

Chapter1 Introduction

1.1 Background

The memory is indispensable for the embedded system. However, it will occupy large area. The memory hierarchy is used in the embedded system to improve the performance. The basic component of memory hierarchy is Cache and Translation Lookaside Buffer (TLB). Cache and TLB can store the data which is used recently. They also provide a fast data comparison to reduce the access time from the microprocessor to the memory.

Cache and TLB can store and compare data. They both contain Content-Addressable Memory (CAM), or called associative memory. CAM is also composed of SRAM, so it has reading operation and writing operation. The difference between CAM and SRAM is the additional comparison part for CAM. CAM has additional circuit for comparison. The comparison is made up of XOR logic gates, which are composed of three or four pass transistors. CAM can make comparison in parallel for each word line by dynamic comparison, so its speed of comparison is fast.

According to the match line scheme of CAM, there are two types of match line scheme. In the comparison mode, any one bit of stored data is mismatch with any one bit of search data. The match line would be discharged to ground. This is called NOR type CAM. Another one is called NAND type CAM. The match line of NAND type CAM would be discharged to ground only when all bits of stored data are match with all bits of search data.

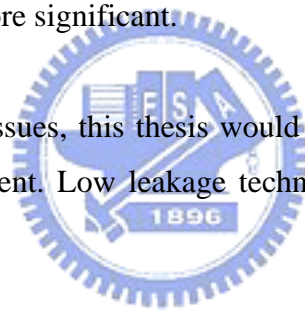
1.2 Motivation

The power consumption gradually becomes the other important key issue for embedded system. As the CMOS technology scales down, the power issue is not neglected. On the other hand, the leakage current becomes a significant problem and can not be neglected in deep submicron technology.

Some statistic data shows that the power consumption of Cache and TLB have a large percentage for embedded system. How to reduce this part of power dissipation is very important. Many power-efficient architectures and methods are provided to reduce the dynamic power for Cache and TLB. The main component of Cache and TLB is CAM. Because CAM can compare in parallel, the power consumption is serious for Cache and TLB. Improving the architecture of CAM is also a feasible way.

On the other hand, Cache and TLB are both storage devices. Their leakage current is a serious problem. In deep-submicron technologies, sub-threshold leakage is the critical component among all kinds of leakage currents. Sub-threshold leakage would increase, as the threshold voltage reduces. The technology scaling down will decrease the supply voltage, and the decreasing supply voltage would make the performance degrade. Therefore, the threshold voltage scales down to satisfy the requirement of the speed. Hence, the influence of the leakage current is more and more significant.

According to the above issues, this thesis would focus on reducing the dynamic power consumption and leakage current. Low leakage technical is used for pre-comparison CAM and applied to TLB.



1.3 Organization

The organization of this thesis is as follows. An overview of memory hierarchy is introduced in Chapter2. Memory hierarchy will introduce Cache. Translation Lookaside Buffer (TLB) and Content-Addressable Memory (CAM) is also presented in this section. Besides, the low power design for Cache and TLB would be described in this chapter, too. Some kinds of CAM would be depicted in this section.

The concept of the pre-comparison is proposed, and the pre-comparison circuit is implemented in Chapter3. Each memory comparison would pre-compare some bits of stored data and searching data. The remaining bits are decided to compare or not through the result of pre-comparison. It can efficient reduce some times of mismatch comparison. The less

comparison can decrease the dynamic power dissipation.

A further improvement of pre-comparison CAM is described in Chapter4. Enhancement of pre-comparison CAM reduces unnecessary XOR gates to lower the match line capacitance. Simulation result will compare different bits of pre-comparison CAM. Two kinds of access time would be measured in this chapter. Two types of power delay product are calculated to make comparison for conventional CAM. The optimal bit of pre-comparison CAM is decided.

A low leakage design technique is implemented in Chapter5. Power gating and dual vdd are used to apply to SRAM, CAM. Using pre-comparison architecture with the techniques of power gating as well as dual vdd is integrated and applied to TLB. Finally, the overall conclusions would be presented in Chapter6.



Chapter 2

Overview of Memory Hierarchy

As the technology scales down, the effect of the power consumption is more and more important. Low power design in the digital integrated circuit becomes a key issue. Especially, large percentage of power is consumed from memory of embedded system. Therefore, how to design low power memory is very important. In Sec. 2.1, memory hierarchy will be described briefly. Cache would be introduced in this section. Sec. 2.2 will discuss the power consumption and low power design for Content-addressable Memory (CAM). Sec. 2.3 demonstrates virtual memory, Translation Lookaside Buffer (TLB) and Memory management unit (MMU). Sec. 2.4 depicts the component of TLB, and the viewpoint how to reduce the power consumption. Sec. 2.5 would introduce some kinds of CAM. Sec. 2.6 is the summary of this chapter.

2.1 Memory Hierarchy



As the speed of CPU increases faster, the latency between the CPU and memory is more and more severe. Fig. 2.1, because the speed of the memory access increases very slower, the gap of speed between the CPU and memory is larger. When the CPU speeds up with 35%, the memory only speeds up 7%. So, the difference of between the CPU and memory is more and more serious from 1980 to 2005. In order to reduce the difference, there is a new method to be provided. It is called memory hierarchy. [1]

Its concept is to use the principle of locality. The principle of locality means that the programs access a relatively small portion of their address space at any instant time. The programs tend to reuse the data and instructions where they have been used recently. There are two different types of locality. One is temporal locality or locality in time. If an item is referenced, it would be referenced again in the near future. For example, the program encounters a loop instruction. Some instructions and data would be reused repeatedly. This is called temporal locality. The other is called spatial locality or locality in space. If an item is referenced, the neighbor of the item could be referenced. This is spatial locality. Fig. 2.2 is a

multi-level memory hierarchy, which includes typical sizes and access speed.

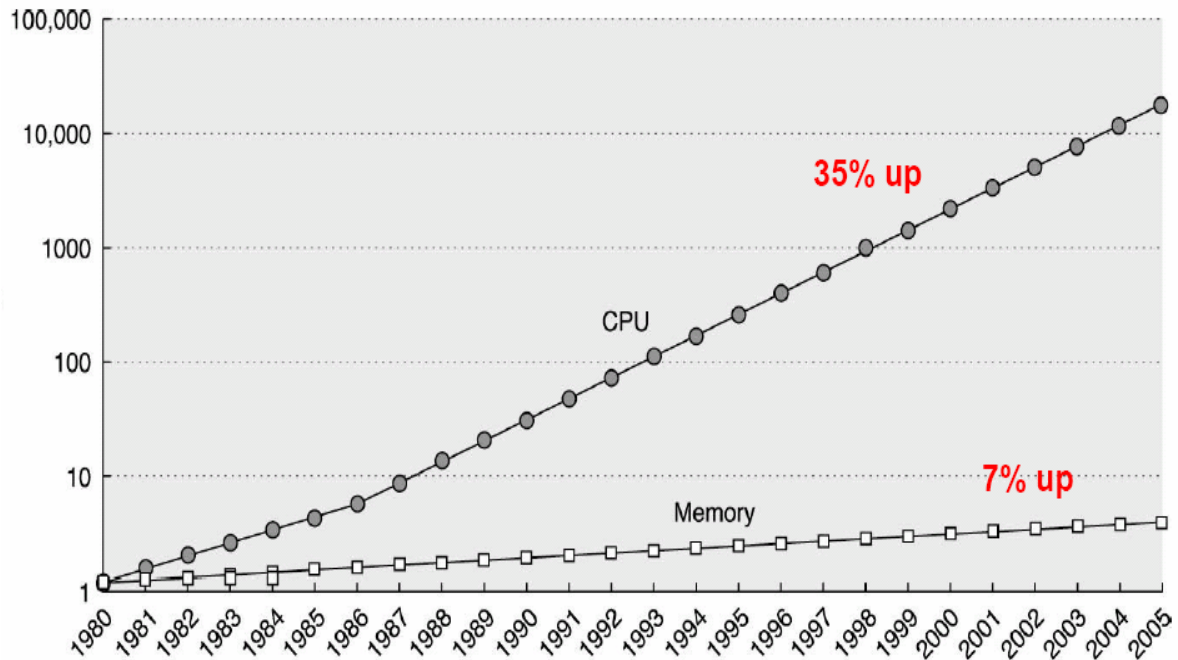


Fig. 2.1 The performance of CPU and memory from 1980 to 2005

The memory hierarchy is organized into several levels. The upper level is smaller, faster, and more expensive per byte than the lower level in Fig. 2.2. We can know that the level of the hierarchy is usually a subset of one another. All data in the first one level are also found in the lower level. It divides memory into several sections, such as register, cache, main memory, hard disk, and etc. From the CPU to memory, if it is closed to the CPU, its speed is fast and its capacity is small, such as Cache and register. They are composed of SRAM. SRAM has fast speed for operation. If the level is far away the CPU, its speed is slower and its capacity is larger, such as main memory and hard disk. They are usually composed of DRAM because of its large capacity. Using the memory hierarchy can get more efficient for the memory access. Availing the special locality and temporal locality can make the memory hierarchy more efficient and reduce the difference of the latency between CPU and memory [1].

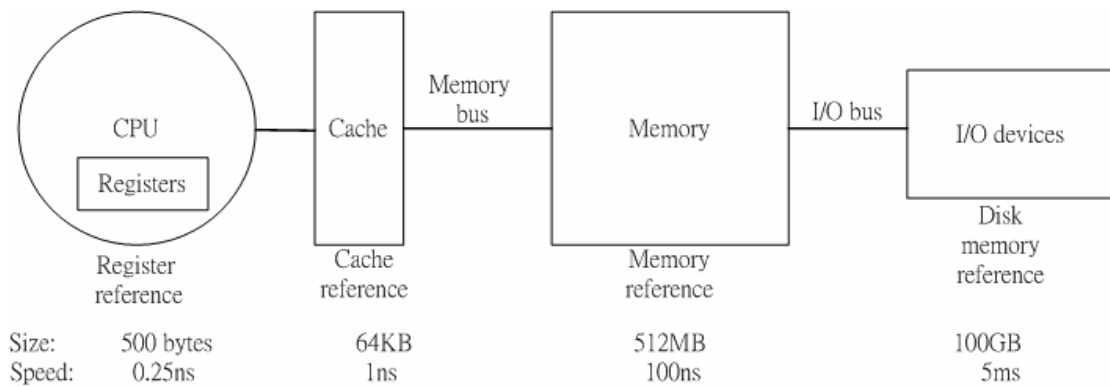


Fig. 2.2 Levels in a typical memory hierarchy in embedded, desktop, and server computer

2.1.1 Cache

Cache was the name chosen to represent the level of the memory hierarchy between the CPU and main memory, in Fig.2.2. Its function is to refer to any storage data. Cache serves as a method for providing fast reference to recently used instructions and data. When CPU finds a wanted data in the cache, it is called a Cache hit. On the contrary, if CPU does not find a wanted data in the cache, it is called Cache miss. Temporal locality means that the requested data is likely to be used again in the near future. It is useful to place the requested data in the Cache where it can be accessed quickly. A fixed-size collection of data which contains the requested data is called block. For the other data in the block, they would be needed soon for spatial locality [2].

Figure2.3 is an example for cache. The address has 32bits, and it divides three parts. One is byte offset, and it occupies two bits. Second part is index, and third part is tag. The numbers of index can tell us the capacity of cache. If there are N bits for Index, the Cache has 2^N entries which can be stored. The action is firstly to find the corresponding position of index. When the corresponding position is found out, the tag in the corresponding position would be taken out. The tag would be compared by the third part of tag. If the tags are the same and valid bit is logic one, a hit signal and the corresponding data would be sent back to CPU. The valid bit is used to indicate whether an entry contains a valid address. If they are not the same, a miss occurs. The requested data is not in the cache. The wanted data may be stored in the lower level memory. When Cache miss occurs, the operating system (OS) would take over

this condition and stall the operation for memory access. OS would find the wanted data in the lower level. If it is found in the lower level, it would be written back to the Cache and be sent to CPU.

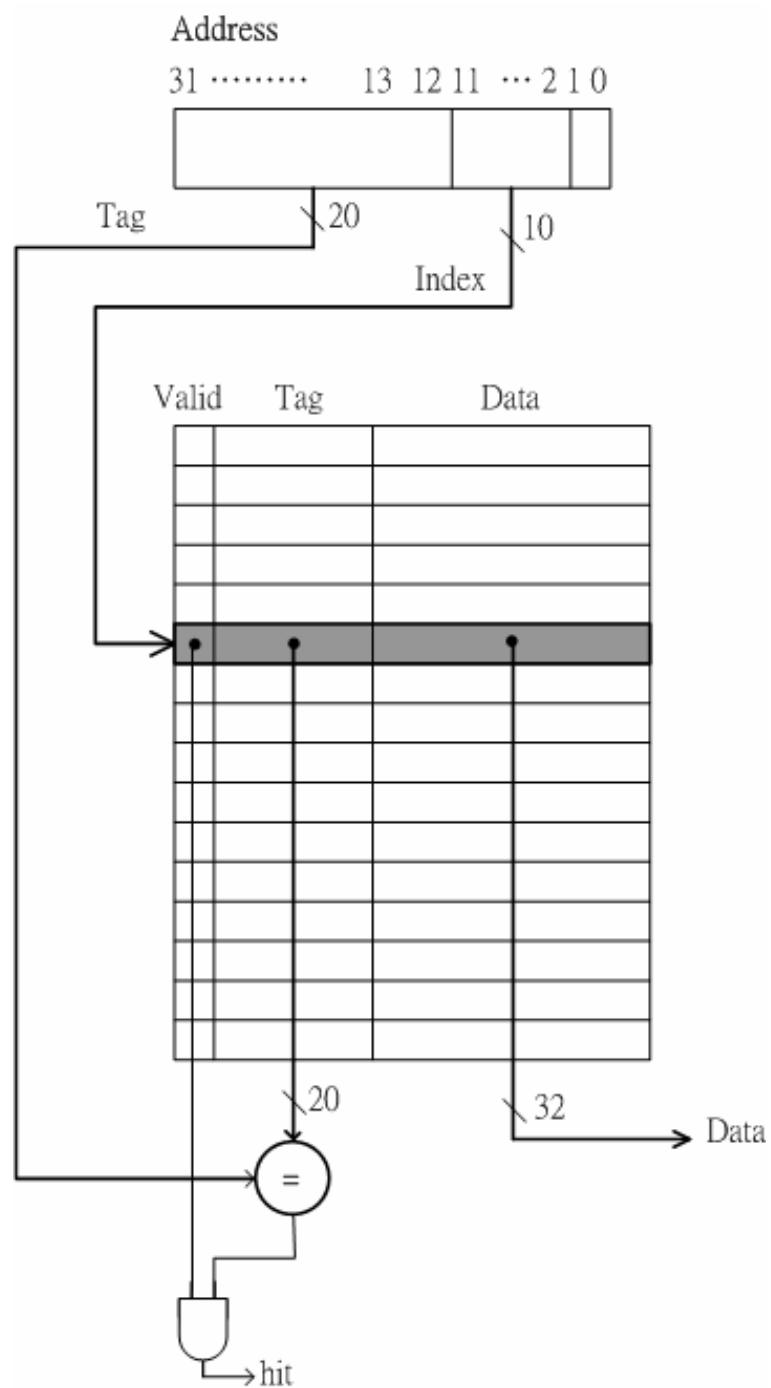


Fig. 2.3 The block diagram of a direct mapping Cache

2.1.2 Utilization of Spatial Locality

In Fig.2.3, the method of data mapping is direct mapped. Direct mapped method means that each address has one corresponding position. The index can find out the corresponding position in the cache. Another method is that each address can be placed in any position. It is called fully associative. The advantage of fully associative is to decrease miss rate, but it would increase the hit time. Because the each position must be compared to seek the wanted data for each memory access.

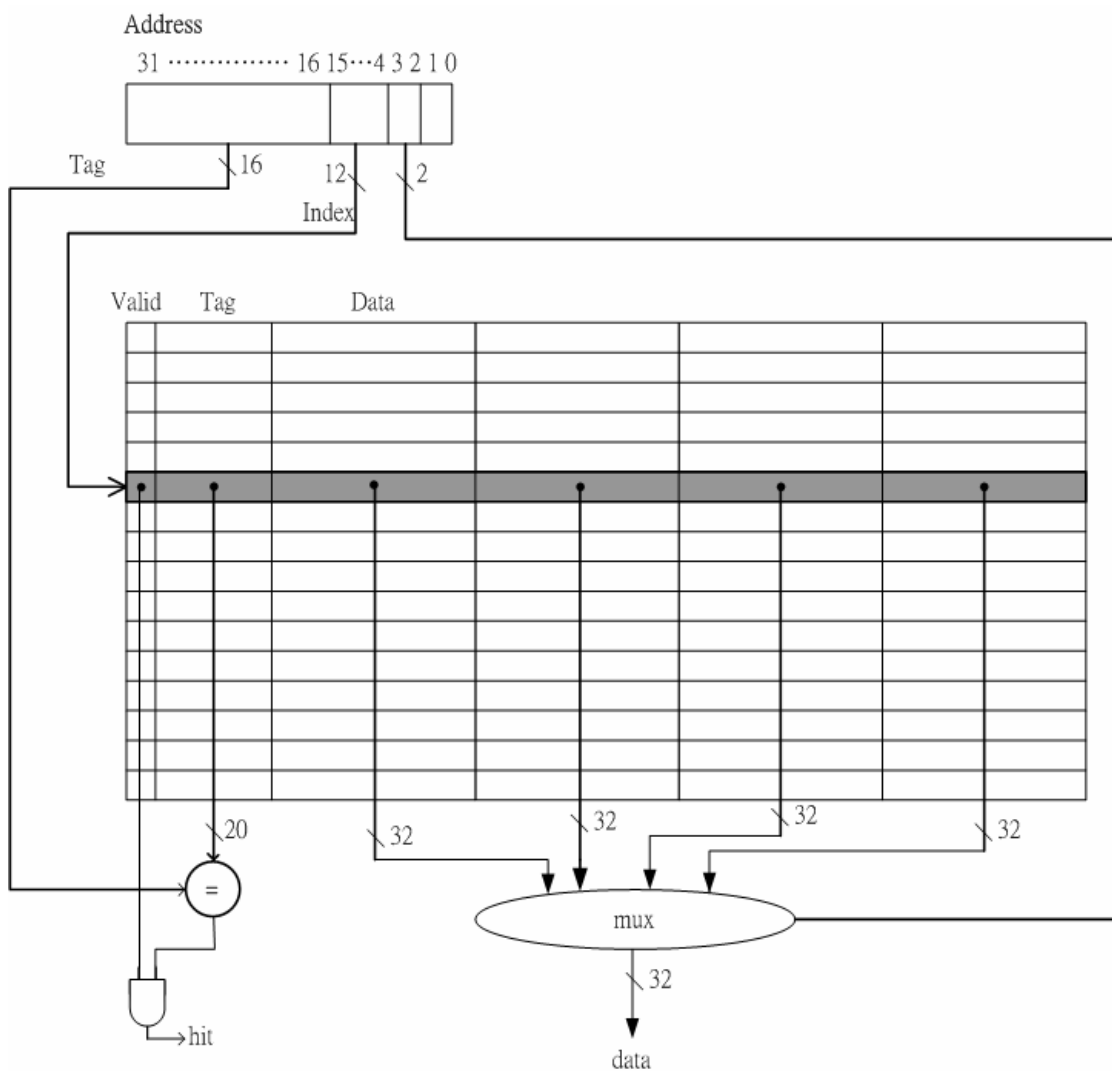


Fig. 2.4 Cache with four blocks data

We know that the cache, Fig.2.3, utilizes the temporal locality. Besides temporal locality, there is spatial locality to be utilized. In Fig.2.4, that is an example to make use of spatial locality. Each tag has the corresponding data. The corresponding block has four data. These four data are neighbors. If there is one data accessed by CPU, its neighbors of data could be accessed by the spatial locality. This kind of Cache needs the additional two bits to distinguish. Only one data is exact needed by CPU through a four to one multiplexer. This method can increase the performance actually, because the miss rate could be reduced obviously and the storage of data would be much.

2.1.3 Placement Scheme

The placement scheme is called direct mapped because there is a direct mapping from any block address in memory to a single location. That is to say, there is only one location to be replaced for each data. There are the other schemes for position replacement. Full associative is a contrary scheme. Full associative means that any one position in the Cache has the probability to be replaced. It is that each time of memory access must search all the entries of Cache to find out the wanted data. Any one location could have the chance to be placed. Its advantage is the convenient for replacement, but the search time is longer than the direct mapped placement scheme.

There is another method that is between direct mapped and full associative. It is called set associative or n-way set associative. Take two-way set associative for example, it means that there are two places to be replaced for any one data. So n-way set associative has n position to be placed. Fig 2.5 would explain these three kinds of scheme. For example, the wanted position is address four in Fig. 2.5. For direct mapped, the location is in address three. For two-way associative, its location is $(4 \bmod 3)=1$. So it is in the address zero, and it has two positions to be put. Fig. 2.6 is example for a m-way set associative Cache. For fully associative, any position is possible placed.

