are one time programmable. Second one uses small static RAM cells to customize routing and logic functions. LUTs are placed in this RAM cells. Compared to conventional ICs, FPGAs are slower and are less efficient in area and power. However, the programming cost is virtually zero. The cost is moderate to high. As initial manufacturing costs for custom ICs increase and FPGA costs decrease, FPGAs become attractive. Example of FPGA based on SRAM technology is XILINX Spartan series, FPGA.

XILINX Spartan series FPGA

These FPGAs are implemented with a flexible, regular and reprogrammable architecture of configurable logic blocks (CLBs), interconnected by programmable routing channels and surrounded by a perimeter of programmable input output blocks (IOBs).

FPGAs can be programmed in two modes. They are master serial mode and slave serial mode. FPGAs can read its configuration data from an external serial PROM in master serial mode. External device such as computer can write data into FPGA in slave serial mode.

able logic block (cells) :-



In each CLB there are three look up tables which are used here as F, G and H function generators, two D flip flops and group of multiplexers. Each F and G function generators are 16x1 memory look-up tables having 4 inputs one output and can implement any boolean function upto four independent inputs. The values stored in memory cells during configuration determine logic function and thus propagation delay is independent of the logic function realized. H function generator is a three input function generator which can implement any boolean function of its input. Two inputs of H function generator are the outputs from F and G function generator. The other input (H1) is coming outside of the CLB.

signals from function generators can exist through two outputs for  $H$  which  $X$  output and or or  $H$  through  $Y$  output. Thus a CLB can implement  $4 \times 2$  input function up to 4 variables each (using  $F$  and  $G$ ), or any function of four variables with some function of 5 variables (using  $F(G)$  one and combination of  $H$  with  $G(F)$ ), or it can realize some function to nine variables (using combination of  $F, G$  and  $H$ ). Thus realization of various functions in a single block reduces delay in signal path and no minimizes number of blocks required, thereby increasing both speed and density. Each CLB contains two edge triggered D flip flops with common clock and clock enable inputs. Clock can enter directly inverted through a MUX into flip flops, thus configuring D flip flops either positive edge triggered or negative edge triggered. The clockable / flip enable input of the flip flops is active high. It can be made ways high (enable) or it can be user defined. There is another set/reset input which can be configured to set or reset each flip flop independently. There is a provision of a global set reset signal which set or clear the flip flops globally during configuration in the same way. The D input of the flip flops is also programmable. It can be used to store any of the function generator outputs  $F, G$  or  $H$ . DIN can be used as a direct input to the flip flops. A 4:1 MUX before each of the flip flops selects which signal to drive the flip flops. The flip flops derive CLB outputs  $X_0$  and  $Y_0$ . Each output is driven (output  $\times 2^{-4}$ ) by 2:1 MUX, for  $X$  output either  $F$  or  $H$ , and for  $Y$  output either  $G$  or  $H$ . Thus each CLB can perform the following types of logic operations.