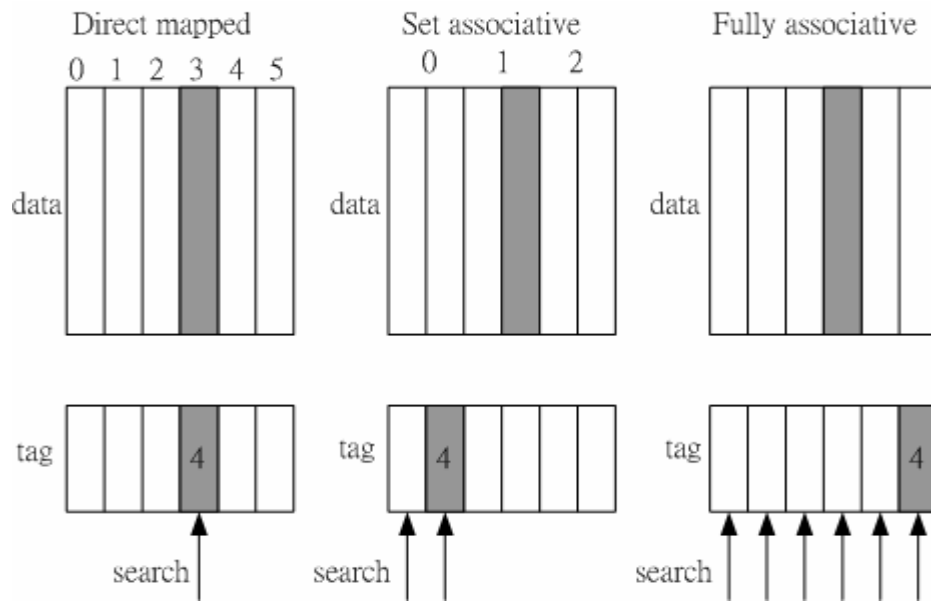


Fig. 2.5 Placement scheme for directed mapped, set associative and fully associative

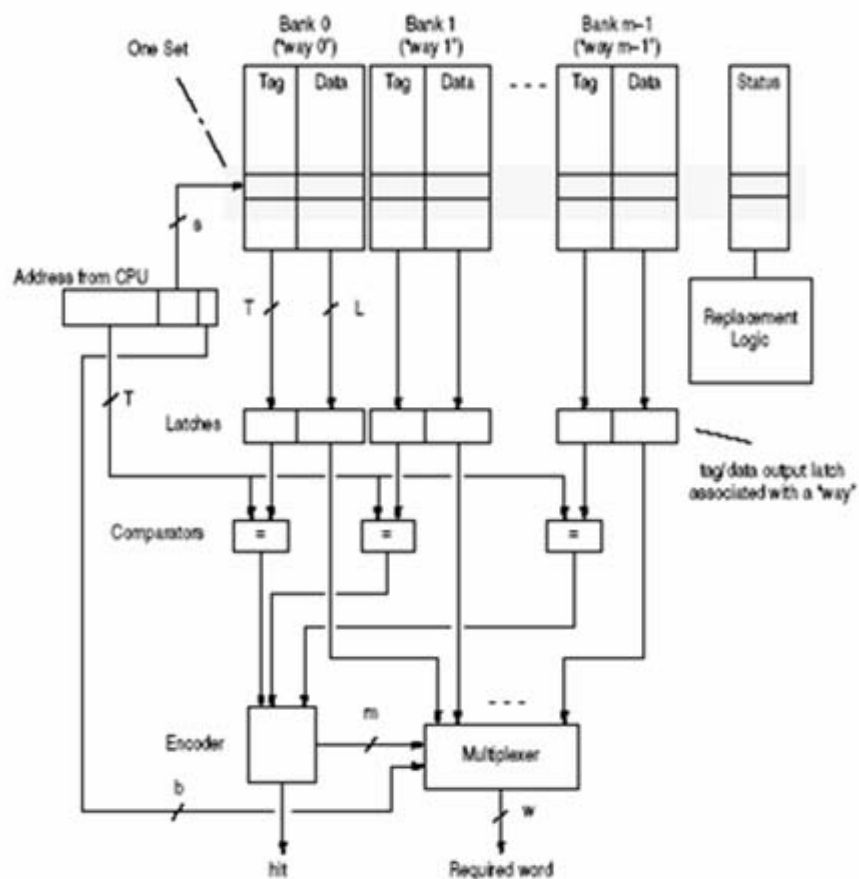


Fig. 2.6 A m-way set-associative Cache

2.2 Power Consumption about Memory Hierarchy

In Fig. 2.7, there are two different computer architectures, Strong ARM and Power PC, for power dissipation analyses. For Fig.2.3 (a) and (b), we can get the information for the percentage of the power consumption about memory section of the embedded system. The power consumption for the section of memory occupies the forty percent in Fig. 2.7 (b). For Fig. 2.7 (a), the power dissipation of memory occupies forty-two percent, ICACHE and DCACHE. From these two figures, the power dissipation for memory occupies very large percentage. So, the power consumption for embedded memory system is a key issue. How to reduce the power consumption of the Cache is important.

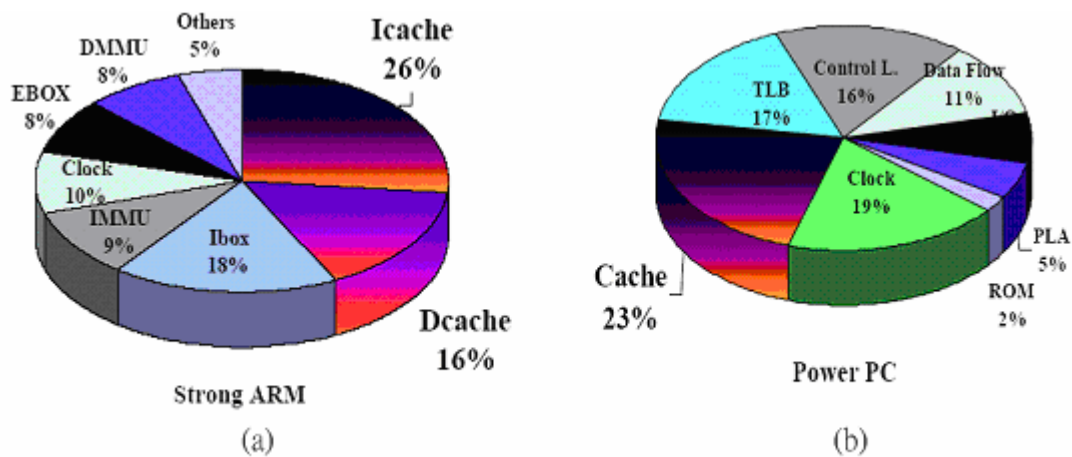


Fig. 2.7 Power consumption of CUP dissipation for (a) Strong ARM and (b) Power PC

2.2.1 Low Power Design for Cache

Memory has large portion of power consumption. How to reduce power consumption becomes an important issue. Cache is mainly composed of SRAM, and it has some peripheral circuit such as decoders, sense amplifiers, comparator, and so on. Whatever SRAM or peripheral circuit can also improve for power saving. Some methods would be introduced in the following section for lower power design.

Divided Word Line (DWL) is a technique that divides the memory into smaller blocks of cell [3]. The long word line of a conventional array is broken up into k sections, and each block is independent. In this way, its RC delay can be reduced to realize high speed and lower power dissipation. Decoder for Cache is necessary and indispensable, but it has large area and consumes much power. Many different decoders are low power, such as pre-decoder, dynamic decoder, tree-based decoder, and etc [4]-[5]. Sense amplifier is also another way. Some high performance and low power sense amplifiers are low power. [5]~[7]

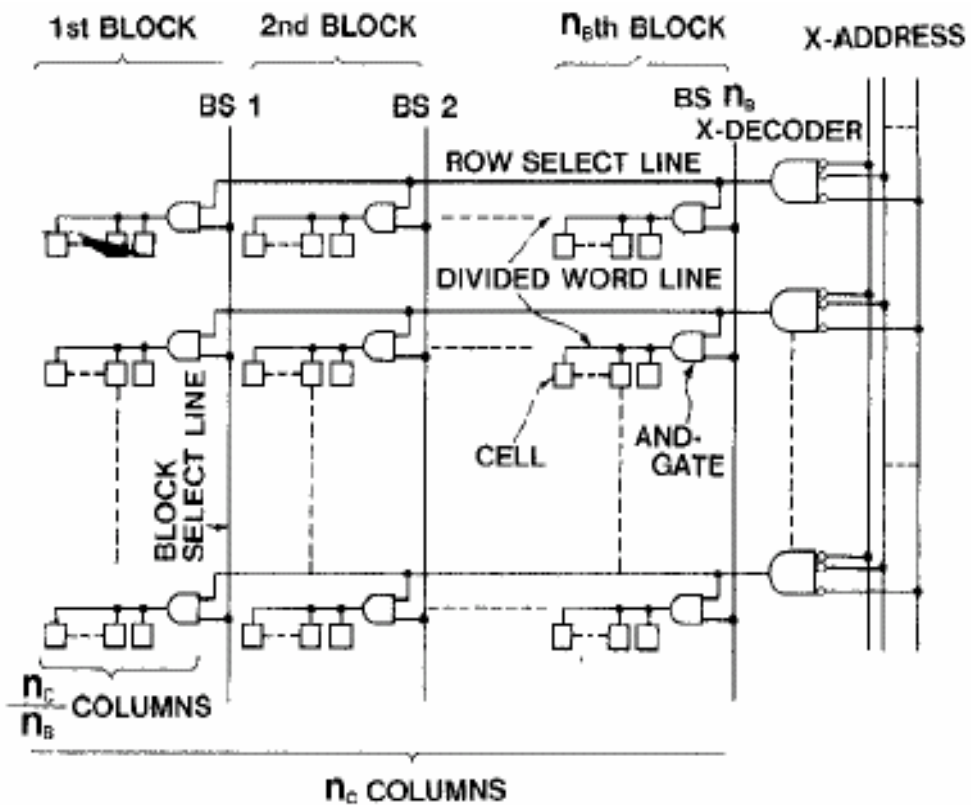


Fig. 2.8 Architecture of divided word line

The above description focuses on some circuit like decoder, sense amplifier, and word line. There are many methods to improve the architecture of Cache [8]-[9] Leakage issue is also another important problem for deep submicron design for memory devices. Some papers focus on reducing leakage power and make some improvement in leakage power consumption [10]-[11].

2.3 Virtual memory

Virtual memory is a technique that uses the main memory as a Cache for the next level memory. The motivation of using virtual memory is two reasons. First one is to let multiple programs share the memory efficiently as well as safely. Second reason is to move the small programming into the next level memory and reduce the burdens of main memory that is a limited amount of size. There are many programs that run in a machine, and these programs need much memory to be stored. The needed capacity for memory may be larger than the actual capacity for the real memory. But not all programs would be executed at the same time, in Fig. 2.9. So if the main memory only contains the programs that run at the moment, the sharing memory for many programs can have better efficient.

The concept of virtual memory is the same as cache, but the virtual memory block is called a page. When a virtual memory address is a miss during the search, it is called a page fault. With virtual memory, the CPU produces a virtual address, which must be translated to a physical address by hardware or software to access the main memory. In virtual memory, the address can be divided into two parts, Fig. 2.10, one part is virtual page number and the other part is page offset. The page offset field determines the page size.

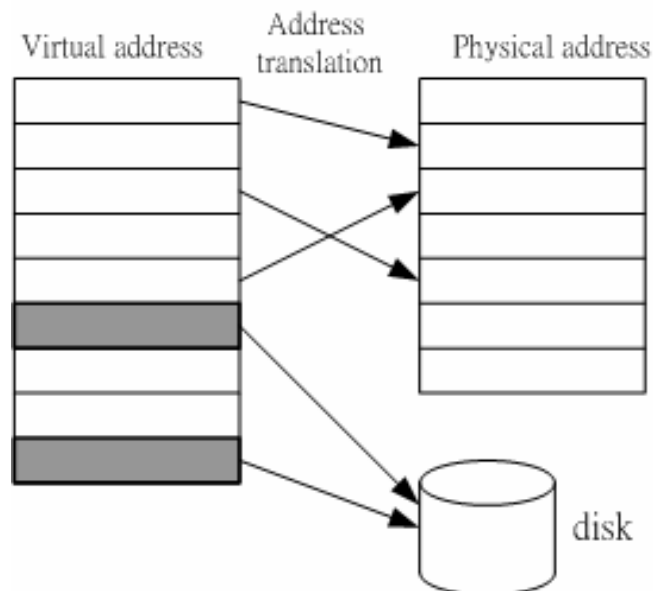


Fig. 2.9 Concept of virtual address memory

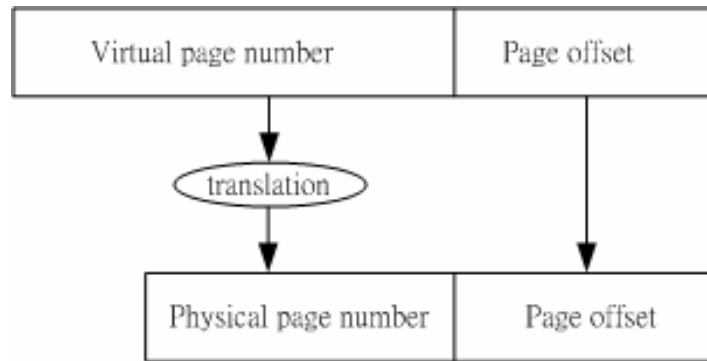


Fig. 2.10 Mapping from a virtual address to physical address

2.3.1 Page table

Page table is a structure that can locate pages by using a table to index the memory. Page table can index from the virtual address and contain the corresponding physical page number. Each program has its own page table, which can map the virtual address space into physical address in main memory. Page table is different from Cache. Page table does not need tag because page table contains a mapping for every one possible virtual page. When page fault occurs, the operating system would take over the condition to control. Operating system must find out the page in the next level of the hierarchy and send back to place the requested page in the main memory. There are some placement schemes to do, such as, least recently used (LRU) replacement scheme [2].

2.3.2 Translation Lookaside Buffer

Page tables are usually so large that getting physical address would waste much cost. One compensation method is to using the principle of locality. By keeping the address translations in a special Cache, the memory access would not need the second time access to translate the data from page table. This special address translation Cache is called translation lookaside buffer, is abbreviated TLB. A TLB entry is like a Cache. The tag holds portions of the virtual address and the data holds a physical page numbers. Each reference looks up the virtual page number in TLB. If it gets a hit, the physical page number is used to form the address. If a miss in TLB occurs, it must be checked whether it is a page fault or merely a TLB miss.

From Fig 2.11, we can know the relation of virtual memory, TLB, and Caches. When a virtual address produces by CPU, it contains virtual page number and page offset. It would be compared by TLB firstly. Virtual page number would be compared with tag of TLB. If the valid bit is on and the page number is the same as tag, then TLB is hit. It will send the corresponding physical page number out. The physical page number and page offset of virtual address will be translated to physical address tag and index. Physical address tag would be used to compare the Cache tag. If the tag is the same, then Cache is hit and sends the stored data out.

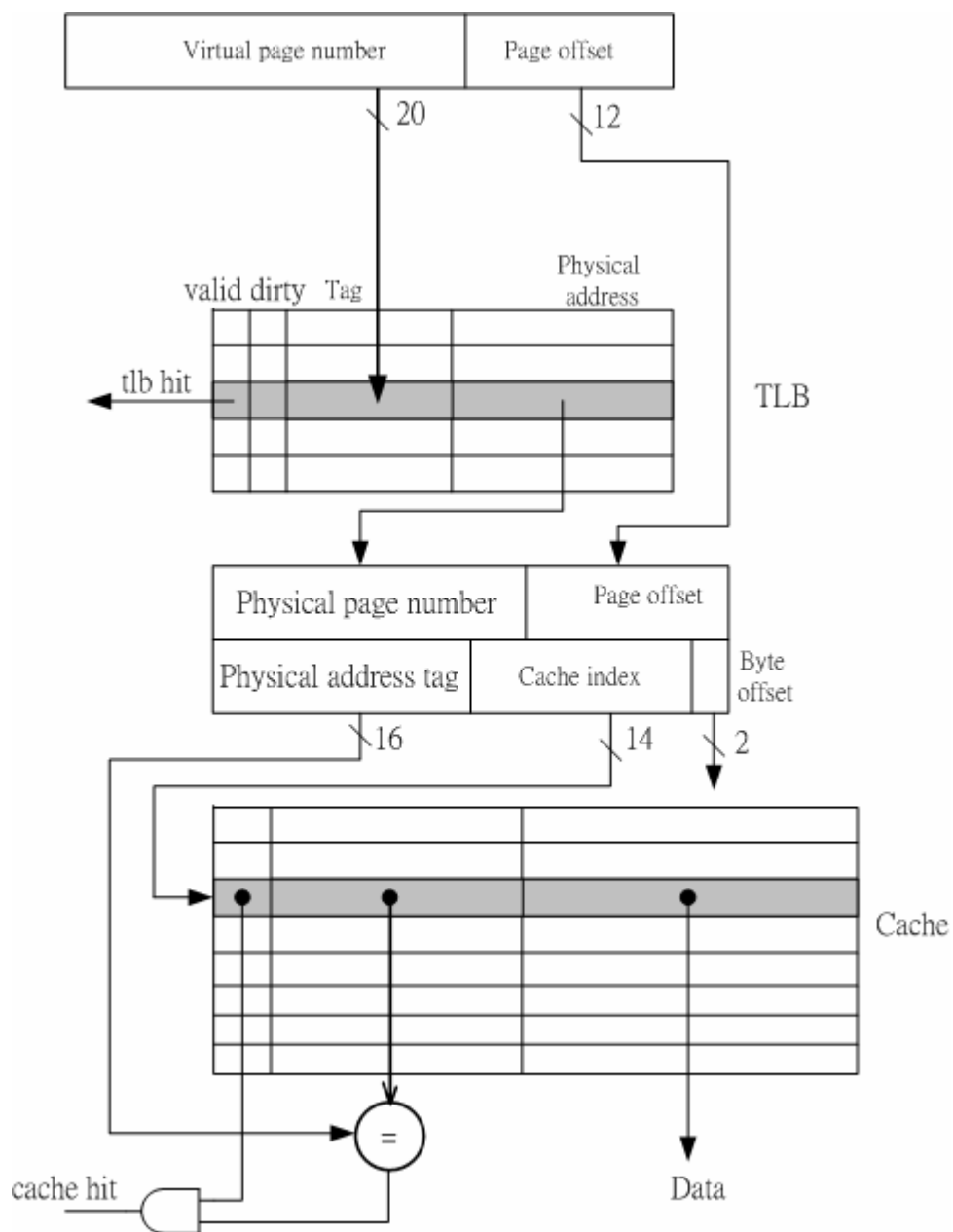
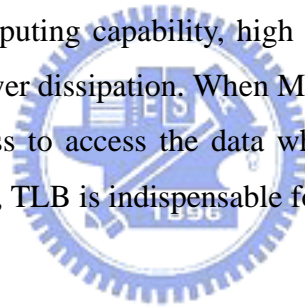


Fig. 2.11 The block diagram of integrating virtual memory, TLB, and Cache

2.3.3 Memory Management Unit

Memory management unit (MMU) is responsible for the memory access. MMU is important, especially, in the system-on-a-chip (SOC) design or embedded system design. The goal of SOC design or embedded system design is high performance and low power, so the efficient MMU is essential. For the SOC or embedded system, there is multiprocessor, such as, CPU, DSP, DMA and etc [12]-[13].

In the computer architecture, function of TLB is to provide fast address translation from the virtual address to the physical address. Any embedded processors that support the virtual memory system through the MMU would need TLB. The need of TLB is more and more important. Those processors are widely used for multimedia and communication applications which require high-speed computing capability, high memory bandwidth, effective memory hierarchy support, and low power dissipation. When MMU get the physical address from TLB, it can use this physical address to access the data which is stored in the memory (such as Cache or DRAM or DISK). So, TLB is indispensable for MMU.



2.4 Organization of TLB

TLB is mainly composed of SRAM and CAM. We can see the figure, Fig. 2.11, and have more concepts about the TLB. We know that TLB provide address translation. When there is virtual address to be send to TLB. The virtual address that are composed of virtual page number and page offset would be an index to compare all the indexes. Those indexes are stored in the CAM. If the virtual page number can find a corresponding index in the CAM, then, the corresponding physical address that is stored in memory will find out the corresponding index. So TLB needs SRAM to store the physical address. The object of CAM provides the comparison about the index addresses and the compared addresses which are stored in the CAM. The role of the CAM is like SRAM to store the addresses and must have the ability to compare their index addresses. When there is a match in one word line of CAM, CAM will send the match information to SRAM, and get the address which is stored in the

corresponding word line of SRAM. There are some other circuits, such as the row decoder, sense amplifier, precharging circuit, and etc.

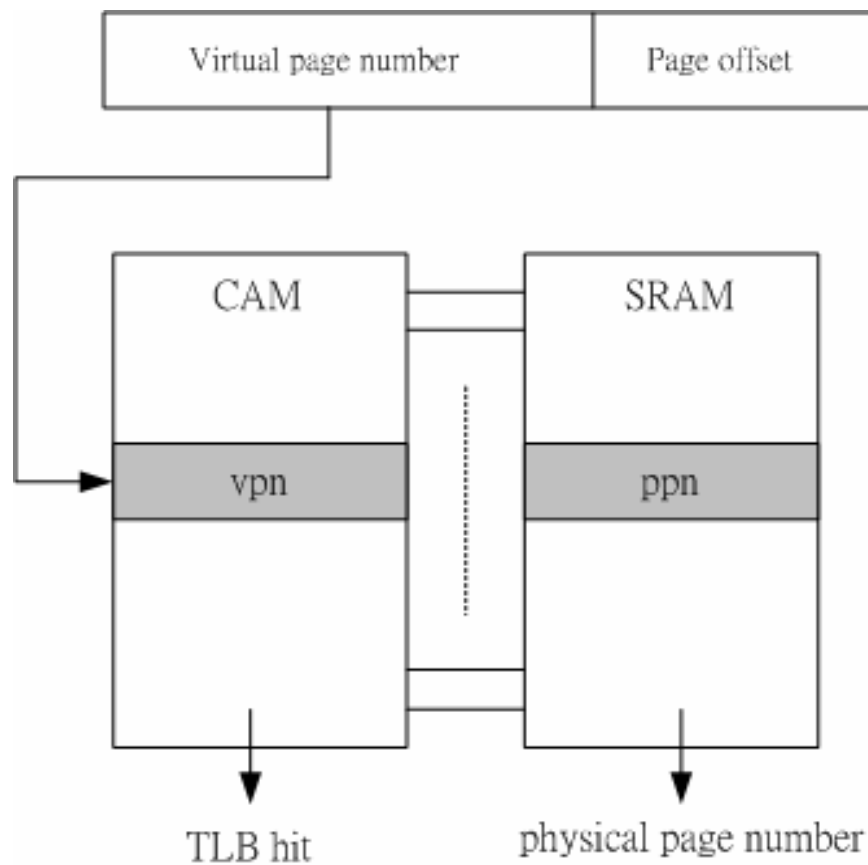


Fig. 2.12 Organization of TLB

2.4.1 Some Viewpoints to Reduce the Power Consumption

Since the TLB and Cache have mainly effects for the power consumption. Reducing the power consumption for TLB is very important issue. There are three viewpoints to achieve the goal of low power consumption. First one is the software level. Using software can replace the hardware. It can reduce the overhead of hardware and decrease power consumption. Second one is to improve the architecture level. Third one is to improve the circuit level. Improving the architecture and the circuit is a method to reduce power. The following section would introduce these three levels more detail and how to reduce the power consumption with these three levels.

2.4.2 Software Level

At the software level, the principle is very simple. If the address translation runs without going through the TLB, the original power consumption of TLB would be saved [14]-[16]. This is the ideas about the software level. If the software can directly provide the physical address of the page being referenced, then we do not need the TLB for that instruction and reference. The CAM will not be used and the power consumption of the CAM would be saved, because the software can replace the hardware. In fact, there are four kinds of possible combinations for the address. They are virtually-indexed, virtually-tagged (VI-VT), virtually-indexed, physical-tagged (VI-PT), physically-indexed, physically-tagged (PI-PI), physically-indexed, virtually-tagged (PI-VT). These four kinds of combination may be in the Cache or in TLB. PI-VT is not really in much use, so it would not be introduced in the following text.

PI-PT means that the physical address needs to be obtained before the TLB or Cache can even be indexed. There are no aliasing problems across different virtual address spaces for this kind of scheme. VI-PT can remove the TLB from critical path because the index uses the virtual address, and TLB is concurrently looked up to obtain the physical address. That is to say, the tag from the physical address is used for the comparison with the corresponding tag bits. So, TLB is not in the critical path. VI-VT means that both index and tag use virtual addresses. It implies that TLB or Cache is not required at all until the TLB miss or Cache miss. This method can look up TLB and Cache in parallel at the same time. But this strategy has aliasing problems, and the solution is to add a few most significant bits to differentiate between address spaces.

There is another method. It can operate without going through TLB. TLB is a translation from virtual address into physical address. If first level Cache (L1) don't use the physical address to access, it directly accesses virtual address. Then, it does not need the translation from virtual address to physical address. It uses the virtual address to access in the first level Cache. This way can save the power consumption of address translation in TLB. Only the letter of the first level Cache could need the physical address to compare. The TLB also needs to work. The above methods need the support of operating system.

2.4.3 Architecture Level

At the architecture level, different architecture of TLB can have different effect. The simple method is to use the smaller size of TLB. Because the larger size of memory is, the more power dissipation occurs. On the other hand, it needs longer searching time for fully associative placement scheme. As the size of memory is increase, the power consumption will increase fastly. Because the loading of the memory also increases, the power dissipation would degenerate. But if the size of TLB is smaller, its performance will degenerate. The stored entries decrease would increase the miss rate. It will make the performance worse. So, there is a tradeoff about size of TLB between power and performance. In the following section, some kinds of TLB would be described, such as multi-level TLB, banked TLB, dynamic resizing TLB, and so on.

2.4.3.1 Multi-level TLB



For instance, multi-level TLB can be viewed as a hierarchy TLB in Fig.2.13. It has two parts. First part is a smaller TLB, called filter TLB. The other part is a bigger TLB, called main TLB. These two TLB will also store the recently used data. When the memory access occurs, the first level TLB, small TLB, would be searched before going to the second level. If it is found in the smaller TLB, it is match and does not need to search the second level. If it is not found in the filter TLB, then, it would search the main TLB. It is like the memory hierarchy.

Using this mechanism, it would consume less power consumption in this condition. When there is a match for filter TLB, the main TLB would not be searched. Because the circuit and area of filter TLB is smaller, its power consumption is less than the power consumption in the main TLB. This kind of architecture is like the principle of memory hierarchy [17].

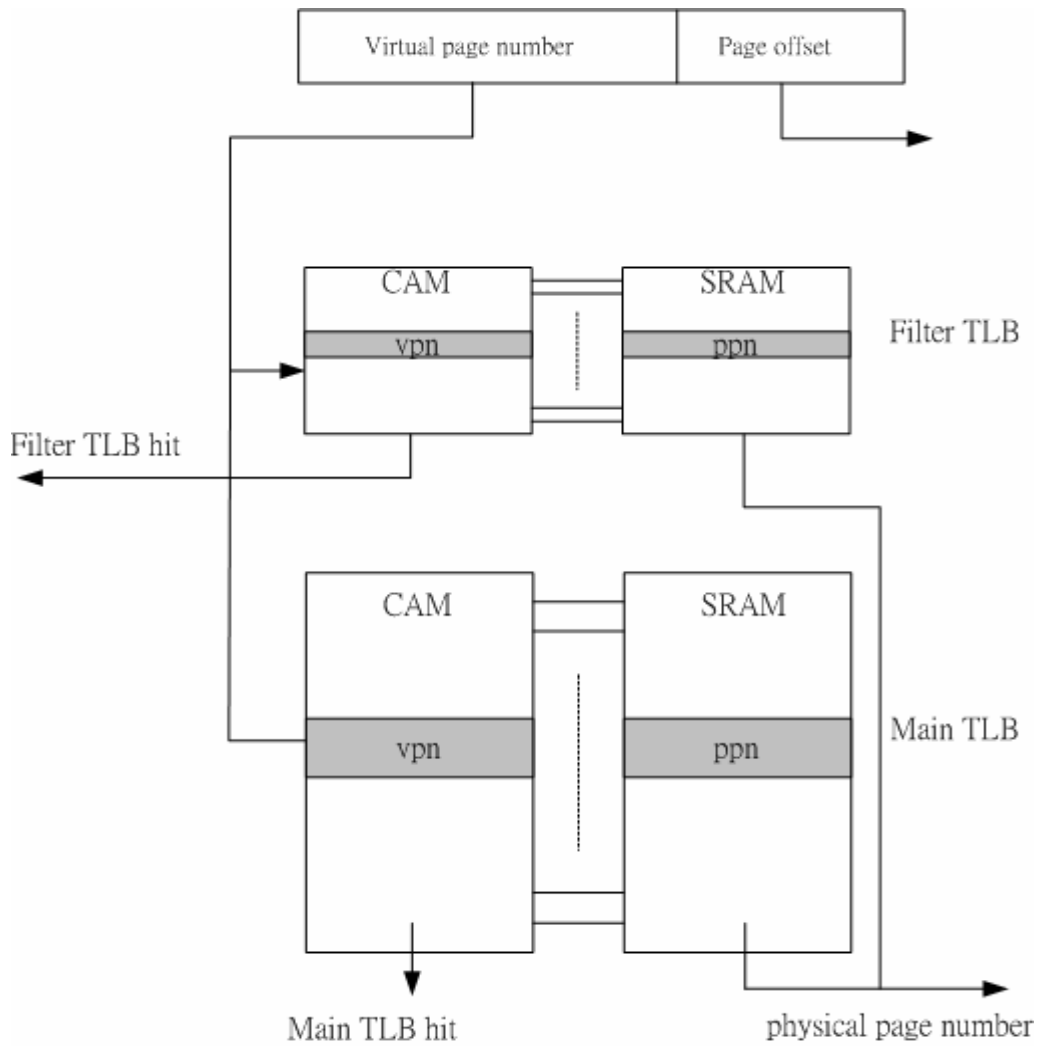


Fig 2.13 Multi-level TLB

2.4.3.2 Banked TLB

Another method splits the TLB into many banks and each bank has the same area and entries. It is called banked TLB in the Fig.2.5. The principle of the banked TLB is that each bank is smaller than original TLB. Some bits are used to index the bank. For each time comparison, it will index one bank and only search the chosen bank. It would reduce the power consumption because each bank is smaller than the original TLB. The power consumption only costs in the bank. There are also many other different methods, such as re-configurable TLB. The re-configurable TLB can dynamic change its size of TLB to reduce the power consumption and miss rate. This method needs another hardware unit to detect change of the miss rate and decides the adaptive size of TLB [18]-[19].

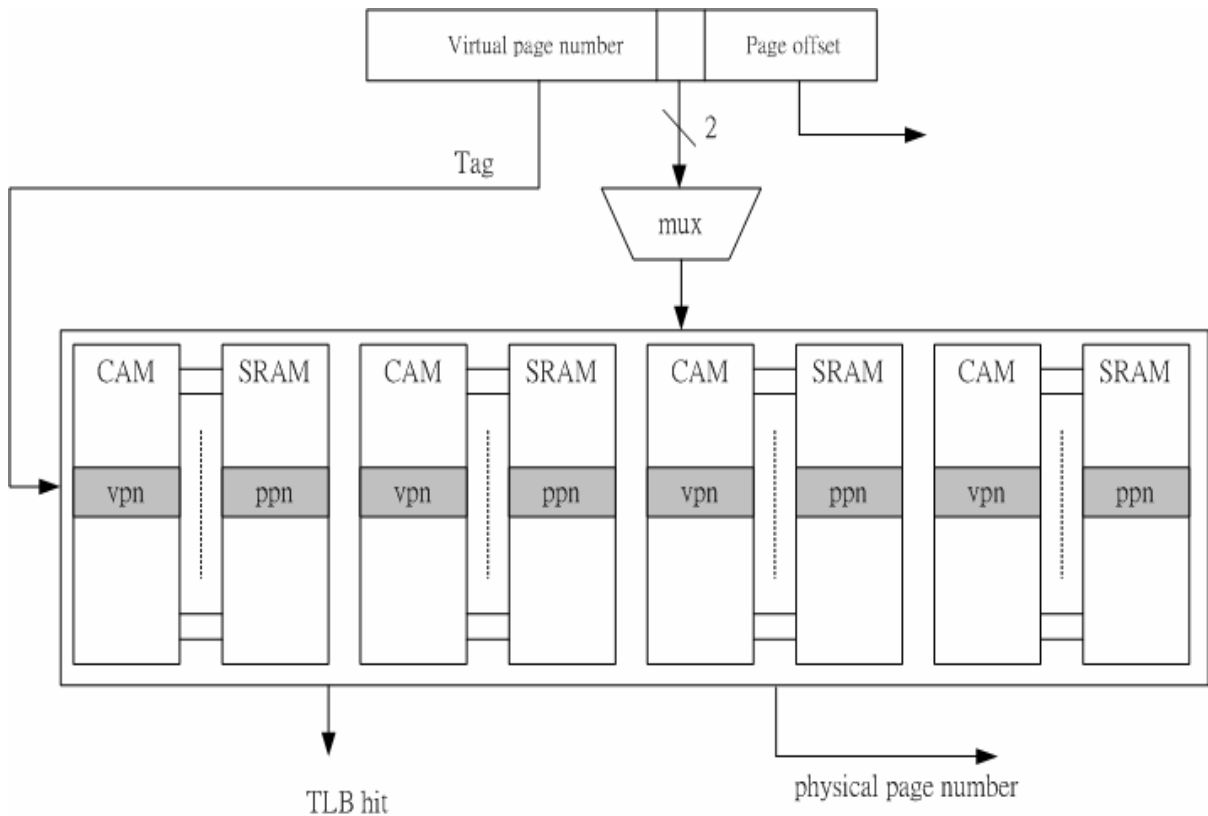


Fig. 2.14 Banked TLB

2.4.3.3 Dynamic Resizing TLB

Dynamic resizing TLB means that its size can adjust dynamic by some algorithm [20]-[21]. Its principle is to use the counter to count the times that is not accessed or the times of miss for all entries of TLB. If the counter exceeds some designated number, the hardware will judge to reduce the size or increase the size. This method only adds some additional hardware, counters and some logic circuit to adjust the size of TLB, and it can use the optimal size of TLB to reduce the redundant power consumption from the extra entries of TLB.

2.4.4 Circuit Level

At the circuit level, designing the CAM is the most important thing. The CAM is very wasteful in the power consumption of the TLB. Because only one match would not discharge, the others would discharge. So this action would increase more power consumption.

Changing the CAM can achieve very significant effect about reducing power consumption and this part would be described in the following text.

On the other hand, it is also to reduce the power consumption by changing the circuit of the SRAM. With gated vdd is an example, the following figure, Fig.2.6. With gated vdd is mainly to reduce the leakage current power consumption. The leakage effect is more and more obvious in the deep submicron technology. So using the method of power gating can reduce more power consumption with deep submicron technology. The extra transistor which is controlled by signal control is turned on in the used sections and is turned off in the unused sections. In this way, it can effectively turn off the supply voltage and eliminate the leakage for the unused sections.

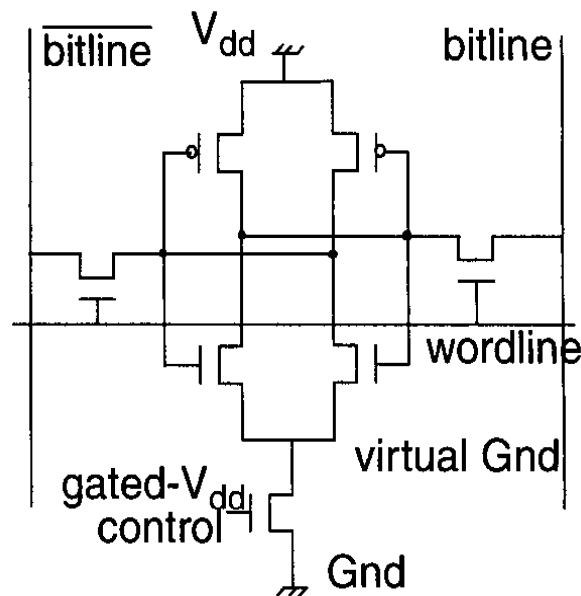


Fig. 2.15 Gated-Vdd SRAM

2.5 Content Addressable Memory

The contents addressable memory, or called associative memory (CAM) is a special type of memory device. Its characteristic is that it can not only store the data like the memory but also can compare the stored data with input data. It is composed of SRAM and XOR logic gate. So it has three modes of operation: read, write, and comparison. Fig. 2.17 is the

architecture of CAM. From the Fig. 2.17, CAM array is its main part obviously.

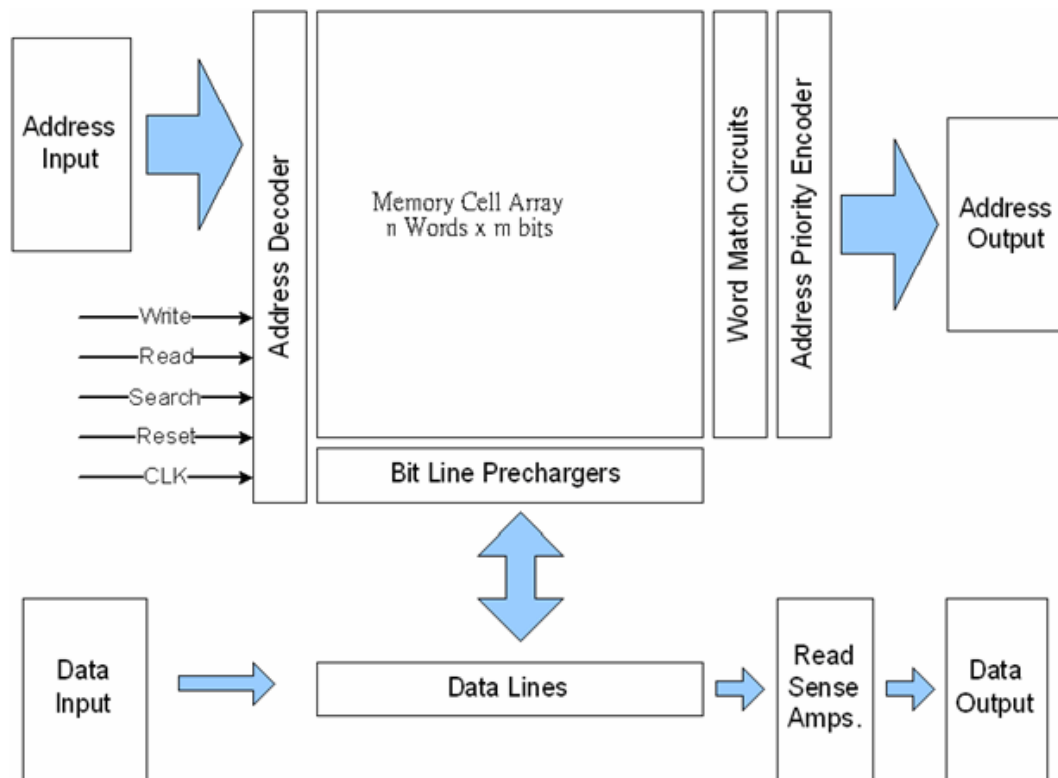


Fig. 2.16 Architecture of CAM

The other peripheral circuits are such as I/O buffers, address decoder, sense amplifier, equalizer and etc. Although CAM is composed of SRAM, its peripheral circuit is some different from the peripheral circuit of SRAM. Address decoder is an example. For CAM, it only needs row decoder and does not need column decoder. This is because CAM stores data, read data, and compared in parallel for a unit of one word line. This characteristic of parallel comparison is very special and important. So, it is often used in many places to be compared and be stored, such as, router, CAM-tag Cache, TLB, and any application which is needed to be compared. Recently, in high speed networks such as gigabit Ethernet and asynchronous transfer mode (ATM) switch are also used it. Its advantage of the comparison is very fast and efficient. But the tradeoff of CAM is to waste much power consumption because of its parallel comparison. In the following text would introduce some kinds of conventional CAM, and describe their components and functions [22]-[23].

2.5.1 Conventional 9T CAM

There are many kinds of CAM. From the CAM, they can be classified by its transistor numbers. The conventional CAM is composed by nine transistors and it is called 9T CAM. Its component is in the Fig. 2.17. CAM is mainly composed by two parts. First part is the SRAM part and the other is the comparison part, XOR logic circuit. The SRAM is composed by six transistors, and there are two inverters which are cross coupled each other. The other two nMOS are composed of pass transistors. The operation of SRAM has two kinds of operation. One is reading operation and the other one is writing operation. When one of the word line W is high, it means that it is in the writing operation. There is one word line that would be chosen. The writing data is put in the wanted bit line and its complementary data value is put in the bit line bar. The writing operation could be completed and all the bits of the chosen word line will be written. In the reading operation, the bit line would be put half of v_{dd} . Choosing one word line makes W high. Through the comparison, if the stored value is bigger than half of v_{dd} , it is the v_{dd} . If the stored value is smaller than half of v_{dd} , it is the zero.

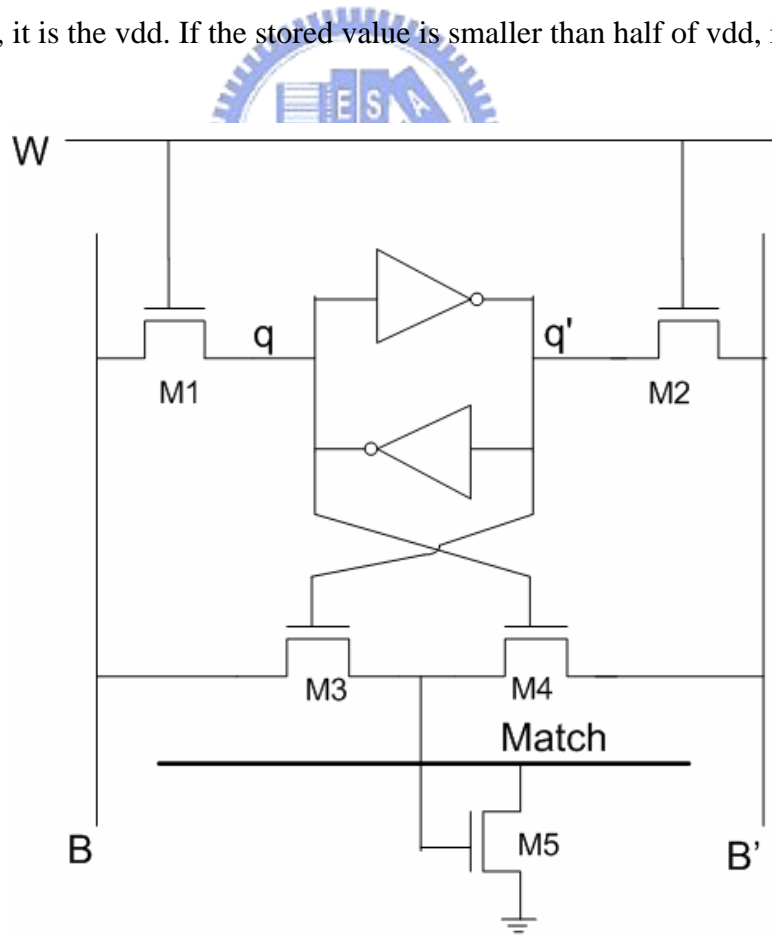


Fig. 2.17 CAM cell

Its comparator part uses XOR. When the two inputs, q' and B (bit line value), are the same values, for example both logic one, the gate of M5 can get logic one. Then, M5 is open and discharges the match line to logic low. From the Fig. 2.17, M5 is open and discharged when the gate of M5 is logic high. M5 will turn off when the gate of M5 is logic zero. That is to say, B and q' are both logic one or B' and q are both logic one. They could pass the logic one to the gate of M5 and turn off the M5. This condition is mismatch. On the contrary, if the q and B are the same values, its gate will get the logic zero. In that way, it can not turn on M5. M5 would remain in logic high and do not discharge the match line. Whatever the match line is logic high or logic low, it would be precharged periodically to maintain its value to vdd, and prepare to do next comparison in the next evaluation phase.

The above description is the operation of the conventional 9T CAM. From the diagram, it is one bit CAM. Take a 8x8 CAM memory array for example. There are sixty-four one bit CAM. Each word line has eight bits and these eight bits use the same match line. There are all eight match lines. Their match lines are precharged periodically. When the comparison occurs, all the eight word lines will be compared. So, if each one of the eight bits is mismatch, the match line would be discharging to zero. Only these eight bits are all match, the match line would stay at vdd. Only one word line will be match in general, and the other match lines will be discharged for mismatch. So, the conventional 9T CAM consumes much power. Besides, its memory density is low. The XOR is not easy to layout tightly. Therefore, the capacitance of bit line and capacitance of bit line bar are heavy. They cause much power consumption. Although the performance of conventional 9T CAM is fast, it also has many disadvantages. Especially, the power dissipation is very big and this is its main problem.

2.5.2 Nor-type 9T CAM

There are two types of CAM. First one is called NOR-type CAM. It connects the comparison nMOS, M5 in Fig. 2.17, in parallel for one word line. The drain terminals of M5 are connected to the same point, match line, and the source terminals are connected to the ground. When any one transistor turns on, the match line would be discharged to ground. In Fig. 2.18, the circuit is a word line of NOR-type 9T CAM. This kind of CAM has an

advantage, its efficiency. Its speed of comparison is very fast. Because any one bit is mismatch, the match line would be pulled down to ground. So it is suited to use in the application for high speed comparison. However, its cost for NOR-type CAM is its power consumption. Fig. 2.18 is an example for n bits CAM of a word line.

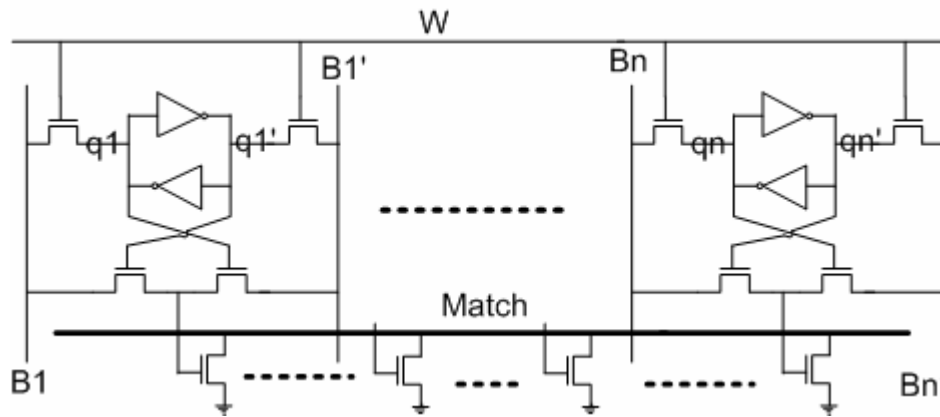


Fig. 2.18 NOR type 9T CAM

We assume that there are n bits for this word line CAM. There are 2^n cases for this n bits word line. For each CAM comparison, there is only one case to be match. That is to say, there are 2^{n-1} cases mismatch. For each time comparison in m word lines, there is only one word line match at most. Sometimes, the comparison may be mismatch for all word lines. Besides, the mismatch word lines would discharge the match lines. NOR-type CAM will discharge more for each time comparison. Therefore, we can know that conventional NOR-type 9T CAM would consume much power dissipation in the comparison operation. Even though the conventional NOR-type 9T CAM has very good advantage about its performance, its worst power consumption is also a big problem for many circuit designs.

2.5.3 Nand-type 9T CAM

The second type 9T CAM is called NAND-type 9T CAM. The discharging transistors for NOR-type CAM are connected in parallel. The NAND-type CAM is opposite. Its discharging transistors are connected in series, in Fig.2.19. Because the discharging transistors are in series, only one case would let all discharging transistors turn on. When all the bits are match

for the compared data, all transistors would be turned on to pull down the match line to ground.