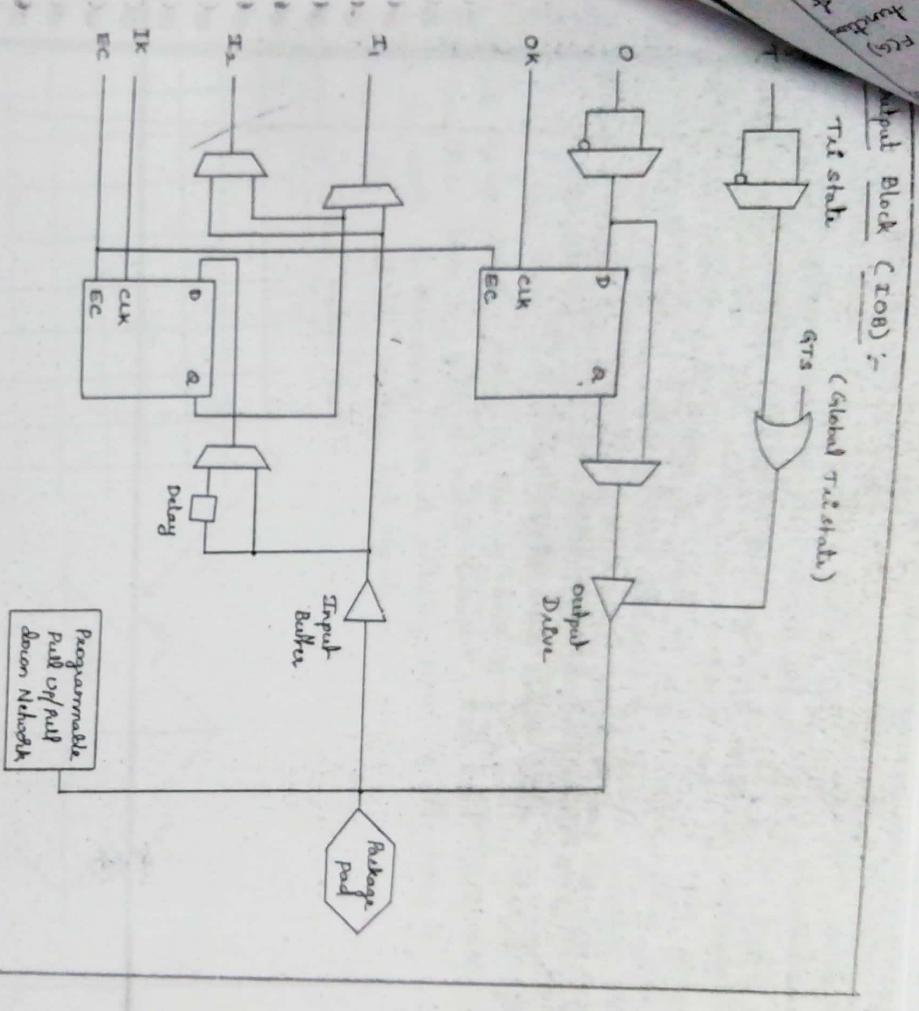
- Combinational functions ( $X$  and  $Y$  output)
- Combinational followed by sequential function ( $X_0$  or  $Y_0$  output)
- Sequential function ( $X_0$  &  $Y_0$  output, when flip flop input is DIN).

In addition each of the CLB contains dedicated carry logic for fast generation of carry signal, therefore saving function generator resources for carry generation.

Output Block (TOB) :-

Tri state

(Global Tri state)



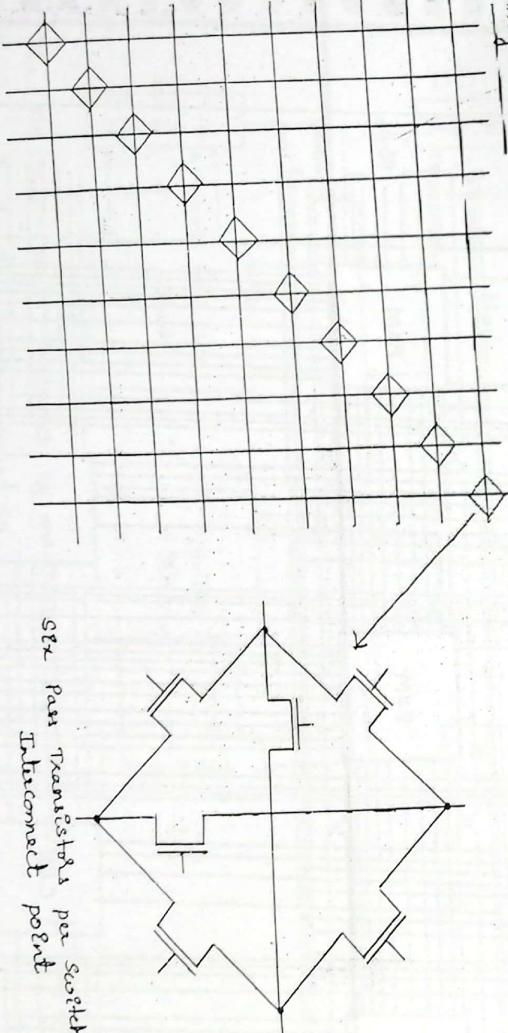
There are two types of TOBs, one type is dedicated for configuration of FPGAs, others are user configurable. These user TOBs provide connection for internal CLB with external package pins. Each TOB has one package pin which can be configured as input, output or bi-directional signals. The input signal to the TOB can directly go to routing channels and connect CLB inputs or it can go to the register before going to the routing channels. The register can be configured as an edge sensitive flip flop or level sensitive latches. The clocks can be direct or inverted. The input data path has one tap delay element which can treat partial delay to adjust set up a hold time of flip flops, if necessary. The Tx and Rx signals travel through TOBs each carrying the repeat after direct or inverted. Output signal (Q) need to connect from CLB outputs to go as output to the external package pin can be inverted inside TOB, can pass directly to output buffer, as it can be stored in an output register and then pass to the output buffer. The tri state input, when high places the output buffer