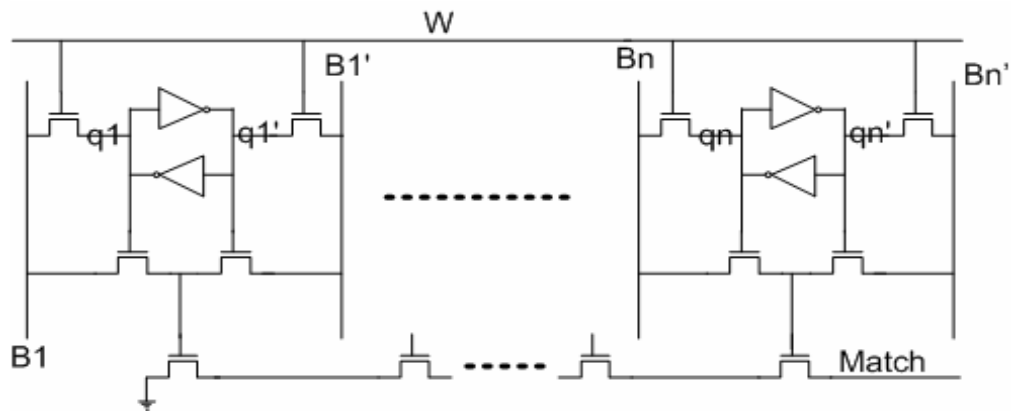


Fig. 2.19 NAND type 9T CAM

When any one bit is mismatch, it will not turn on all discharging path for NAND-type CAM. For each comparison, there is one word line match or no word lines match. So, one word line could be discharged at most. Most of match line will be keep the voltage. For this reason, the power dissipation for NAND-type CAM is much less than the NOR-type. But NAND-type has a disadvantage. That is its performance. When the discharging path produces, the discharging current would flow through all the discharging nMOS. This would cause the comparison speed to slow down.

There is a tradeoff between power consumption and performance for using NOR-type CAM or NAND-type CAM. In order to having better performance and lower power, there are some methods of compromise to be used. The serial nMOS increases more, then, the access time would be longer. This would make the performance down. There is one method to increase the performance of the NAND-type CAM. It adds repeaters in the middle of the precharging path. The serial nMOS would be separated into some parts by repeaters. If the match line behind of the repeater is match, the repeater would stay at logic high to connect to the next serial nMOS. If the match line that is behind of the repeater is mismatch, the repeater would pull the next serial nMOS down to vdd. Repeater can decrease the access time and enhance the signal of match line, but its penalty is its area. Repeater would occupy some area for CAM.

2.5.4 Ternary CAM

For the CAM circuit design, the ternary CAM (TCAM) performs a more powerful data search function [24]-[25]. Different from binary CAM which has two states: one (1) and zero (0) state. It needs sixteen transistors overall for one bit ternary CAM. The ternary CAM (TCAM) cell has an additional state: don't care (X) state. Alike binary CAM, TCAM would be classified into two kinds: NOR-type TCAM and AND-type TCAM. Both of them would be introduced in following sections.

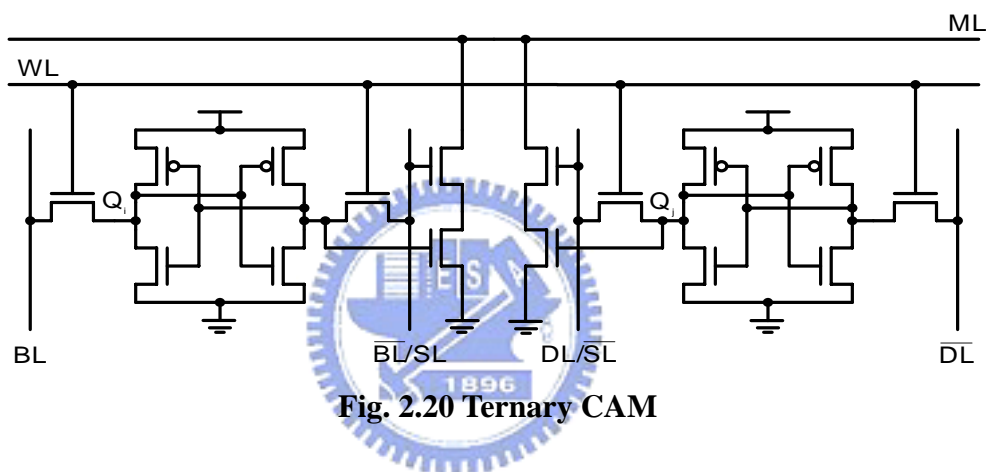


Fig. 2.20 Ternary CAM

2.6 Summary

This chapter makes an overview about memory hierarchy. The memory hierarchy is indispensable for high performance design. The component of memory hierarchy has Cache, TLB and so on. Cache and TLB also contain CAM. CAM can provide high speed comparison to support the high performance. Accept for performance, the power consumption becomes another important issue for embedded system. CAM is very power wasting because it compares all the word line in parallel for high performance. How to reduce power for embedded system is also a significant issue.

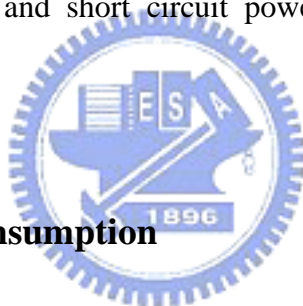
Chapter 3

Low power Pre-comparison CAM

In this chapter, some kinds of power consumption would be discussed. Conventional NOR type 10T CAM would be introduced. The concept of pre-comparison would be depicted in the following section and the pre-comparison CAM is low power.

3.1 Power Consumption

In recent years, as the CMOS technology scales down, the power consumption becomes a very important issue in CMOS design. There are three kinds of power consumption, dynamic power, static power and short circuit power. They would be introduced in the following section.



3.1.1 Dynamic Power Consumption

Dynamic power consumption is due to charging and discharging capacitances. The formula of dynamic power consumption is as follow [4]. It occurs only during transients, when the gate is switching.

$$P_{\text{dynamic}} = C_L \cdot V_{\text{DD}}^2 \cdot f_{\text{CLK}} \quad (3.1)$$

From the formula, dynamic power is dominated by three components. They are C_L , V_{DD} , and f_{clk} . C_L is the loading capacitance, or called switching capacitance. V_{DD} is supply voltage, and f_{clk} is clock frequency. Clock frequency is also called switching frequency. In the low power design, scaling down V_{DD} can have a quadratic reduction of power. As the supply voltage scales down, it has some bottlenecks. First of all, the performance of the circuit will have degradation that is due to lowering the supply voltage. Secondly, the static noise margin would limit the scaling down V_{DD} in the memory design. In general low power design, V_{DD} would not be change for these two reasons. But there are still some methods that are used to

reduce the VDD.

How to reduce C_L and f_{clk} is also important in low power design. Separating large switching capacitance into several banks can reduce the switching capacitor. Lowering switching capacitance is also to improve the performance of the circuit. The loading capacitance is usually due to transistor capacitance in the combinational circuit. To keep the minimal size is necessary for low power circuit design in the transistor level. Decreasing unnecessary times of switching can reduce power consumption, too. Reducing the switching activity can only be accomplished at the logic and architecture levels.

3.1.2 Static Power Consumption

The static power consumption is present when no switching occurs and is caused by static conductive paths between the supply paths or leakage currents. It is always present for the circuit in standby mode. The expression is as follows [4]:


$$P_{\text{static}} = I_{\text{static}} \cdot V_{\text{DD}} \quad (3.2)$$

I_{static} is the current that flows between the supply paths in the standby mode. Ideally, the static current would be zero. Actually, there are many kinds of leakage current in the MOS. There are subthreshold, band-to-band tunneling, gate tunneling, pn junction reverse bias, DIBL, GIDL, and punchthrough leakage, in Fig. 3.1.

In generally, the leakage current is very small and can be ignored. However, the junction leakage currents are caused by thermally generated carriers. It would increase with the increasing temperature exponentially. Subthreshold leakage would be smaller with a higher threshold voltage (V_t). It would become important power consumption because the CMOS technology scaling down causes the threshold voltage to shrink. As the process steps into the deep submicron, gate leakage will dominate the leakage current and maybe exceed the dynamic power consumption.

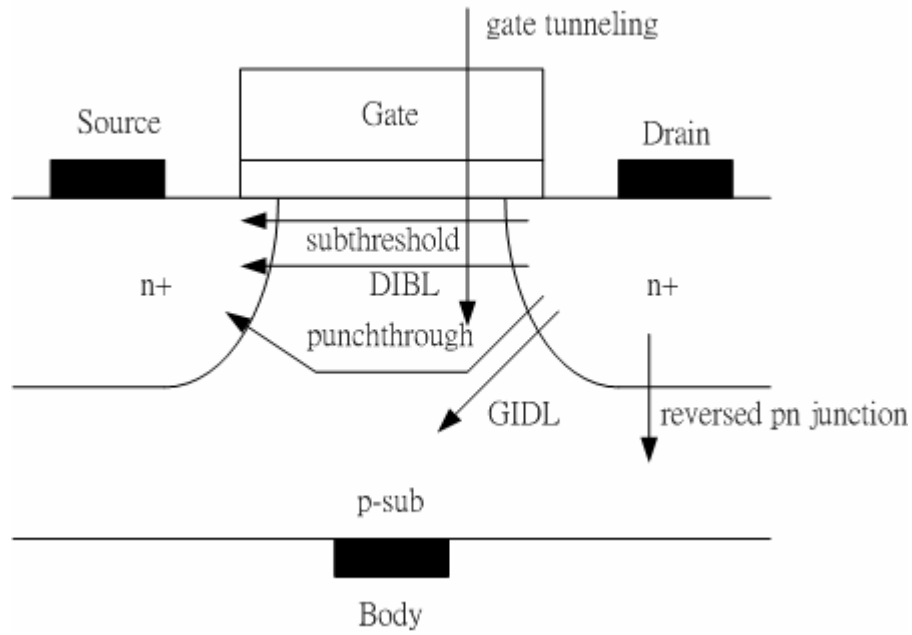


Fig. 3.1 Leakage current sources in a MOSFET device

3.1.3 Short Circuit Power Consumption

In the ideal logic circuit, there are zero rise time and fall time for the input wave form. Actually, the input wave form has a finite slope. The slope would cause a direct current path between VDD and GND for a short period of time during switching. In the short time, the NMOS and PMOS are conducting simultaneously. From the following Fig. 3.2, the current would occur in its rising and falling response. The energy and the average power consumption can be computed per switching period as follows [4]:

$$E_{sc} = V_{DD} \left(\frac{I_{peak} t_{sc}}{2} \right) + V_{DD} \left(\frac{I_{peak} t_{sc}}{2} \right) = t_{sc} \cdot V_{DD} \cdot I_{peak} \quad (3.3)$$

$$P_{short-circuit} = t_{sc} \cdot I_{peak} \cdot V_{DD} \cdot f = C_{sc} \cdot V_{DD}^2 \cdot f \quad (3.4)$$

The short-circuit power consumption is proportional to the switching activity. [1] The t_{sc} is the time that NMOS and PMOS are conducting. I_{peak} is determined by the saturation current of the devices. The peak current is also a strong function of the ratio between input and output slopes. There are two cases. First case is that the load capacitance is very large. The output would response very slower than the input signal. In this way, the input signal would move

through the transient region before the output signal starts to change. Because the output almost does not change, the device would not deliver any current. Another case is that the load capacitance is very small. The response time of output would be substantially smaller than the response time of input. In this situation, it is the maximal short-circuit current and is equal to the saturation current of the PMOS. So, making the output response time larger than the input response time can minimize the short-circuit power consumption.

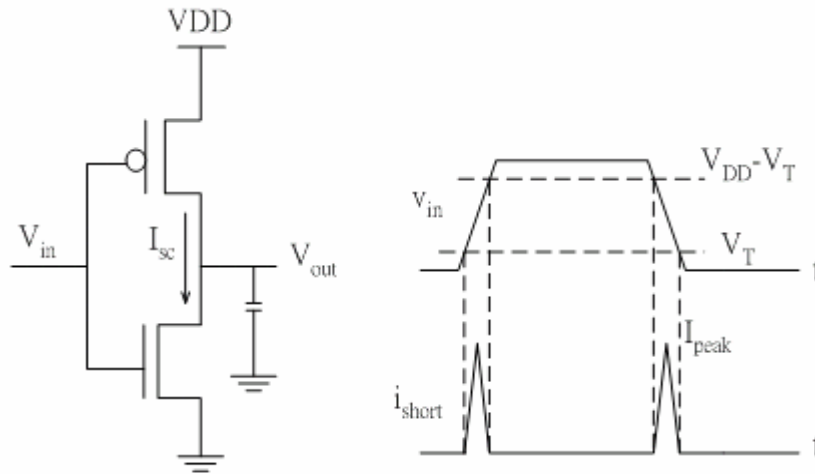


Fig. 3.2 Short-circuit currents during transients



3.1.4 Total Power Consumption

The total power consumption can be expressed as the sum of dynamic, static, and short circuit power consumption. The formula is as follows:

$$\begin{aligned}
 P_{\text{tot}} &= P_{\text{dynamic}} + P_{\text{static}} + P_{\text{short-circuit}} \\
 &= C_L \cdot V_{DD}^2 \cdot f_{\text{clk}} + V_{DD} \cdot I_{\text{leak}} + V_{DD} \cdot I_{\text{peak}} \cdot t_s
 \end{aligned}
 \tag{3.5}$$

In typical CMOS circuits, the capacitive dissipation is obviously the dominant factor. The short circuit current and leakage current are small. The leakage current would be an important issue in the deep submicron design. So, the dynamic power dissipation is the dominant component for the power consumption in the circuit design at present.

3.2 Conventional NOR Type 10T CAM

Chapter 2 has introduced two kinds of conventional 9T CAM. But 9T CAM has some disadvantages for NOR-type and NAND-type. First of all, the memory density for 9T CAM is very low. Its layout is not easy to be tight. Secondly, the bit line loading and bit line bar loading are very heavy for CAM array. As the voltage switching for the bit line and bit line bar is frequent, the power consumption would increase because of the large capacitance.

In order to reduce these two kinds of problems, 10T CAM is used. The Fig. 3.3 is the circuit of one bit conventional 10T CAM. It is also composed of SRAM and XOR. The writing and reading operation are the same as the conventional 9T CAM that has been introduced. M3, M4, M5, and M6 compose the part of XOR. There is some difference between 9T and 10T CAM. 10T CAM has separate bit lines and search lines. Bit line loading increases as the word lines increase. The large loading would make more power consumption and performance down. On the other hand, the bit line must read, write, and compare. The voltage switching for bit lines are very frequent for 9T CAM. It will increase the power dissipation. So 10T CAM separates the bit lines into bit lines and search lines. The bit lines are responsible for reading and writing for SRAM. The search lines are responsible for providing the search data. Search line data will compare with the data values which are stored in SRAM in comparison mode. This way, it can reduce the loading capacitance and reduce the voltage of switching frequency to decrease power consumption.

From Fig. 3.3, M3 and M4 are in series, and M5 and M6 is in series. If the compared data value and the stored value are different value, M3 and T4 would turn on at the same time or M5 and M6 would turn on. It would cause the match line to be discharged to logic low that is mismatch. If the compared data value and the stored value are the same value, M3 and M4 or M5 and M6 would only one transistor turn on for each one of these two teams. When the match occurs, the match line would not be discharged, and it can remain logic high. Fig. 3.4 is one word line conventional 10T CAM.

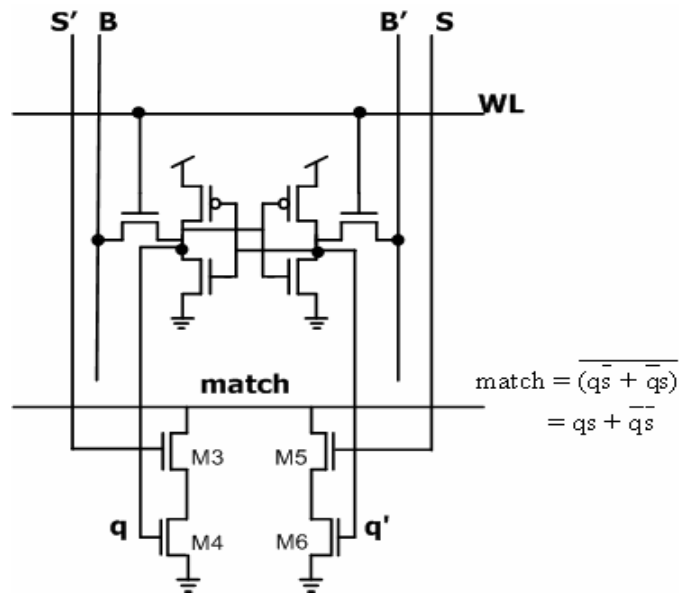


Fig. 3.3 Conventional 10T CAM

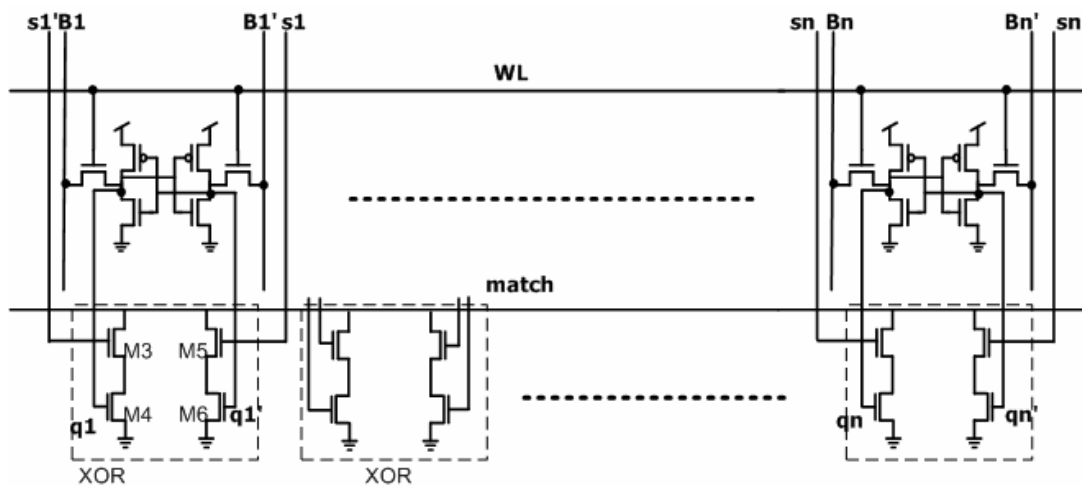


Fig. 3.4 One word line of n bits 10T CAM in parallel

3.2.1 Precharging Circuit

Fig. 3.5 is the precharging circuit. It is mainly composed of a pMOS that is controlled by signal “pulse”. The signal “pulse” will make the match line precharge periodically. There are two operations for pre-charging circuit. One is match line precharging to let floating node, a, be charged to VDD, and another one is match line evaluation when clock goes high. So there is an additional nMOS to connect the precharging point and match line to deal with evaluating. These two MOS, nMOS and pMOS, are controlled by the same signal clock, “pulse” in the

Fig. 3.5. When clock is low, pMOS turns on to precharge the node. When clock goes high, pMOS turns off and nMOS turns on. The precharging point connects the match line circuit to do evaluation at this time.

The right part of the Fig. 3.5 is a keeper. The precharging node connects the input of inverter, and the output of inverter connects the gate of a pMOS to be a feedback path. When the precharging node is precharging, the voltage of the precharging node will go to vdd. Node b will be pulled to logic low. When node b goes to logic low, the pMOS would be turned on to enhance the precharging action. So, the keeper can provide a positive feedback loop to enhance the precharging action and the judge of comparison for evaluation. The final comparison result would be through the keeper and be connected an inverter to send the result to output.

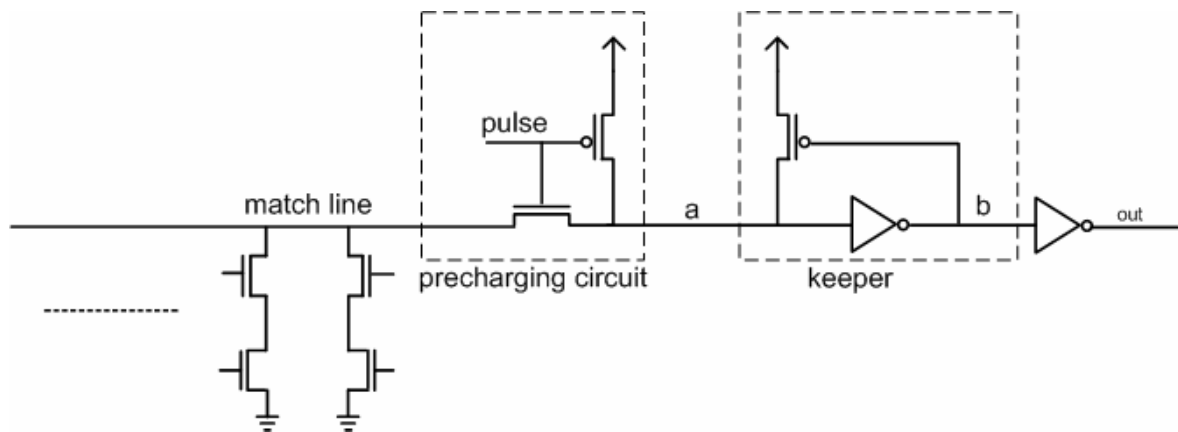


Fig. 3.5 Precharging circuit and keeper for CAM

3.3 Low Power Design for CAM

We know that the dynamic power dissipation is the main factor to cause the most power consumption. If we can reduce the supply voltage, switching capacitance, and switching frequency, we can have low power design circuit. In the following text, we would introduce some low power circuit designs for CAM. There are many methods to be used. Some people change the circuit of CAM to lower the power dissipation. Some people change the match line and precharging circuit to lower the power consumption. There are also many methods to modify the architecture of CAM array to reduce power.

3.3.1 Active Low CAM

The structure is in the Fig. 3.6. The Active Low CAM has an active low match line as opposed to the active high match line of the conventional 9T CAM. The match line would be discharged to GND periodically. In the evaluation, the discharge path would be closed and evaluated. If a read miss occurs, the match line would pre-charge to $(VDD-V_t)$. If a read hit occurs, the match line would stay low state. The full swing of conventional CAM is VDD. So, its function is contrary to the conventional 9T CAM [26].

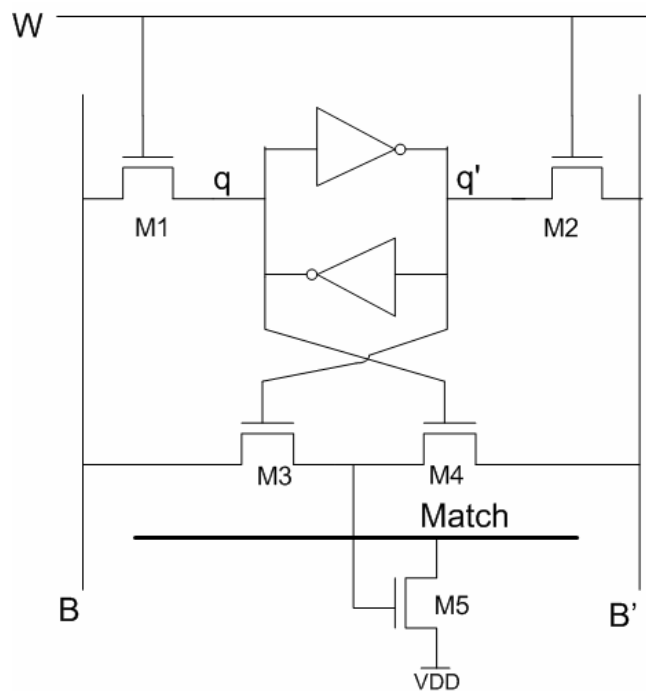


Fig. 3.6 Structure of Active Low CAM

There are two differences between this kind of CAM and the conventional CAM. Firstly, this Active Low CAM requires a pre-discharge circuit instead of a pre-charge circuit, because it functions the contrary to the conventional 9T CAM. So, it needs to pre-discharge the match line periodically. Secondly, the transistor M5 is connected to VDD instead of ground. M5 is nMOS, and nMOS is used to be pass transistor. NMOS is not good to pass the logic one, but it is good to pass the logic zero. It only delivers full voltage swing when the voltage is GND. The voltage is VDD, and it only can deliver $(VDD-V_t)$, and it is not a full voltage swing. So

the Active Low CAM cell consumes less power consumption than the conventional 9T CAM cell. It can reduce the power of quadratic V_t . So, Active Low CAM cell can save more power than the conventional 9T CAM cell.

3.3.2 Control-gatedd CAM

The circuit of Control-gated CAM cell is as the following Fig. 3.7. Its difference is that it has another nMOS logic gate, transistor M6. M6 is added to the circuit in the discharging path and is connected by M5 in series [26], [28]. It is used a control signal to control the discharging path when it is really needed on. It must remain off during the precharging region and remain on during the evaluating region.

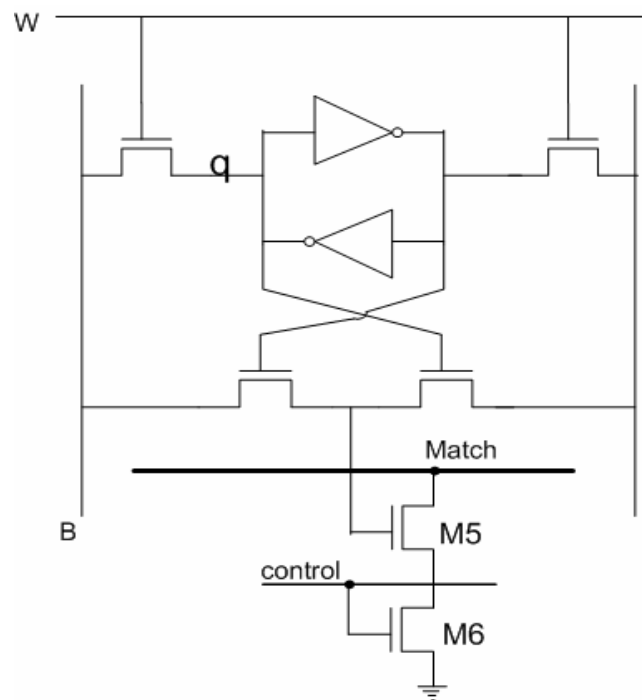


Fig. 3.7 Structure of control-gated CAM

Because the bit line and bit line bar of conventional 9T CAM must set to logic zero during the precharging region. During the evaluation region, half the lines must be set to logic one as the search lines and search bar lines are inverse. There are N (numbers of bit in a word) lines switching at one time, so the additional transistor T6 is added to prevent a path to ground

during the precharging region for the Control-gated CAM. The Control-gated CAM can allow the match line to be precharged without having to zero both the bit lines and bit lines bar. The bit lines only change due to bit changes in the value of consecutive accesses. So the Control-gated CAM cell can reduce the power consumption by reducing the times of switching of bit lines.

3.3.3 Selective Precharging CAM

For the fully associative CAM, all of the match lines are always precharged and discharged for each time memory access. The switching of the precharging and discharging would consume significant power consumption. Let CAM divide several blocks. If each time of memory access can only search one small block of the CAM, it can reduce some portion of power consumption because of smaller portion comparison. The Selective Precharge CAM is used by this concept, in Fig. 3.8. It uses a small subset of inputs that is used to be compared firstly. This action will check which one block is the comparison target. If a partial match is obtained in a given row, it is match. The best case is that none match lines are precharged. Its worst case is that all the match lines to be precharged with a slight delay penalty. On the other hand, the circuit of the Selective Precharge CAM is the same as the conventional CAM [27], [29].

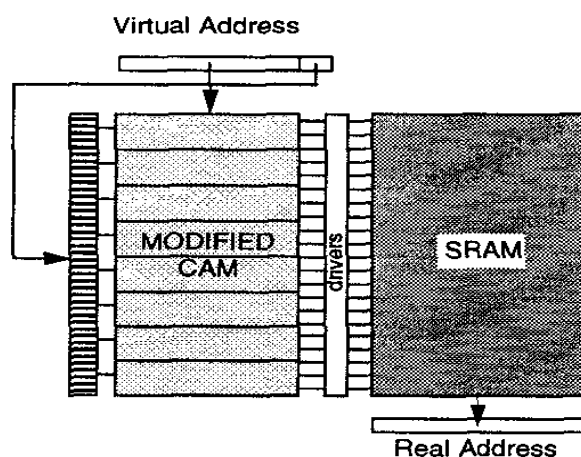


Fig. 3.8 Selective Precharging CAM

3.3.4 Precomputation-based CAM

The concept of Precomputation-based CAM is to use an additional hardware, parameter extractor, from the input data, in Fig. 3.9. Through the parameter extractor, it can decide the precharging circuit to precharge or not to reduce the power. The parameter extractor is like as the counter to count how many ones are. Each data which is stored in CAM is also pre-counted to store in P_0 to P_{m-1} . When any access data comes, it would be counted by parameter extractor to compare others that are stored. If the numbers of one stored in CAM is the same as the output of parameter extractor, the word line would turn on to evaluate and compare, in Fig. 3.10. Its precharging circuit is some different from the conventional precharging circuit. Traditional circuit adopts the dynamic circuit to improve the overall system performance and hardware cost. This architecture uses static circuit design, in Fig. 3.10. The author lists some drawbacks about the traditional dynamic circuit. First of all, the dynamic circuit needs an extra precharging time for each data searching operation. Secondly, dynamic circuit has charging sharing and noise problems. So, this architecture adopts a static circuit design for the precharging circuit [30]-[32].

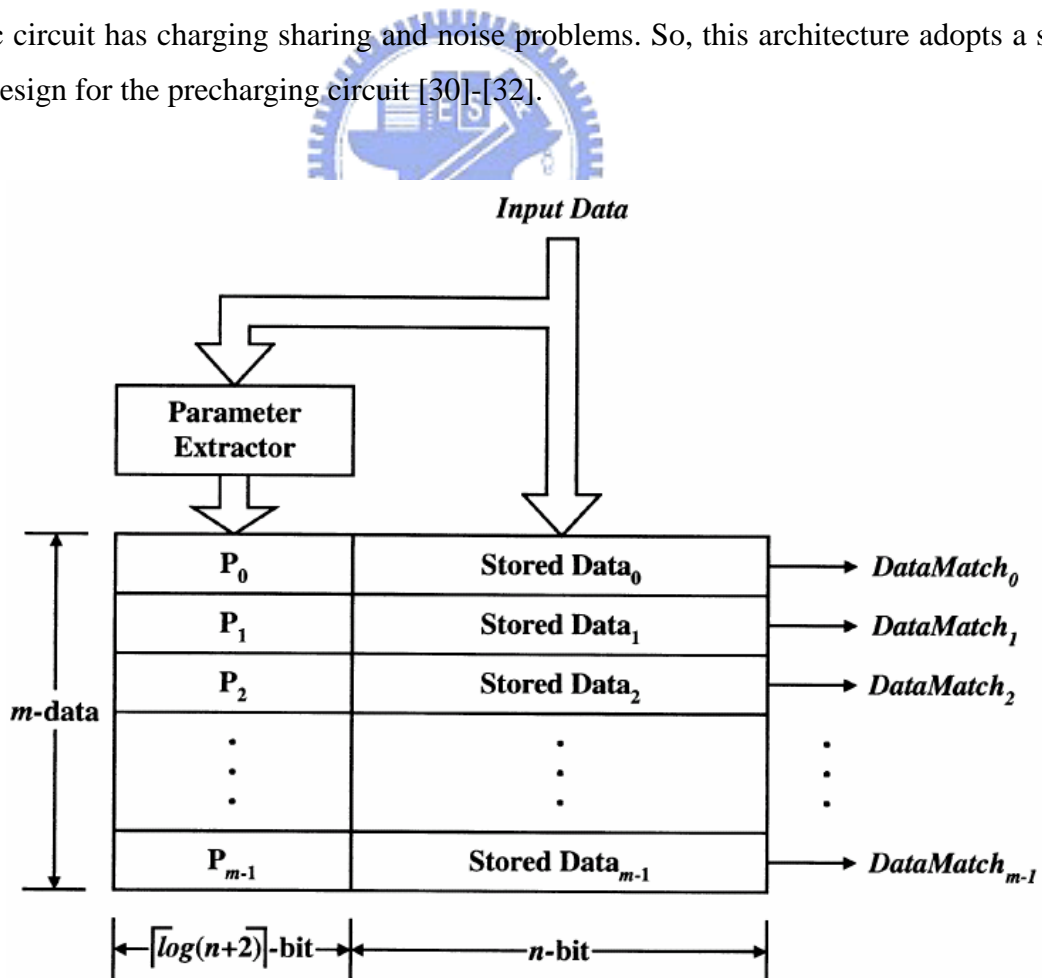


Fig.3.9 The memory organization of the precomputation-based CAM

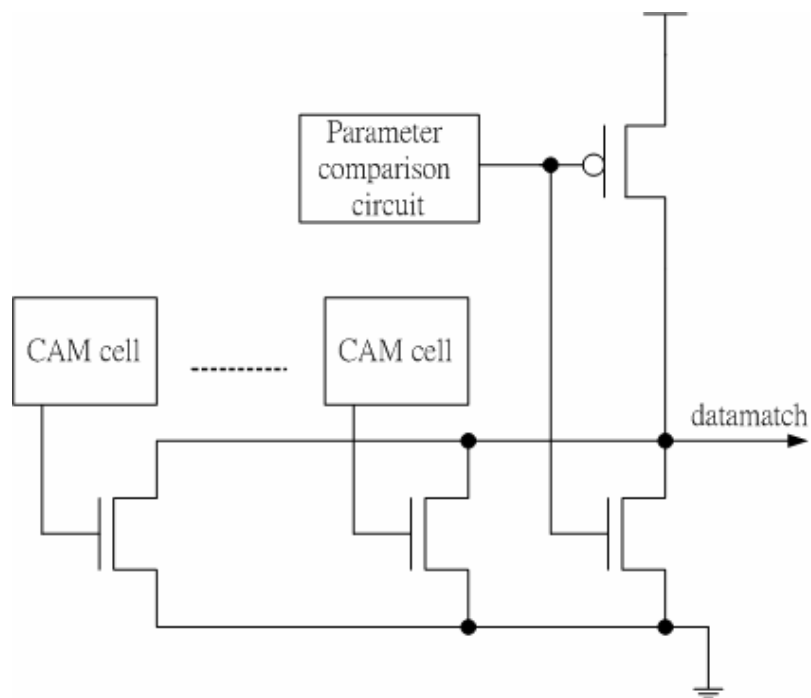


Fig. 3.10 Static pseudo-nMOS CAM circuit



3.3.5 Serial-parallel Comparison CAM

Serial-parallel comparison for CAM is a hybrid mode CAM, in Fig. 3.11. It combines the NOR-type 10t CAM as well as NAND-type 10t CAM, and it needs twice comparison. First part is called serial CAM which is composed of NAND-type CAM. Second part is called parallel CAM which is composed of NOR-type CAM. This architecture executes the comparison by comparing the serial part firstly. If the serial part is match, then the parallel part would be compared. The logic high for serial part would make the virtual ground discharge to ground. After virtual ground goes to ground, the parallel CAM would start to compare the remaining bits. When the other bits are match, then the match line is match. If the serial part has any one bit mismatch, the virtual ground of parallel would be still maintained and no discharging path produces in the parallel CAM. This result means that it is mismatch. The serial part CAM is composed of four bits, and two bits consists of a serial group in Fig. 3.11. Why this architecture takes four bits to be compared first? The authors have run the simulation about the different of bits. They find out that four bits is the optimal decision to find out the majority mismatch [33]-[34].

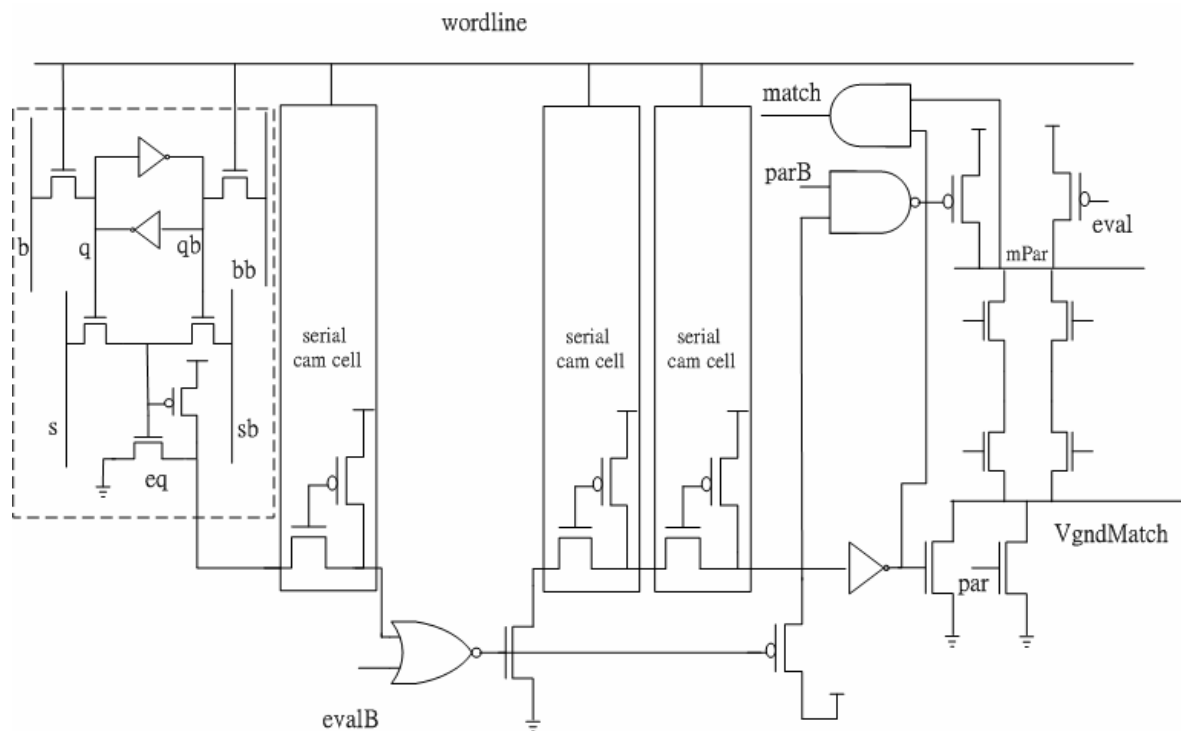


Fig. 3.11 Serial-parallel comparison CAM

3.3.6 Divided Matching Line Circuit for CAM

Fig. 3.12 is another kind of serial-parallel comparison CAM. Usually, the match line capacitance is very heavy. The match lines will precharge and discharge very often, so its dynamic power is big. We know that reducing the loading capacitance can reduce the dynamic power. This paper uses this concept, and it divides the match line into two comparison parts. It has twice comparisons, one is serial comparison and the other is parallel comparison. This kind of architecture can eliminate some drawbacks for traditional circuit, such as charge sharing, leakage current, lower noise margin, and heavy loading capacitance. Only the first comparison process is match, the precharging node will be pulled low. The logic low signal will be through an inverter to let second comparison process go to logic high. It will start to do evaluation. If the second comparison process is also to remain high after evaluation, then the comparison is match, otherwise it is mismatch [35].

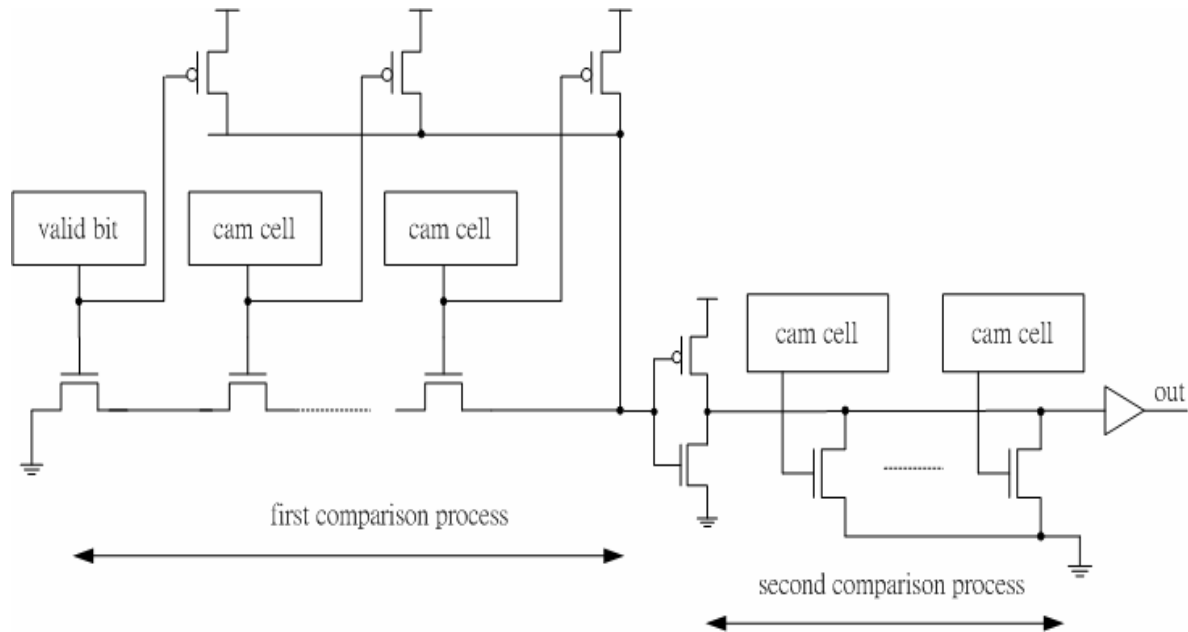


Fig. 3.12 Divided matching line circuit for CAM

3.4 Low Power Pre-comparison CAM

A lot of dynamic energy consumption is because of repeatedly precharging and discharging for all match lines. So we want to have power saving and change the architecture from the match line circuit. My idea is very simple. The principle is similar to Pre-computation-based (PB) CAM and serial-parallel comparison CAM that are mentioned in above section. Our idea is to utilize the additional logic circuit to do comparison firstly. The additional logic circuit would pre-compare some parts of all bits, before all bits that are stored in CAM cell are compared. If the small bits are mismatch, the match line will not be match. So, the remaining bits of CAM cell will not be compared, if the small bits are not match.

The pre-comparison circuit will be added between precharging circuit and match line circuit. If a mismatch occurs in the pre-comparison part, the pre-comparison circuit disconnects the path that is a passageway to connect the precharging circuit and match line circuit. When the small bits are match for the pre-comparison circuit, the all bits of CAM cell would be compared in the following step. Only the other bits are all match, it is just a match. If the other bits have any one bit mismatch, the match line would be discharging to logic low and be a mismatch. When a mismatch occurs on the pre-comparison circuit, only small bits are compared. The power dissipation is reduced because smaller bits are compare. The

Pre-computation-based CAM cell needs an additional parameter extractor to do the precomputation. It needs many full-adder (FA) to compute the value for extractor parameter. In the pre-comparison CAM, it also needs an additional logic circuit, but the additional logic circuit is less logic gates. The low power pre-comparison circuit is just to cut off the discharging path, so the circuit is simple and smaller.

3.4.1 Traditional Precharging Circuit

Fig. 3.13 is the match line circuit of NOR type 10T CAM. Match line is connected a pMOS that is controlled by pulse signal to precharge the match line periodically. Another nMOS is responsible for evaluating the match line periodically. The right side is keeper which provides a feedback loop to enhance the function of precharge. When pulse signal goes high to evaluate the match line, any one bit is mismatch that will cause the match line to discharge from VDD to GND. In parallel comparison for CAM, most of match lines are mismatch. So NOR type CAM cell is power wasting.

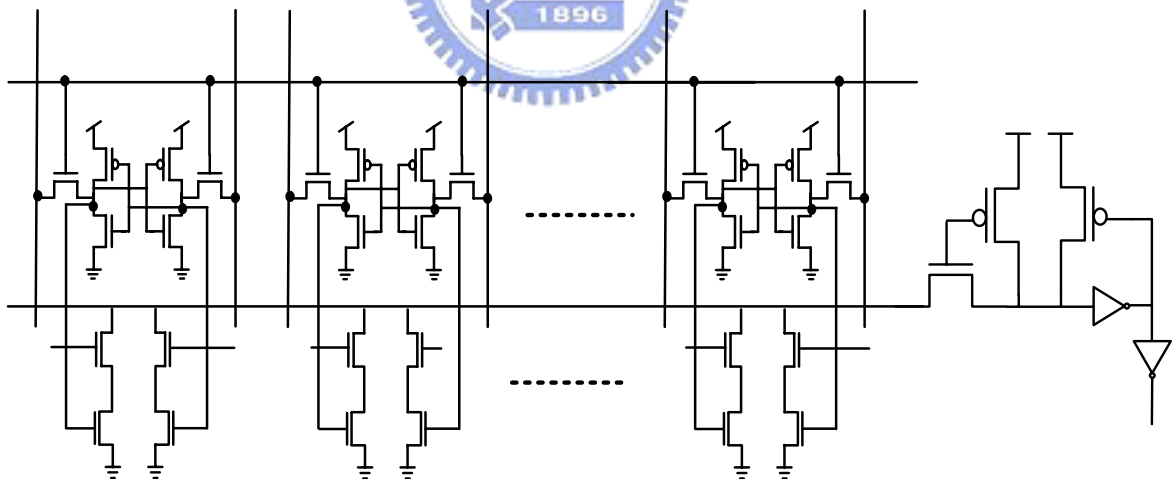


Fig. 3.13 Conventional NOR type 10T CAM for one word line

3.4.2 Low power Pre-comparison Circuit

Fig. 3.14 is a block diagram of the low power pre-comparison circuit. The low power pre-comparison circuit has an additional pre-comparison circuit for the low power CAM. The pre-comparison circuit is between the match line circuit of CAM as well as precharging circuit. Take an m bits pre-comparison circuit for example. When a comparison occurs, the m bits search data, $s_1 \sim s_m$, as well as stored data, $q_1 \sim q_m$, would be sent into the pre-comparison circuit to do comparison firstly. After the comparison is completed, the result of the pre-comparison will decide to turn on or turn off the discharging channel. The result of pre-comparison is also sent to the output logic circuit. Output logic circuit sends out the final result that is decided by two signals. First one is from the pre-comparison circuit and the other is from the precharging circuit. If the pre-comparison circuit is mismatch, then the final output value would be logic low regardless of the value of precharging circuit. Only the pre-comparison circuit is match, then the final output value is decided by the precharging circuit. So the output logic circuit can be described in a two bits truth table, in Table 3.1. From the truth table, the output logic circuit can be simplified to a two inputs and one output AND logic gate.