high impedance state. The user controlled T input can go directly or muted to the tri-state control after being connected with global tail state (GTS) and GTS signal, which is common to all user I/Os, is made high during configuration of FPGA in order to put all the package pins, kept program pins, in high impedance state. Then GTS automatically goes low after configuration is over in order to ensure that FPGA connects with external circuitry only after configuration is done. Programmable pull-up/pull-down network is used for keeping unused pin to ground level in order to minimize power and increase noise immunity. Separate clock signals are given to input and output flip-flops. Common enable signal is used.

channel around the cib have three types of interconnects: single length, double length and long lines. Each psm consists of pass transistors having programmable gate voltage and driving on, establishing connection between lines connected at its source and drain. Each connection between lines connected at its source and drain contains six pass transistors in series to route

signals in all possible directions. signals have the greatest interconnect facility since it runs through each psm block, but due to this phenomenon it also induces highest amount of delay among the three types of interconnects. Double length lines are twice longer than single length lines and they bypass alternate psms. so the delay is less but interconnect facility is also less. long lines from a grid of metal intersect, they run along the length or width of the array, without entering psms. so its delay is minimum and used for routing fast signals.

Programmable switch matrix (PSM):



6x6 Pass Transistors per Switch Matrix

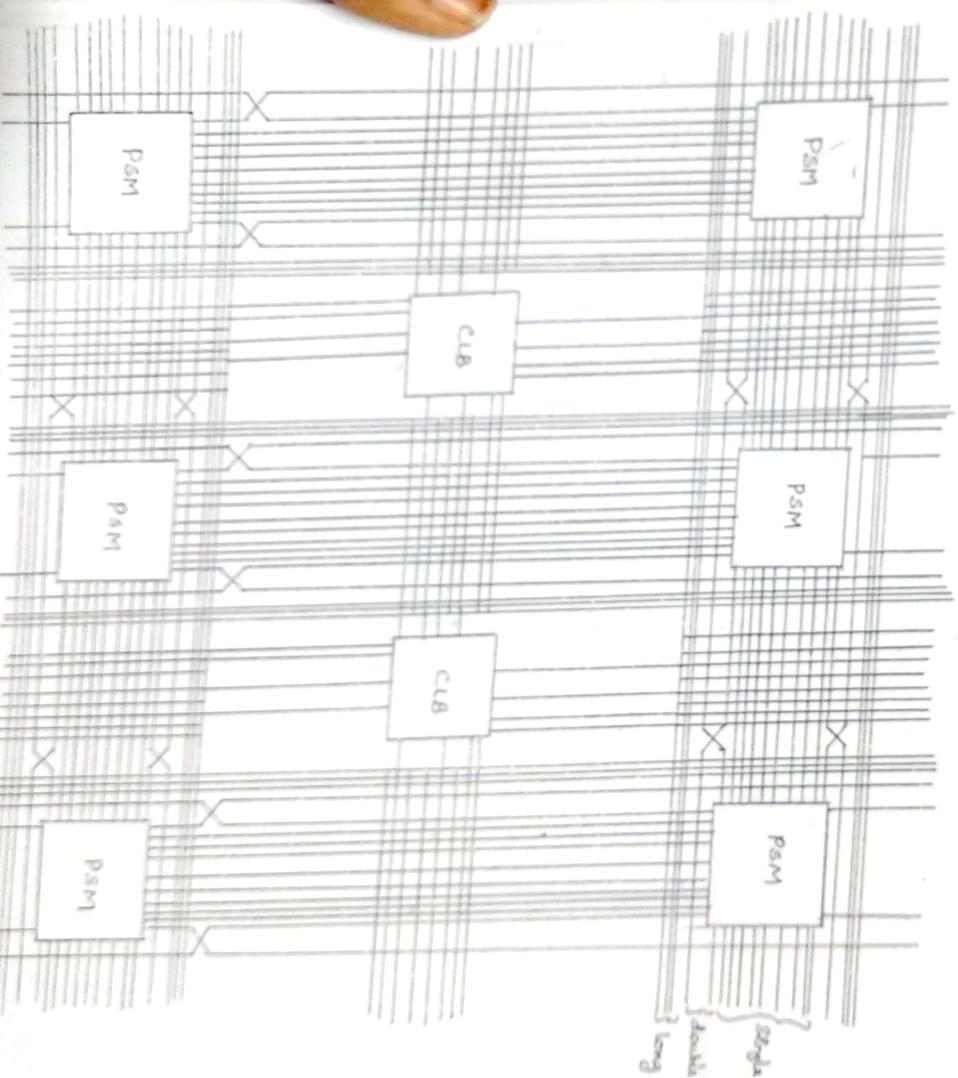
### CLA Routing channels

Routing channels provide paths for interconnecting inputs and outputs of CLAs and IOBs. The internal routing channels are metal segments with reprogrammable pins and programmable switch matrices (PSM).

There are three types of routing channels.  
1) CLB routing channels which run along each row and column of CLB may.