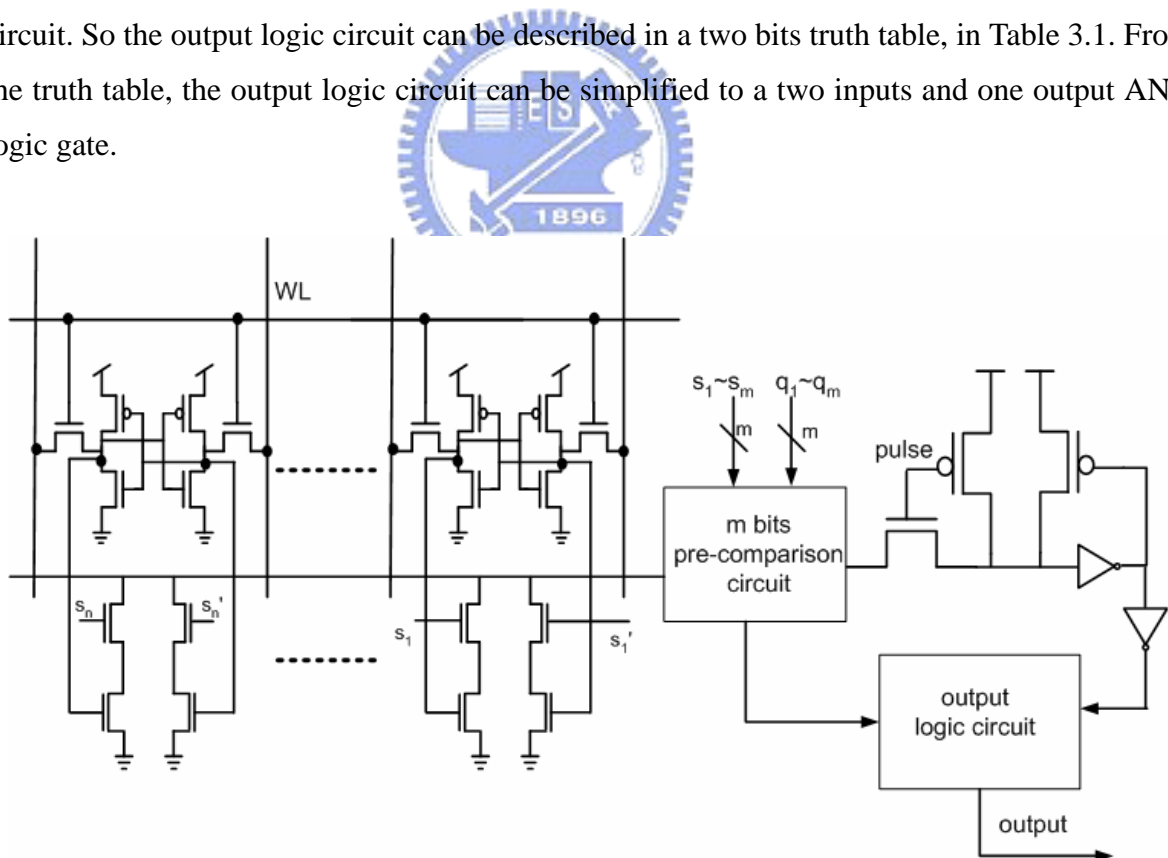


Fig. 3.14 Block diagram of low power pre-comparison CAM

Table 3.1 Truth table for Pre-comparison CAM

	0·(precharging)	1·(precharging)
0(pre-comparison)↵	Mismatch↵	Mismatch↵
1(pre-comparison)↵	Mismatch↵	Match↵

Fig. 3.15 is one bits of my low power pre-comparison circuit. Utilizing the logic function is as follows:

$$\begin{aligned}
 T &= s_1q_1 + \overline{s_1q_1} = \overline{\overline{s_1q_1 + s_1q_1}} \\
 &= \overline{[(s_1q_1)(s_1q_1)]}
 \end{aligned}
 \tag{3.6}$$

From the logic function, we adopt three NAND logic gates of complementary circuit. Only when the s_1 and q_1 are both logic zero or logic one, the compare bit would be match. The T would be just logic one when s_1 and q_1 are the same. When T is logic one, the transmission gate between precharging circuit, node A, and match line would be turned on due to T and T' to do evaluation. If the s_1 and q_1 are different, the T would be logic zero and turn off the transmission gate. On the other hand, the output logic just adopts an NOR gate, two inputs are T' and B. Fig. 3.16 is m bits pre-comparison circuit. The logic function is as follows:

$$\begin{aligned}
 T &= (s_1q_1 + \overline{s_1q_1}) + (s_2q_2 + \overline{s_2q_2}) + \dots + (s_nq_n + \overline{s_nq_n}) \\
 &= \overline{[(s_1q_1 + s_1q_1) + (s_2q_2 + s_2q_2) + \dots + (s_nq_n + s_nq_n)]} \\
 &= \overline{[(s_1q_1 + s_1q_1) \cdot (s_2q_2 + s_2q_2) \cdot \dots \cdot (s_nq_n + s_nq_n)]} \\
 &= \overline{\{[(s_1q_1)(s_1q_1)] \cdot [(s_2q_2)(s_2q_2)] \cdot \dots \cdot [(s_nq_n)(s_nq_n)]\}}
 \end{aligned}
 \tag{3.7}$$

The function can be implemented by using NAND logic gates to implement this function. It is the same as one bit pre-comparison circuit. The logic value, T, would decide that if the transmission gate turns on or turns off. The output logic circuit is also composed of by one NOR logic gate, signal \overline{T} and signal B are its inputs. Only the node A and T are both logic one, the output value would be logic one. The original output can be indicated as:

$$output = T \cdot A = \overline{\overline{T \cdot A}} = \overline{\overline{T} + \overline{A}} = \overline{\overline{T} + B} \quad (3.8)$$

So my output logic circuit is to use the NOR logic gate and two input signals are \overline{T} and B.

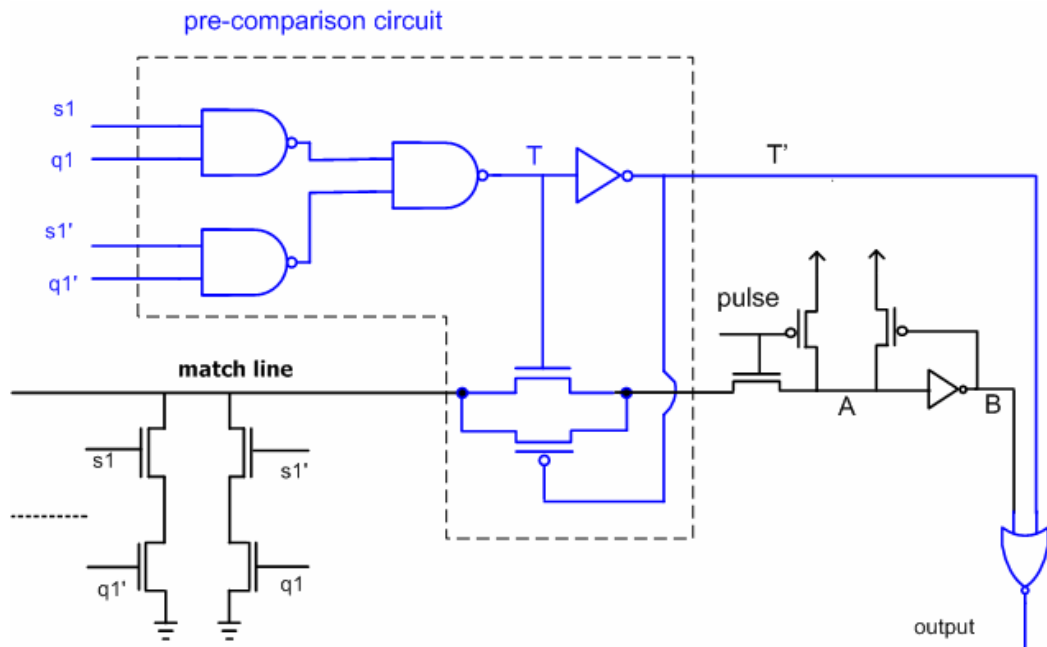


Fig. 3.15 One bit pre-comparison NOR type 10T CAM

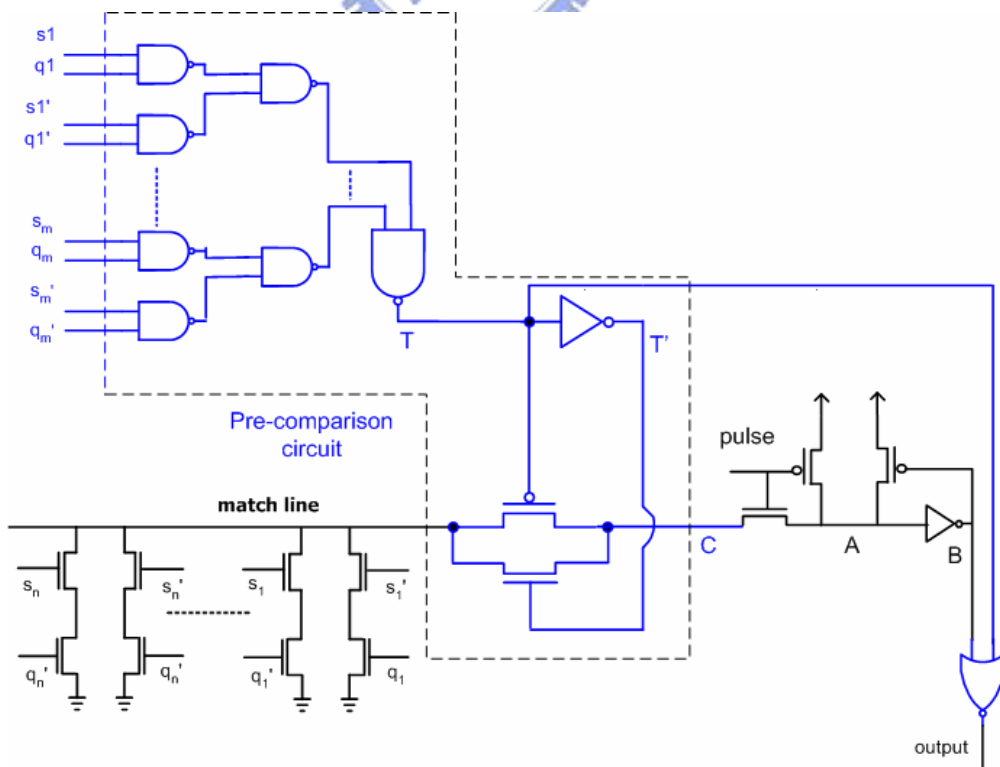


Fig. 3.16 M bits pre-comparison NOR type 10T CAM

In generally, as the bits for the pre-comparison circuit increases, the area of pre-comparison circuit increases too. Each additional bit would need additional transistors. Although more bits for pre-comparison circuit, the precision would be more accurate. But it would increase the circuit area and the delay for the pre-comparison circuit would also increase. It would make the performance down because of waiting the pre-comparison result. So, there is a tradeoff for the bit numbers of pre-comparison circuit, the speed of the pre-comparison circuit, and the area.

3.4.3 Power Analysis

From 3.1.1, we know the formula of dynamic power consumption. We do the power analysis to prove that the pre-comparison CAM is actual power saving. Firstly, we assume that any one bit for CAM has a 50% probability to be logic one or logic zero [33]-[34]. Each time comparison for CAM that has the match probability and mismatch probability for one bit is as follows:

$$P_{\text{one bit match}} = P_{\text{both1}} + P_{\text{both0}} = \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) + \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) = \frac{1}{2} \quad (3.9)$$

$$P_{\text{one bits mismatch}} = P_{\text{one 0, the other 1}} + P_{\text{one 1, the other 0}} = \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) + \left(\frac{1}{2}\right)\left(\frac{1}{2}\right) = \frac{1}{2} \quad (3.10)$$

We assume that each word line of CAM cell has n bits. The probability of word line match is

$$\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\dots\dots\left(\frac{1}{2}\right) = \left(\frac{1}{2}\right)^n \quad (3.11)$$

On the contrary, the probability of word line mismatch can be written as

$$\left(\frac{1}{2}\right)^1 + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^3 + \dots + \left(\frac{1}{2}\right)^n = 1 - \left(\frac{1}{2}\right)^n \quad (3.12)$$

We know that the dynamic power is proportional to supply voltage, loading capacitance, and switching frequency. The following discussion is to focus on the match line power consumption. We know that probability for circuit mismatch is $1-(1/2)^n$ for my assumption. The power dissipation would only occur for discharging the match line for mismatch. So the match line of power consumption can be simplified as follows [36]:

$$P = f \cdot (C_{\text{match}} + C_{\text{pre}}) \cdot VDD^2 \cdot (1 - (\frac{1}{2})^n) \quad (3.13)$$

C_{match} is the loading capacitance of match line and C_{pre} is the loading capacitance of the precharging node. When a mismatch occurs, the discharging path is turned on. So, the loading capacitance is the sum of match line capacitance and precharging node capacitance. My low power pre-comparison circuit has three nodes about the discharging path. First node and second node are the match line capacitance and precharging node capacitance, C_{match} and C_{pre} . The third node capacitance is because of the transmission gate in fig, node C and its capacitance is marked as C_c . For my low power pre-comparison circuit, the loading capacitance is the sum of C_{match} , C_{pre} , and C_c . It assumes that there are m bits for the low power pre-comparison circuit. If the pre-comparison bits are some mismatch, the transmission gate would not turn on. Only when the pre-comparison circuit is all match, the nMOS between precharging circuit and match line circuit would be turn on to evaluate. So, the probability for pre-comparison circuit match is $(1/2)^m$. And the probability for match line circuit mismatch is $(1-(1/2)^n)$. So the power consumption can be written:

$$P = f \cdot (C_{\text{match}} + C_{\text{pre}} + C_c) \cdot VDD^2 \cdot (1 - (\frac{1}{2})^n) \cdot (\frac{1}{2})^m \quad (3.14)$$

The value of C_c would not be very large, so the capacitance does not increase more. The power supply and switching frequency are the same. The other parameter is the probability. The value of $(1-(1/2)^n) \cdot (1/2)^m$ must be smaller than the value of $(1-(1/2)^n)$. We take $m=1$ and $n=10$ for example. The former value is $(1-(1/2)^{10}) \cdot (1/2)^1 = 1023/2048$. The latter value is $(1-(1/2)^{10}) = 1023/1024$. We can obviously obtain that $1023/2048$ is smaller than $1023/1024$, so my pre-comparison NOR type 10T CAM would be more power saving than the conventional NOR type 10T CAM.

3.4.4 Power Simulation

There are some simulations about the CAM comparison. Using some test patterns run simulation. The CMOS technology is TSMC 0.13um, in Table 3.2. The voltage is 1.2V, and temperature is 25°C. The frequency of the clock that is used to precharge the match line periodically is 500 MHz, and its duty cycle is 50%.

There are some test patterns on to observe the different condition of CAM and different size of CAM array. The simulation result of Fig. 3.17 is the comparison of conventional 10T CAM cell and one bit pre-comparison CAM cell for different CAM array, 32 (word size) x 4 (bit size), 32x8, 32x16, and 32x32. From the result, power dissipation for smaller array size, 32x4 and 32x8, is not only reduced, but also increases. This is because the penalty of pre-comparison circuit that produces the additional power dissipation is greater than the reducing power dissipation by pre-comparison mechanism. When the bit size is bigger sixteen, the power consumption has obviously reduced in Fig. 3.17. So, if the pre-comparison CAM wants to have better power efficiency, the array size of CAM must be greater. Otherwise its characteristic of pre-comparison for low power can not have good outcomes. So the following simulation would be based on the 32x32 CAM array for simulation.

Table 3.2 Simulation technology

technology [↗]	TSMC·0.13um [↗]
voltage [↗]	1.2V [↗]
temperature [↗]	25°C [↗]
frequency [↗]	500MHz [↗]
duty cycle [↗]	50% [↗]

power consumption

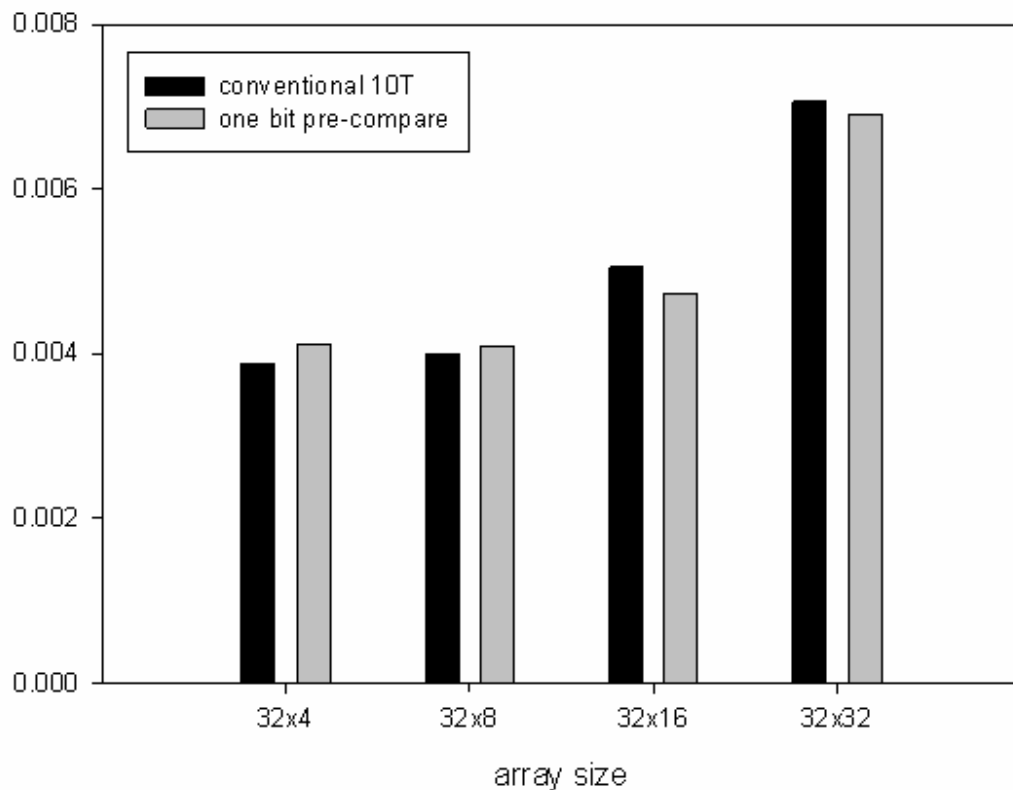


Fig. 3.17 Power consumption of conventional and low power pre-comparison CAM array



Another simulation result is in Fig. 3.18. Three kinds of CAM array are simulated, and their sizes are the same 32x32. First one column is used for contrast. It is the conventional NOR type 10T CAM cell. The other two CAM is my low power pre-comparison CAM, one bit pre-comparison and two bits pre-comparison CAM cell. We can find out a phenomenon. As the pre-comparison bit number increases, the power dissipation is reduced more obviously. For one bit pre-comparison CAM cell, it has 6% power reduction. For two bits pre-comparison CAM cell, it has 19.8% power reduction. The more bit number of pre-comparison CAM cell would reduce more power dissipation because my pre-comparison CAM cell can determine match or mismatch more precisely for more pre-comparison bit in advance.

power consumption

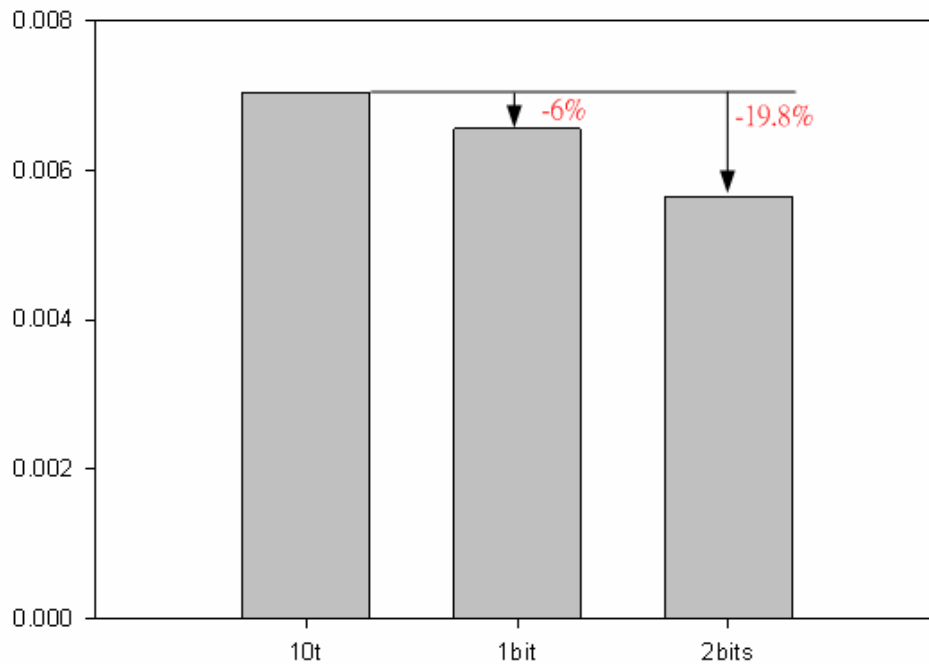


Fig. 3.18 Power consumption of conventional and low power pre-comparison CAM array



3.5 Conclusions

In this chapter, some low power design of CAM has introduced. They maybe improve the architecture or CAM cell to obtain lower power consumption. We also propose a pre-comparison concept. We apply the pre-comparison circuit into NOR type 10T CAM. The low power pre-comparison architecture also has two comparison parts such as divides match line in Fig. 3.12. If the pre-comparison is match, the second part of comparison would be compared to decide if the CAM is match or not. If the pre-comparison is mismatch, it would not compare second part of CAM cell and it is mismatch. My low power pre-comparison CAM cell also can lower the power consumption through the simulation with TSMC 0.13um CMOS technology. For the pre-comparison NOR type 10T CAM, the one bit pre-comparison CAM can reduce 6% power consumption, and two bits pre-comparison CAM can reduce 19.8% power consumption for the array size is 32x32. So this architecture that we propose is a feasible method for low power design.

Chapter 4

Enhancement of Low power Pre-comparison CAM

Chapter3 has introduced low power pre-comparison CAM. In chapter3, simulation shows that my pre-comparison CAM actually has lower power dissipation. But there is still some portion to improve. We would talk about the improvement for the low power pre-comparison CAM in this chapter. And we would make some comparison about its power consumption, match line capacitance, and access time for match line in the following section.