2) IOB routing channels which form a ring, called a via ring, around the CLB array. It connects the IOB with the CLB routing channels.

3) Global routing channels for routing global signals like clocks with minimum delay and skew.



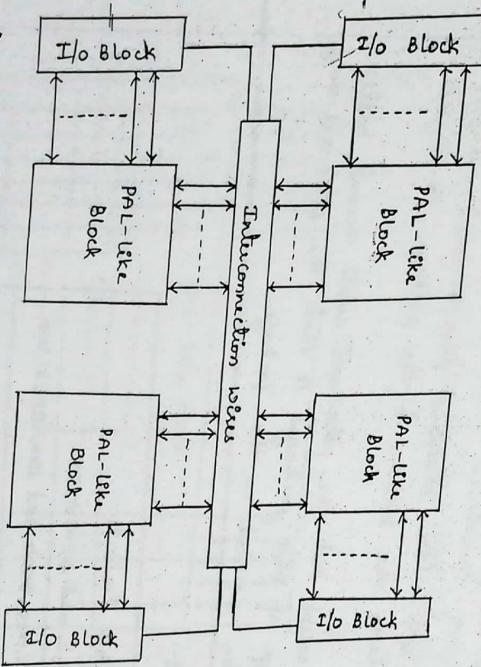
### Programmable logic devices (PLDs)

be programmed

and PALs are useful for implementing a wide variety of small outputs, either multiple PLAs or PALs can be employed or else a more sophisticated type of chip, called a complex programmable logic device, can be used.

A CPLD compresses multiple circuit blocks on a single chip, with internal

similar to a PLA in a PAL. CPLD includes PAL like blocks that are connected to a set of interconnection wires. Each PAL like block is connected to I/O block. I/O block is connected to input and output pins.



The PAL-like block includes 16 macrocells, each consisting of OR gate, with OR gate inputs between 5 and 20.

A section of CPLD:

In the figure below the wiring structure and the connections to a PAL-like

are shown. In this example the PAL-like block includes 3 macrocells, each consisting of a four input OR gate. The OR gate output is connected to XOR gate. One input to XOR gate can be programmed to 1 or 0, then the XOR gate complements the OR gate output, and if 0, then the XOR gate has no effect. The macrocell also includes a flip flop, a multiplexer, and a tristate buffer. The flip flop is used to store the output value produced by the OR gate. Each tristate buffer allows each pin to be used either

### FPGA Packaging:

FPGAs are available in a variety of packages. They are

- 1) PLCC - Plastic Leaded Chip carrier.
- 2) QFP - Quad Flat Pack.
- 3) PGA - Pin grid array
- 4) BGA - Ball grid array.

PLCC package has pins that wrap around the edges of the chip on all four sides.

The QFP package pins extend outward from the package. This package has pins on all four sides. The QFP's pins are much thinner than those on a PLCC, which means that the package can support a large number of pins.

A PGA package may have upto a few hundred pins in total, which extend straight outward from the bottom of the package, in a grid pattern.

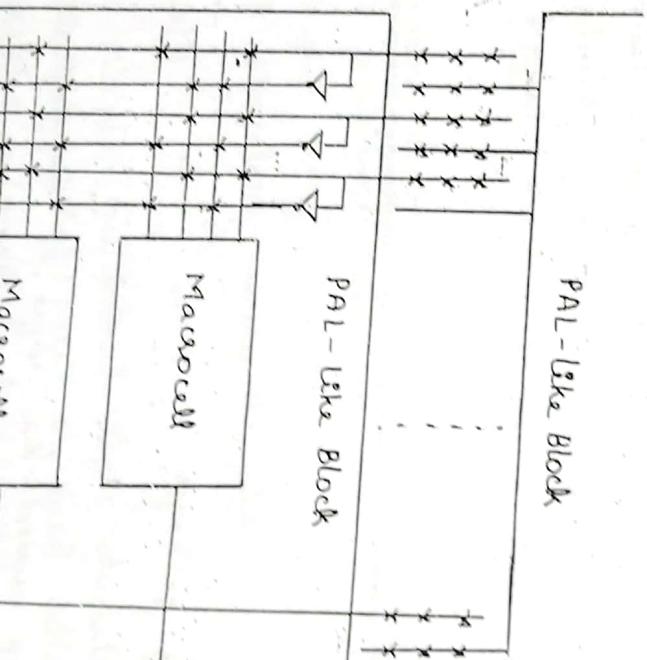
BGA is similar to PGA except that the pins are small round balls, instead of posts. The advantage of BGA packages is that the pins are very small, hence more pins can be provided on the package.

### Commercial FPGA products:

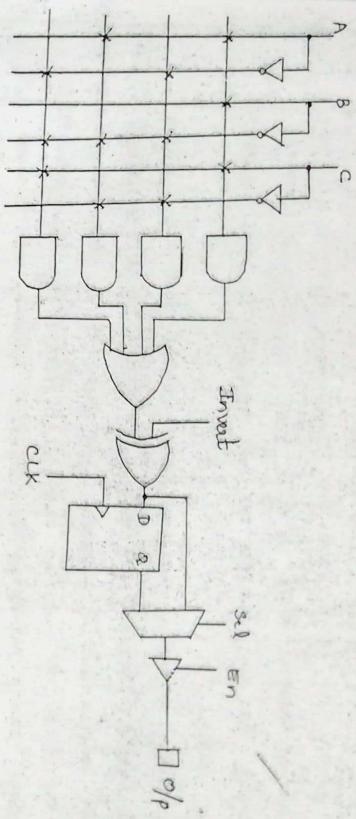
Manufacturer	FPGA Products
Actel	Act 1,2 and 3, MX, SX
Altera	PLEX 6000, 8000 and 10K APEX 20K
Atmel	AT6000, AT40K
Lucent	ORCA 1,2 and 3
QuickLogic	PASIC 1,2 and 3
Vantis	VPI
Xilinx	XC3000, XC4000, XC5200,
	Virtex

output from the CPLD or as an input responding the state buffer is enabled, and if the pin is to be used as an input, then acting as a switch that is turned off. Then can drive a signal onto the pin, which using the interconnection wiring.

The interconnection wiring contains program to connect the PAL-like blocks. Each of the has to some of the vertical wires that it crosses number of switches are shown to provide segments without wasting many switches. In as an input, the macrocell associated with it therefore wasted. Some CPLDs include additie macrocells and the interconnection wiring to in such situation.



### Latches :-



$$F = \bar{ABC} + \bar{B}\bar{C} + A\bar{B}C + A\bar{B}\bar{C}$$

En	Invert	sel	output
0	*	*	=
1	0	0	F
1	0	1	delayed F
1	1	0	F'
1	1	1	delayed F'

### CPLD Packaging :-

CPLDs are available in the following packages.

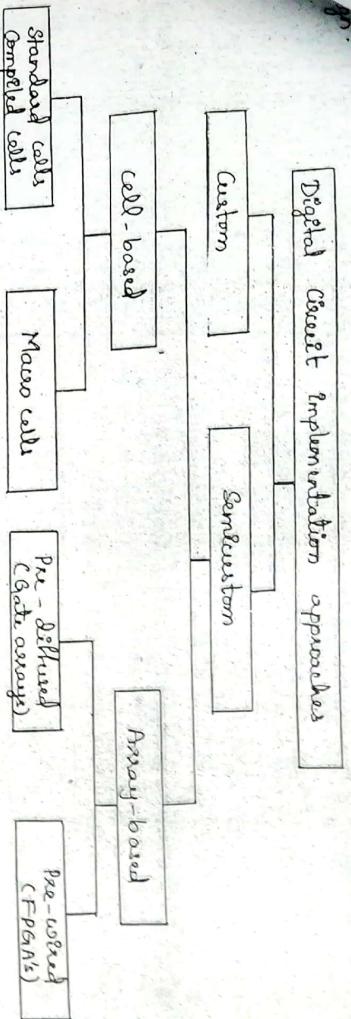
i) Plastic leaded chip carrier - PLCC

ii) Quad flat pack - QFP

The PLCC package has pins that wrap around the edges of the chip on all four sides. The socket that houses the PLCC is attached by solder to the circuit board, and the PLCC is held in the socket by friction. The QFP package has pins on all four sides, and they extend outward from the package, with a downward-curving shape.