4.1 Match Line Capacitance

Table 4.1 is the capacitance of match line for the low power pre-comparison NOR type 10T CAM and conventional NOR type 10T CAM. We can find out one characteristic for the low power pre-comparison CAM. As the bit of pre-comparison circuit increases, the capacitance of the pre-comparison CAM would increase, too. This is because the pre-comparison circuit is a duplicate of XOR logic gate for each bit of CAM cell. On the other hand, my pre-comparison CAM cell needs another transmission gate to be a switch, and it would also increase the match line capacitance. Although, my pre-comparison CAM cell would lower more power dissipation for higher pre-comparison bits However, the higher pre-comparison bits would raise the match line capacitance and degenerate the power consumption. How to solve this problem is important. There will be some improvement for match line capacitance, and simulate for power consumption, access time and so on. And find out the optimal pre-comparison bits for CAM in the following chapter by considering power, timing, and area overhead.

Table. 4.1 The capacitance for different architecture of CAM

capacitance	1bit pre-comparison	2bits pre-comparison
Pre-comparison CAM	41.9292f	42.1368f
Conventional 10T CAM	41.3126f	

4.1.1 Enhancement of Low power Pre-comparison Circuit

From Fig. 4.1, the low power pre-comparison CAM would pre-compare first part to see whether the first part is match or not. If it is match, the comparison of second part would start to activate. From Fig. 4.1, there would be a cost of duplicate comparison portion for pre-comparison circuit. That is, the first part is compared and would be compared in the second part if the comparison of first part is match. We also can obviously find out that the same comparison bits have been compared twice in pre-comparison circuit and XOR logic gate that are printed in blue lines. Actually, if the bits are match in the pre-comparison circuit, then the bits would be match in the match line comparison for second part. For the same reason, if the bits are mismatch in the pre-comparison part, then it would be mismatch in the match line comparison. That is, there is no need for the repeatedly comparison twice for my low power pre-comparison CAM, so we make some improvement of my pre-comparison CAM.

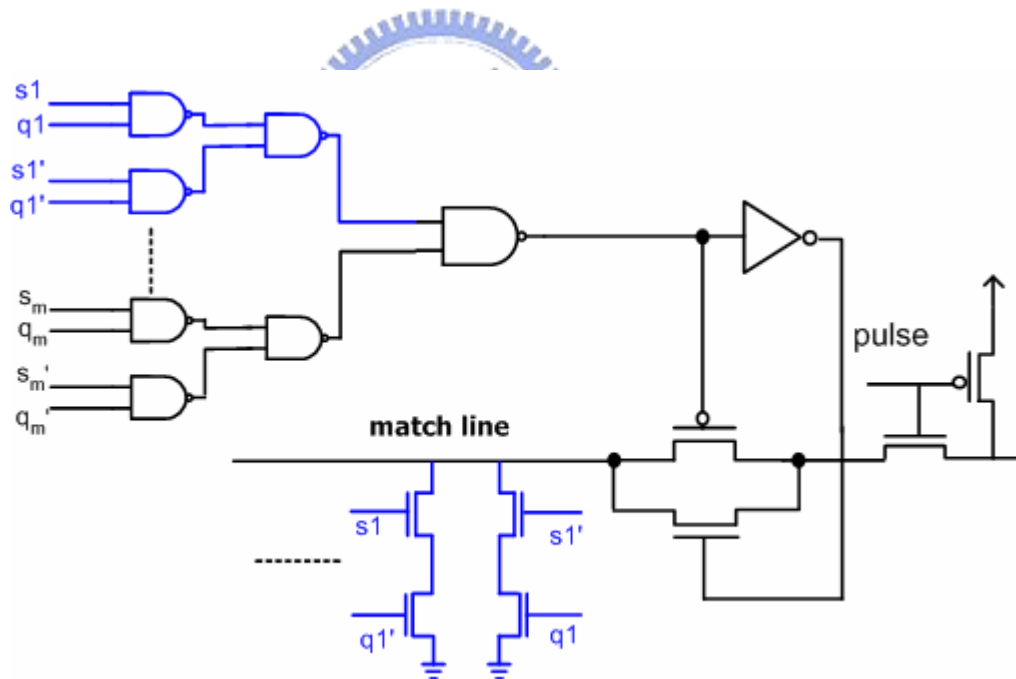


Fig. 4.1 The original pre-comparison CAM

In Fig. 4.2, that is the improvement of one bits low power pre-comparison circuit. Obviously, the difference of Fig. 4.1 and Fig. 4.2 is that match line has reduced one bit comparison, s_1 , \bar{s}_1 , q_1 , \bar{q}_1 . The improvement is used for m bits pre-comparison circuit as well as n bits CAM in Fig. 4.3. In Fig. 4.3, the m bits pre-comparison has been reduced in the

match line circuit, so the bit number in the match line is only $(n-m+1)$ bits and $(n-m+1)$ XOR gates to be used. If the pre-comparison bit number is bigger, the reduced bit number in the match line circuit is more.

There is some advantage for the improved pre-comparison CAM. First of all, the reduced XOR logic gates in the match line circuit can reduce the area. The original CAM is composed of XOR logic gate and SRAM. The reduced m bits XOR logic gates in my low power CAM means that there are $(n-m+1)$ bits CAM and only m bit SRAM. The area is smaller for the original CAM array. Secondly, the reduced m bits XOR logic gates can lower the capacitance of the match line. Dynamic power consumption is proportional to the loading capacitance. If the smaller loading capacitance is, the more power dissipation can be saved. The improvement of the low power pre-comparison CAM can reduce more power dissipation than original pre-comparison CAM because of the smaller loading capacitance of match line. We make some simulation in the following section.

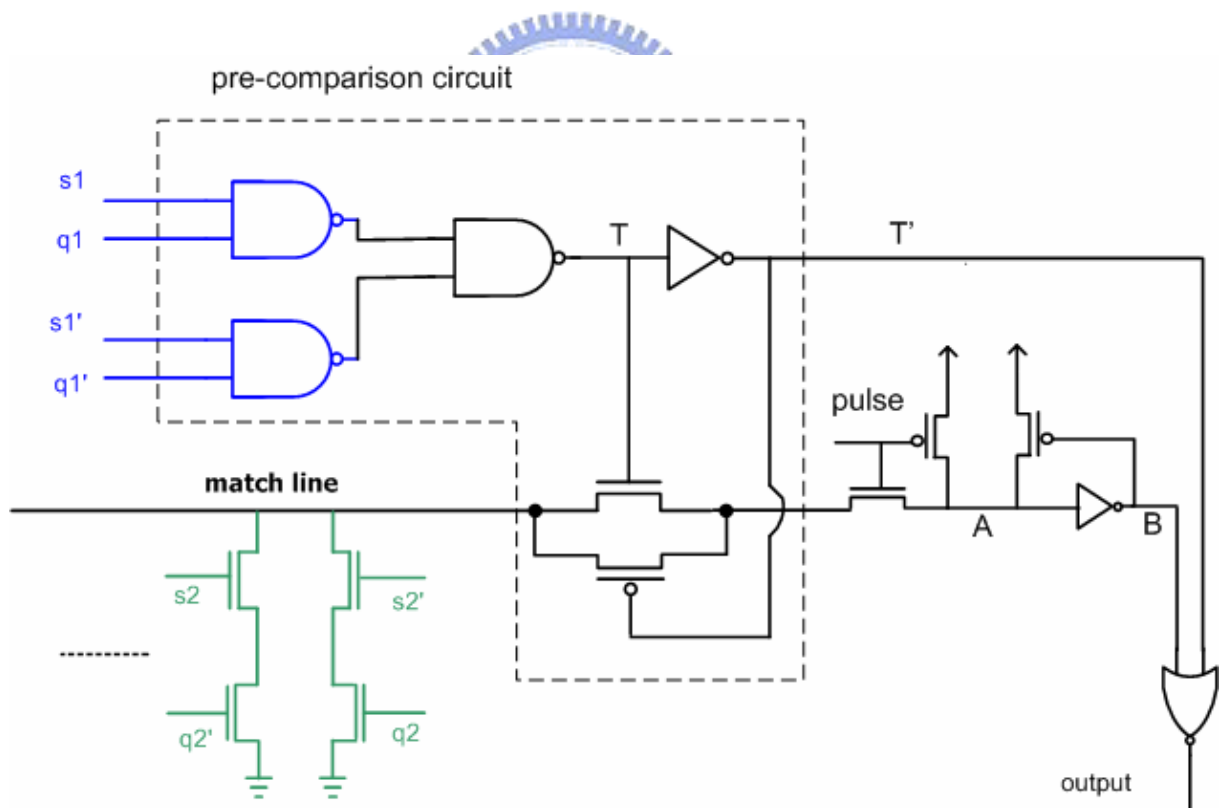


Fig. 4.2 The enhancement of one bit low power pre-comparison CAM

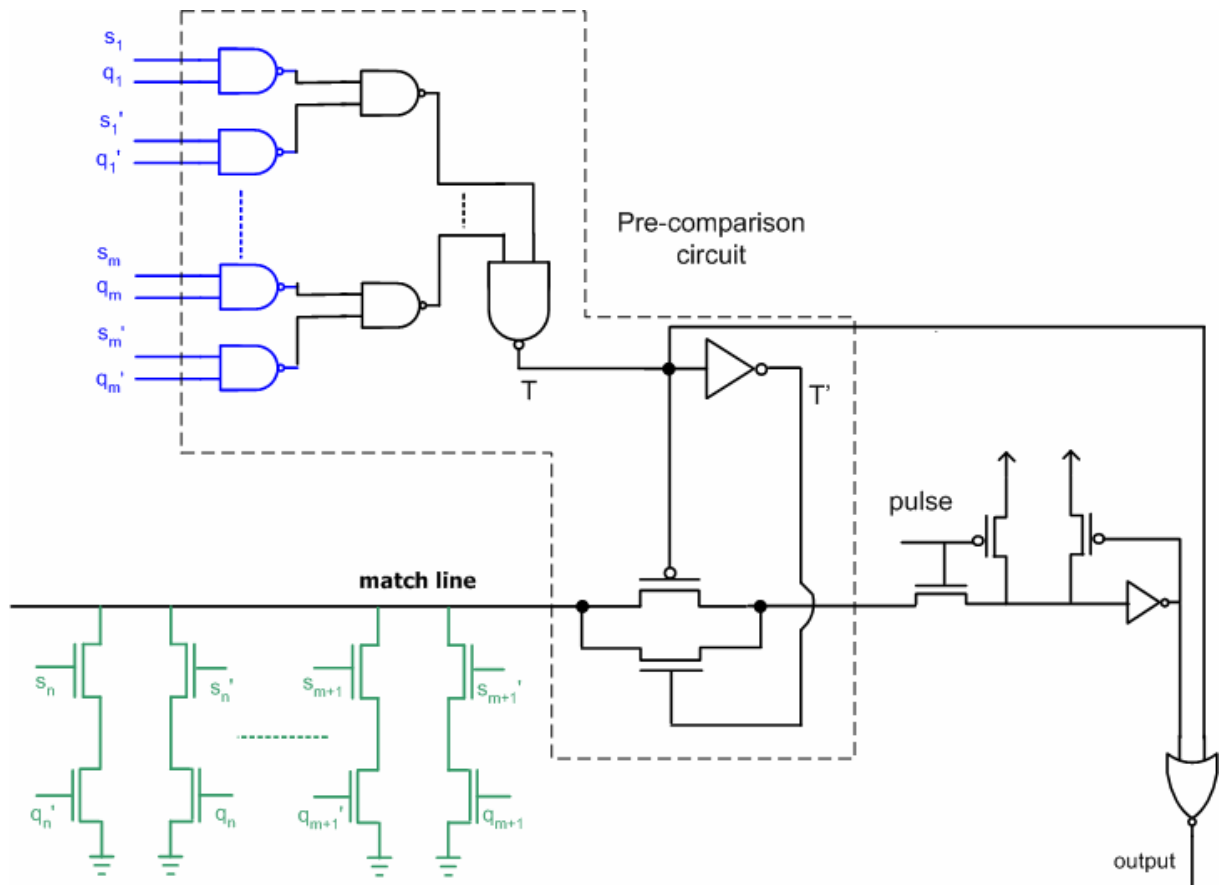


Fig. 4.3 The enhancement of m bit low power pre-comparison CAM



4.1.2 Match Line Capacitance Analysis

The enhancement of pre-comparison circuit has smaller capacitance of match line in section 4.1. Through simulation, the comparison of the match line capacitance is in Fig. 4.4. It is obviously that the enhanced pre-comparison CAM can reduce 3% match capacitance for one bit pre-comparison CAM and reduce 6.5% match line capacitance for two bit pre-comparison CAM. If the enhanced pre-comparison bit increases, the lower match line capacitance would be. We make the capacitance comparison for conventional NOR type 10T CAM and improved pre-comparison CAM, and the result shows in the Fig. 4.5. One bit improved pre-comparison CAM can reduce 1.5% match line capacitance. Two bits improved pre-comparison CAM can reduce 4.5% capacitance. In the later of two bits improved pre-comparison CAM, the more one bit for the improved pre-comparison CAM, the capacitance can be reduced more 3%. For example, five bits improved pre-comparison CAM has reduced 13.7% match line capacitance. So, the lower loading capacitance is, the lower

power dissipation can be reduced.

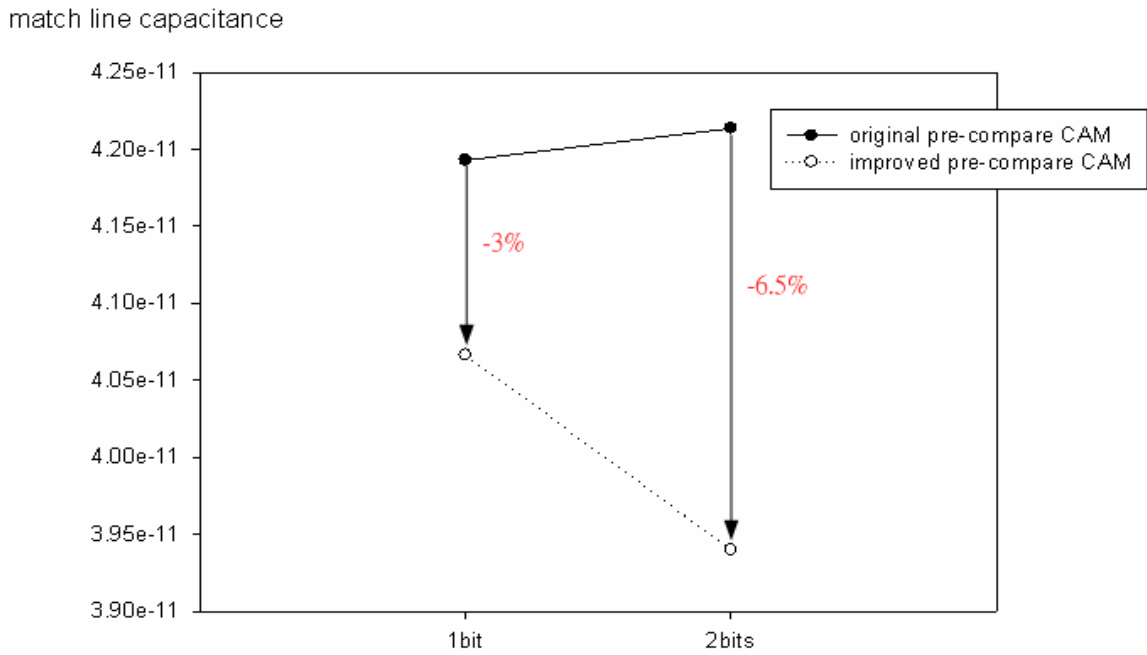


Fig. 4.4 Match line capacitance for original and enhancement pre-comparison CAM

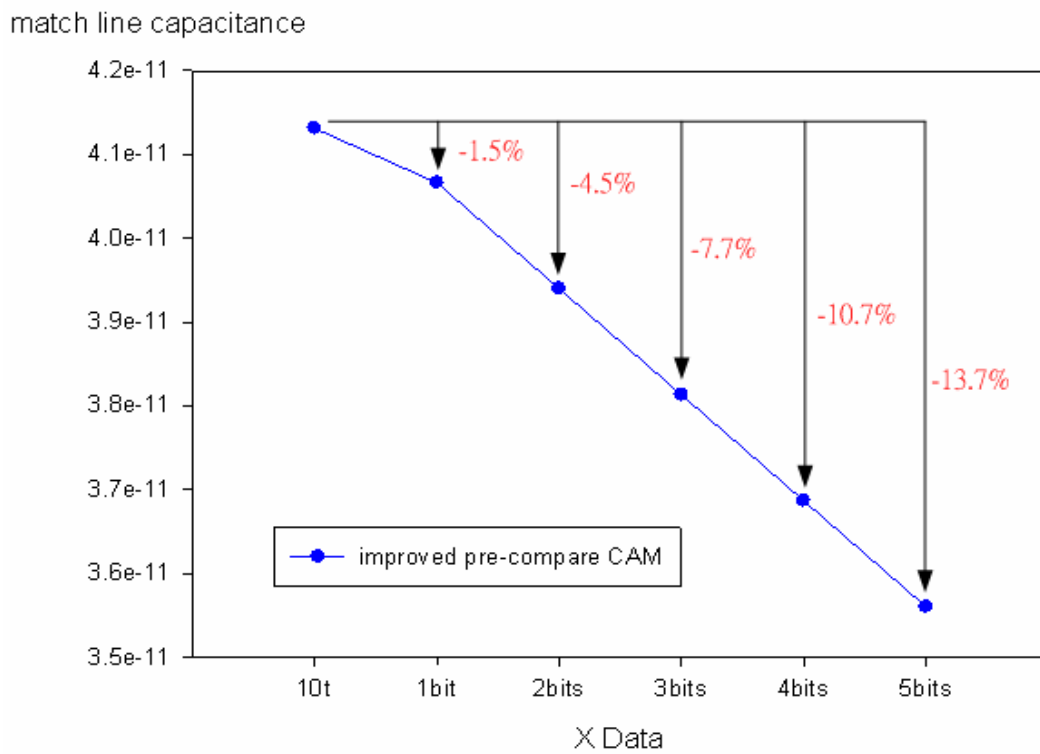


Fig. 4.5 Match line capacitance is related with pre-comparison bits

4.2 Prediction of Precision and Power Analysis

The different bit number of improved pre-comparison CAM not only affect the match line capacitance, but also affect the prediction of precision for the first part of comparison. If the prediction of precision is more precise, the performance of the first part comparison is better. If the efficiency of first part comparison is better, the power dissipation would be reduced more. Although the more pre-comparison bits have better precision, the area penalty of pre-comparison circuit and the access time of pre-comparison circuit are also considered. We have improved the low power pre-comparison CAM to reduce the match line capacitance in the above section. Prediction of precision and power consumption analysis would be discussed in the following section.

4.2.1 Prediction of Precision for Different Pre-comparison Bits

As the pre-comparison bits increase, the pre-comparison result would be more precise for the first comparison. As the pre-comparison bits increase, the pre-comparison circuit would become larger and complex. So, the speed of deciding match or mismatch would be slower. Because more logic circuit must judge that the result of pre-comparison is match or not. On the other hand, the more logic circuit would make the power dissipation not decrease but increase. So there is a tradeoff about the bit numbers for pre-comparison circuit. The bit numbers of the pre-comparison CAM would affect the precision of comparison, access time, and power dissipation

We assume that each bit has the 50% probability for match or mismatch of comparison. For the n bits conventional NOR-type CAM cell, there are 2^n-1 kinds of condition for mismatch and there is only one condition for all n bits match. The probability of all n bits for a word line match is

$$\left(\frac{1}{2}\right) \cdot \left(\frac{1}{2}\right) \cdot \dots \cdot \left(\frac{1}{2}\right) = \left(\frac{1}{2}\right)^n \quad (4.1)$$

On the contrary, the probability of a word line not mismatching is

$$\left(\frac{1}{2}\right) + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^3 + \dots + \left(\frac{1}{2}\right)^n = 1 - \left(\frac{1}{2}\right)^n \quad (4.2)$$

In Fig 4.6, it is the mismatching probability that bit numbers can predict the mismatching probability. If the bit number is one, the probability of match or mismatch is 50%. As the bit number increases, the mismatching probability would increase. When bit numbers have five bits or more, the mismatching probability is very close to 100%. Fig. 4.7 is the really experimental data from paper. It shows a very important thing. Over 90% of the tag line mismatches are determined within the four least significant bits (LSB) of the tags.

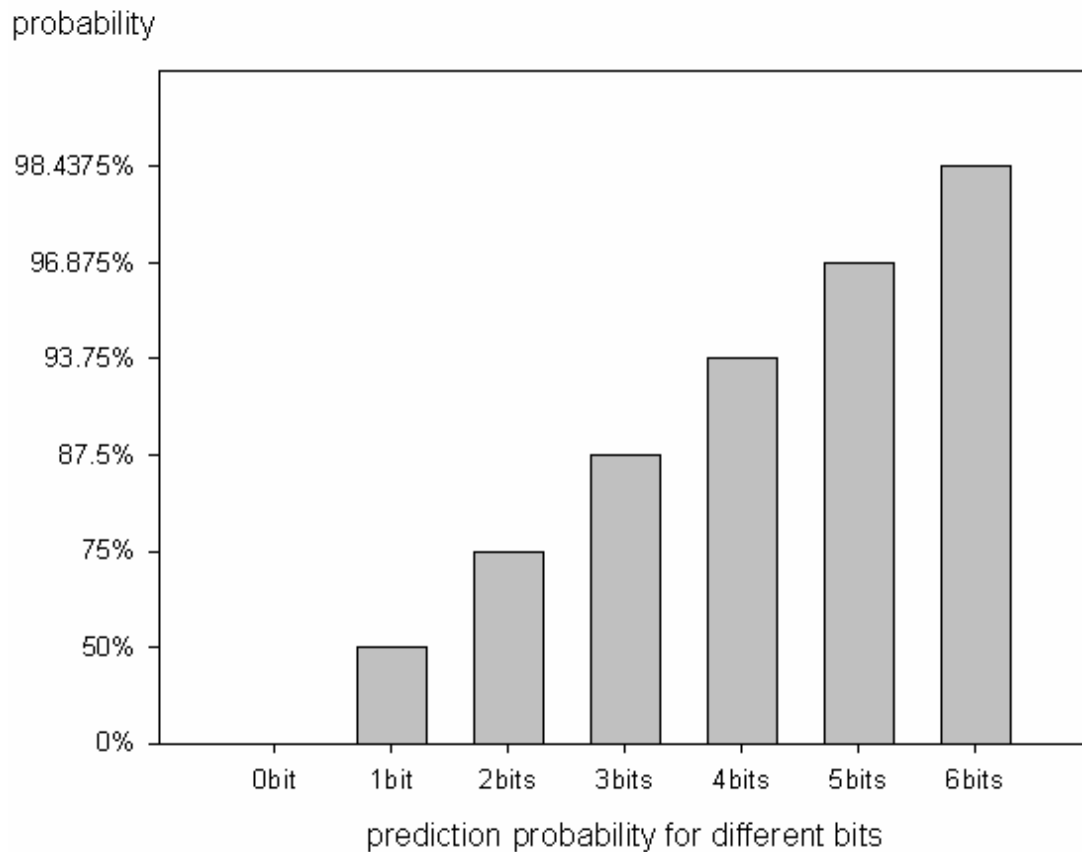


Fig. 4.6 Mismatch probability for different pre-comparison bits

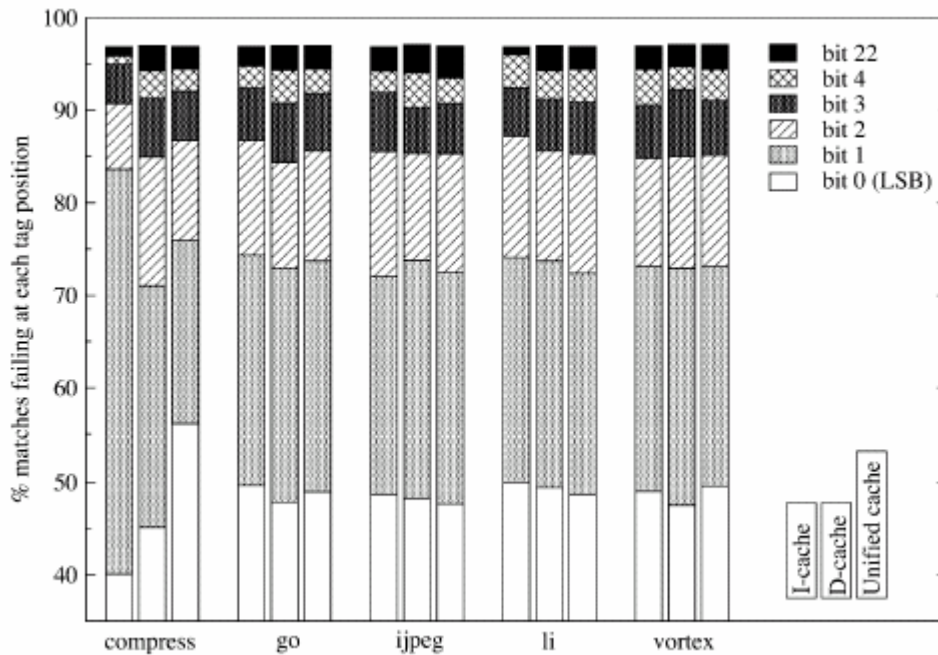


Fig. 4.7 Percentage of tag checks ending at each bit position

So my low power pre-comparison CAM also applies this concept to reduce the times of discharging for the match line circuit. My low power pre-comparison CAM is used to pre-compare the mismatching bits. The pre-comparison efficiency can improve very obviously as the pre-comparison bit increases. We can see that bit numbers go to five or six bits. Its probability is very close to 100%. So maintain the bit numbers smaller six bits is feasible. We would make some comparison for the bits numbers that are used by pre-comparison circuit. Using some test patterns test different bit numbers to find out the optimal bit numbers from the view of power dissipation. Before simulating the different test patterns, we use formula to analyze my low power pre-comparison CAM.