### Approaches:

Following figure shows various implementation approaches of VLSI



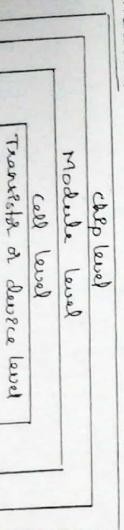
In the custom design approach, each individual transistor is designed and laid out manually. The main advantage of this method is that the circuit is highly optimized for speed, area & power. This approach is also known as full custom design approach. In a full custom design approach the semiconductor chips are Application Specific Integrated Circuits (ASICs) which are designed specifically for a given application or application domain.

In semi custom design approach, the majority of the chip is designed using a group of predefined cells called as standard cells and rest are designed manually. A gate array is an IC chip on which gates are placed in matrix form without connections among the gates.

### Standard cells:

standard cells are pre-defined logic elements used in the circuit. The design methodology that uses standard cells is known as cell based design methodology. A standard cell library is one of the foundations upon which the VLSI design approach is built. A standard cell is designed either to store information or perform a specific logic function. The type of standard cell created to store data is referred to as a sequential cell. flip flops and latches are examples of sequential cells. The type of standard cell used to perform logic operations on the signals presented on its inputs is called combinational cell.

### Level of abstraction:



## Commercial CPLD Products

Manufacturer	CPLD Product
Altera	Max 5000, 4000 and 1000
Atmel	ATF, ATV
Cypress	FLASH310, Ultra 3100
Lattice	EP LSI 1000 to 8000
Philips	XPLA
Vantis	MACH 1 to 5
Xilinx	XC9500

### Comparison of FPGA and CPLD

FPGA	CPLD
1) FPGA architecture contains gate array like structures.	1) CPLD architecture contains PAL like structures.
2) Its density range is from medium to high.	2) Its density range is from low to medium.
3) FPGA speed depends on the application.	3) CPLD speed is very high.
4) For interconnection purpose routing channels are used.	4) For interconnection purpose switchable connections are used.
5) In FPGA power consumption is medium.	5) In CPLD power consumption is high.
6) Once supply is removed, data is losted in FPGA, i.e. they are volatile.	6) CPLD contains on chip non-volatile memory.
7) FPGAs are internally based on look-up tables (LUTs).	7) CPLDs have the logic functions with sea-of-gates.

## Influencing low power VLSI design:

overall power consumption of an integrated circuit can be influenced at all levels of its design. In fabrication technology, circuit optimisation, logic design, control and clocking strategies, architectural partitioning and layout, and the underlying system's algorithm. To reduce the power consumption efficient methods are available.

bd :-

The best strategy is to reduce  $V_{DD}$ . Variation in  $V_{DD}$  leads to quadratic range in the power delay product.

### threshold voltage:

By reducing  $V_t$  to lower levels it would be possible to reduce  $V_{DD}$  even further.

### compromising for low speed:

As the enduring move from the standard 5V process to ones with lower supply voltages, design engineers need to compensate for the loss of performance. There are two architectural approaches, first apply the standard speed optimisation techniques, second use parallelism. Parallelism is the standard technique for increasing the overall speed of a circuit. The idea is using parallelism is simply to have more operations being conducted at the slower speed to achieve the same overall performance.

### voltage scaling:

A final way of reducing power loss connected with supply voltage is to reduce voltage scaling.

### reduce C:

The second strategy is to reduce capacitance. This comes naturally with smaller feature sizes.

### partition blocks:

It is best to partition large blocks into smaller ones. As a general rule each block depends on the capacitance of that block. Power calculation for each block depends on the capacitance of that block. When large blocks are partitioned into smaller blocks, only one block is operated at a time due to their product of activity and capacitance is reduced.

### Clocks and Control:

In architectures with distributed processing there should be global control and clock signals, but the global distribution network has a very large capacitance and is switched frequently. To overcome this several