4.2.2 Power Consumption Analysis for Match Line

In chapter3, the power consumption of match line has analyzed. The power dissipation of match line for the traditional CAM cell and the low power pre-comparison CAM cell are as follows:

$$P = f \cdot (C_{\text{match}} + C_{\text{pre}}) \cdot VDD^2 \cdot (1 - (\frac{1}{2})^n) \quad (4.3)$$

$$P = f \cdot (C_{\text{match}(n)} + C_{\text{pre}(n)} + C_{\text{c}(n)}) \cdot VDD^2 \cdot (1 - (\frac{1}{2})^n) \cdot (\frac{1}{2})^m \quad (4.4)$$

The formula is some different about the C_{match} and C_{pre} . In chapter 3, the original CAM and pre-comparison CAM have the same value of C_{match} and C_{pre} . That is out assumption that these two values are the same for the original CAM and pre-comparison CAM. Actually, the value of C_{match} and C_{pre} for conventional and pre-comparison CAM are different. From my simulation, the pre-comparison match line capacitance, called $C_{\text{match}(n)}$, is a little bigger than the conventional match line capacitance, called C_{match} . The pre-comparison precharging capacitance, $C_{\text{pre}(n)} + C_{\text{c}(n)}$, is a little bigger than the conventional precharging capacitance, C_{pre} . Why are the capacitance some different? Because the size of the transistors of precharging circuit, M1~M4 in Fig. 4.8, are different from the traditional CAM cell and pre-comparison CAM cell.

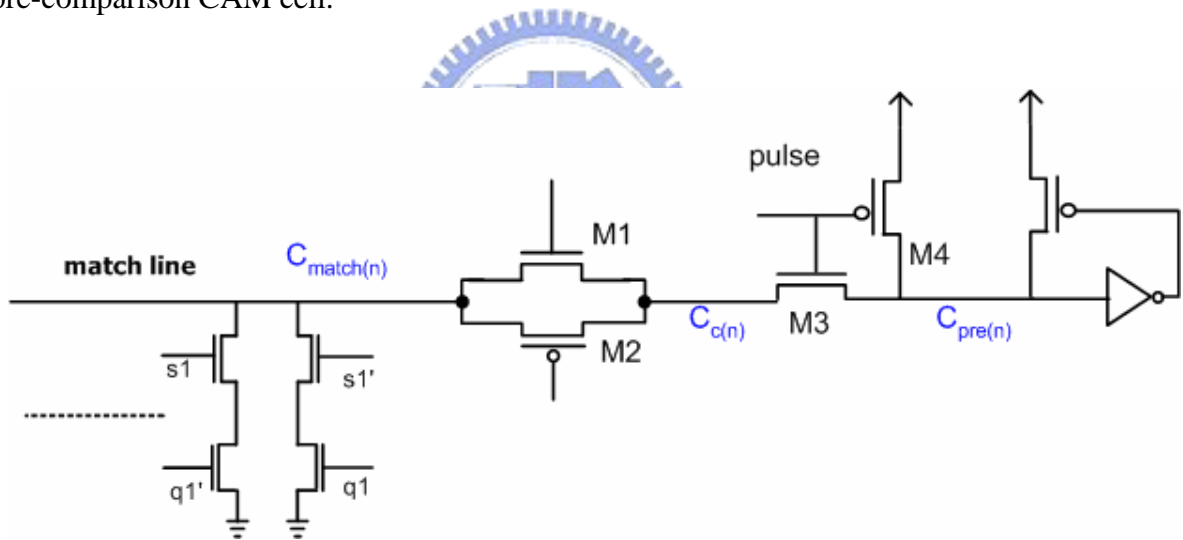


Fig. 4.8 One bit low power pre-comparison circuit

In the following text, we would analyze the power consumption of match line for the enhancement of pre-comparison CAM. The pre-comparison circuit is the same as pre-comparison CAM for improvement of pre-comparison CAM. So we assume that the value of C_{pre} and C_{c} is the same as the former. The match line capacitance of the improvement of pre-comparison CAM is abbreviated to $C_{\text{match}(n-m)}$, in Fig. 4.9. We know that the $C_{\text{match}(n-m)}$ is smaller than the C_{match} , so the value of sum of C_{pre} , C_{c} and $C_{\text{match}(n-m)}$ would be smaller the sum of $C_{\text{match}} + C_{\text{pre}} + C_{\text{c}}$. On the other hand, the XOR logic gates of match line have been

reduced m bits. So, there are remaining (n-m) bits that would be compared. When pre-comparison circuit is all match, the nMOS between precharging circuit and match line circuit would be turn on to evaluate. We assume that a 50% probability of a match at any bits.

So, the probability for pre-comparison circuit matching is the same as $(\frac{1}{2})^m$. And the probability for match line circuit mismatch is

$$\left(\frac{1}{2}\right) + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^3 + \dots + \left(\frac{1}{2}\right)^{n-m} = 1 - \left(\frac{1}{2}\right)^{n-m} \quad (4.5)$$

So the power consumption can be written as follows:

$$P = f \cdot (C_{\text{match}(n-m)} + C_{\text{pre}(n-m)} + C_{c(n-m)}) \cdot VDD^2 \cdot \left(1 - \left(\frac{1}{2}\right)^{n-m}\right) \cdot \left(\frac{1}{2}\right)^m \quad (4.6)$$

The value of $\left(1 - \left(\frac{1}{2}\right)^{n-m}\right)$ must be smaller than the value of $\left(1 - \left(\frac{1}{2}\right)^n\right)$ from formula 4.4 and formula 4.6. The value of capacitance and the value of probability are smaller for the improvement of pre-comparison CAM. In other words, we have improved the pre-comparison CAM further and have lower power dissipation than the original low power pre-comparison CAM from my analysis.

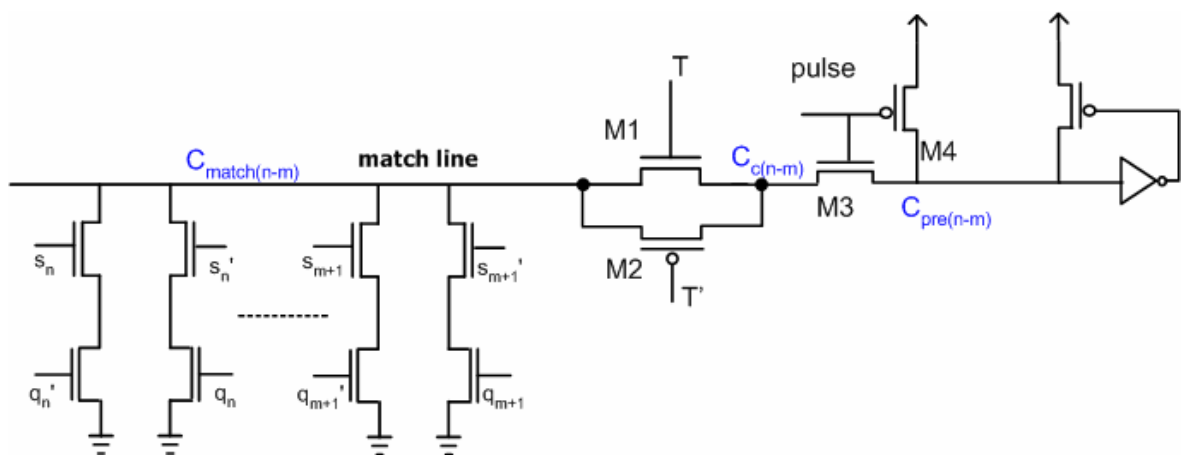
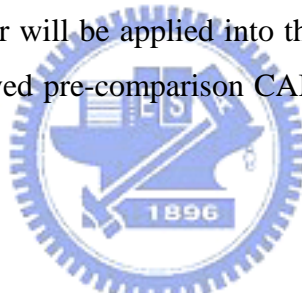


Fig. 4.9 M bits enhanced pre-comparison CAM

4.3 Simulation Result

It has analyzed for section 4.2. Therefore, there are some practical test pattern simulations to prove that the improved pre-comparison CAM consumes less power consumption. It also wants to compare the power consumption for different pre-comparison bit to see which bit number can save much power consumption. The access time is important, too. Of course, as the bit number increases, the access time would increase more. There is a tradeoff for access time and power consumption. The power delay product is a good pointer to find out the better bit number for improved pre-comparison CAM. On the other hand, as the bit number increases, the transistor of the improved pre-comparison CAM would also increase. So, it would also talk about the transistor penalty for different bit number of improved pre-comparison CAM. Finally in this section, we would find out the optimal bit number by considering many reasons, such as reduced power consumption, access time, time delay, transistor penalty, and so on. After we decides the optimal improved pre-comparison bit number, the optimal bit number will be applied into the following discussion, comparison of conventional CAM and improved pre-comparison CAM, layout of improved pre-comparison CAM, and etc.



4.3.1 Power Consumption for Different Pre-comparison Bits

In this section, using some simulations observe the power consumption for different bit number of improved pre-comparison CAM. We also simulate the conventional 10T CAM for comparison. There are five different test patterns to compare, and runs six cases, such as conventional 10T, and one bit to five bits improved pre-comparison CAM. The simulation result is listed in Fig. 4.10 for five test patterns and six cases. From Fig. 4.10, we can discover that improved pre-comparison CAM is less power consumption for one bit to five bits than the conventional 10T CAM. For different test patterns, the each different bit number has different power saving.

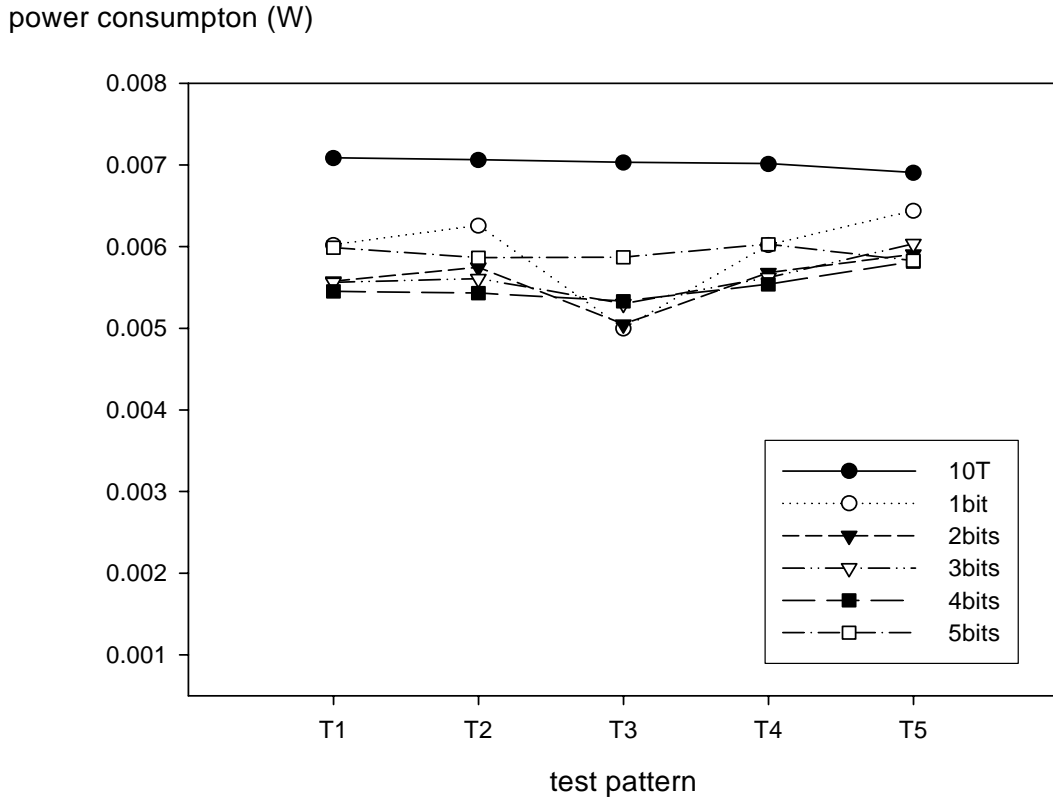


Fig. 4.10 Simulation for different different bit numbers of enhanced pre-comparison CAM



The average power consumption for different cases is listed in Fig. 4.11. From the figure, we can know a simple conclusion. When the bit number is two, three, and four bits, the reduced power consumption is best, 21.6%~22.8%. We also have another conclusion. As the bit number increases, the power consumption is not always reduced more. The five bits is an example. Its power saving, 15.7%, is less than the four bits, 22.8%. This is because the pre-comparison circuit would also occupy some area penalty and consume additional power. On the other hand, the prediction of precision for four bits and five bits are close. For the same test patterns, five bits average power consumption is obviously worse than the four bits. If we want to find the optimal bit number of improved pre-comparison CAM, the candidate is two bits to four bits on the pure consideration for power dissipation.

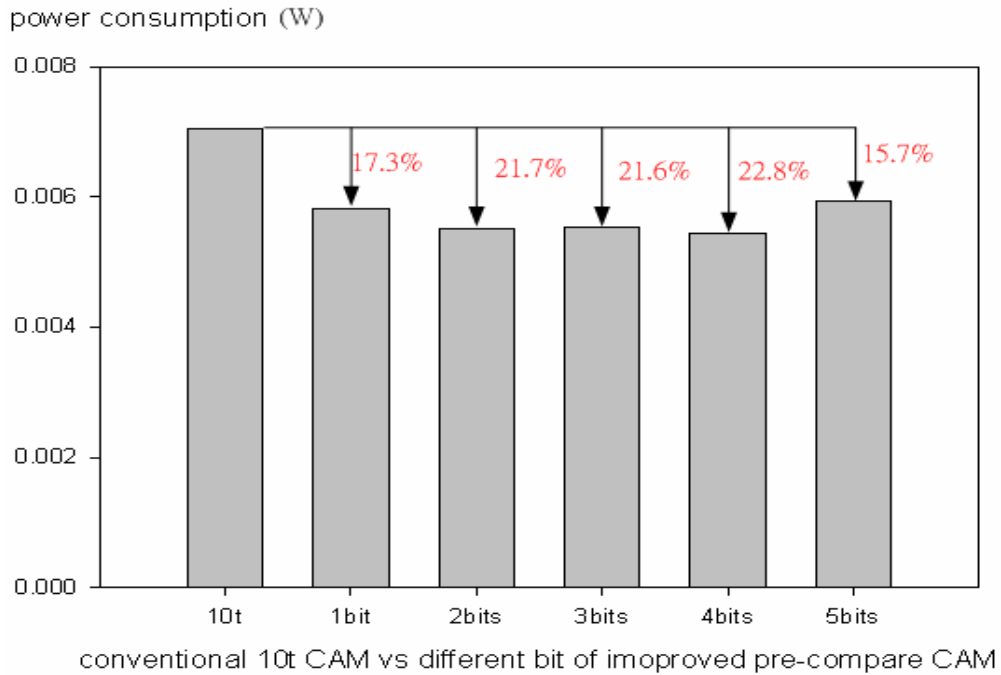


Fig. 4.11 Power consumption for conventional and enhanced pre-comparison CAM

4.3.2 Access Time



Fig. 4.12 shows the definition of access time. The access time is defined as the timing between match precharging signal that goes high and match line output that sends the determination value. Precharging signal goes high is in evaluation operation, and the comparison will proceed. Access time for CAM means that speed of comparison. That is to say, access time can be thought as performance. If the access time is smaller, the speed of the circuit is faster and the performance is better. We know that the match line is precharging periodically. In precharging section, the match line would go to logic high. In evaluation section, there are two kinds of condition. If the match line is match, the match line would stay logic high. If the match line is mismatch, the match line would discharge to logic low. The access time would be measured in the mismatch state. Because only the logic state changes, the access time could be measured.

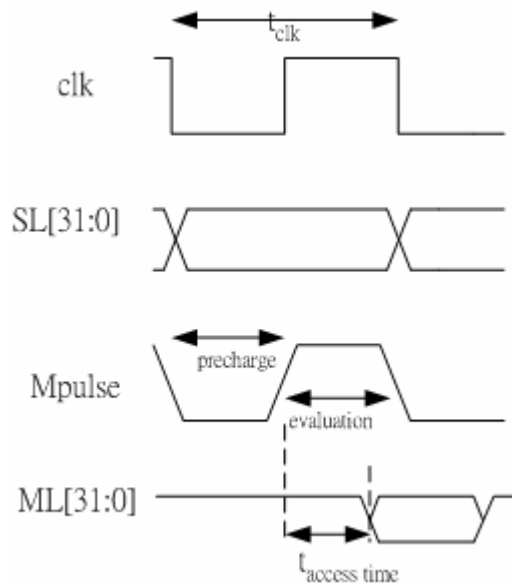


Fig. 4.12 Timing diagram of CAM during the search operation

For my improved pre-comparison CAM, the access time has some different considerations. There are three kinds of condition in my improved pre-comparison CAM. First kind of condition is that the first part comparison is mismatch. When first part comparison is mismatch, the result of second part comparison is not important and the result must be mismatch. Because the first part of comparison is mismatch, the mismatch information would be sent to output logic earlier and let the output be decided more quickly. In this situation, the access time would be faster than the conventional 10T CAM. Second kind of condition is that the first part comparison is match and the second part comparison is match, and the output result must be match. Third kind of condition is that the first part comparison is match and the second part comparison is mismatch, and the output result would be mismatch. The access time of third situation is longer than the conventional 10T CAM. Because the improved pre-comparison CAM has twice comparison, its access time is the addition of first part comparison and second part comparison. So its access time would be longer than the conventional 10T CAM.

The access time is measured for pre-comparison CAM and improved pre-comparison CAM for one pre-comparison bit and two pre-comparison bits. The comparison result is shown in Fig. 4.13. Using the improved one bit pre-comparison CAM can be the basis. From the figure, the access time of original pre-comparison CAM is worse than the improved pre-comparison CAM. Especially, the access time of original n2 has increased 40.6% when

pre-comparison bit is match. So the improved pre-comparison CAM not only reduces the capacitance of match line, but also reduces the access time.

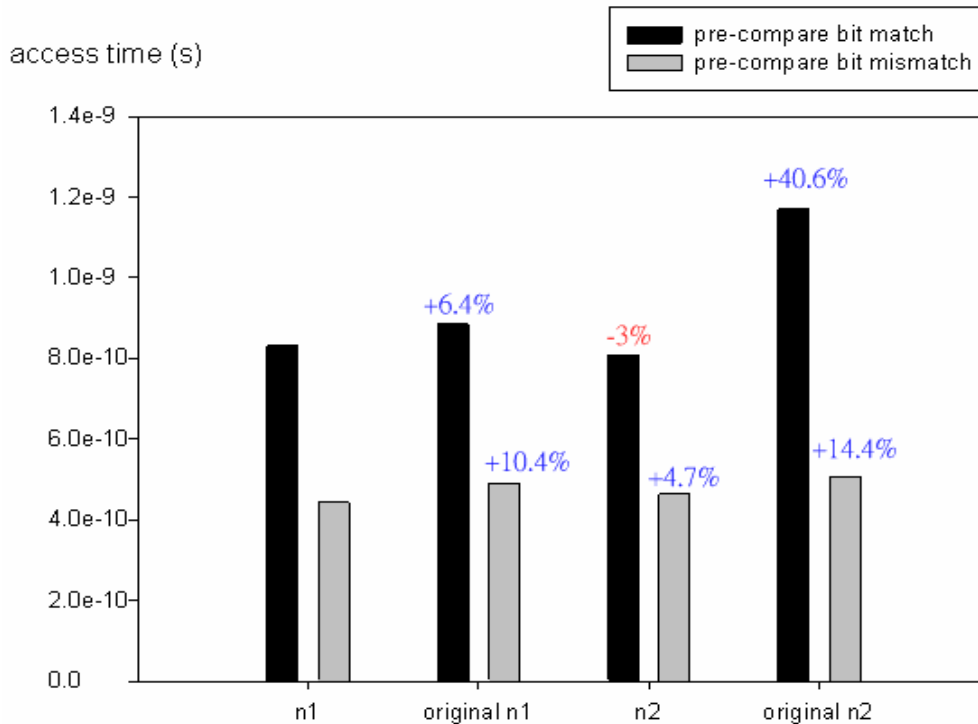


Fig. 4.13 The access time of pre-comparison and enhanced pre-comparison CAM



After comparing the original pre-comparison CAM, the measurement of the access time for conventional 10T CAM and different bit of improved pre-comparison CAM are performed, from one pre-comparison bit to five pre-comparison bits. The comparison result is shown in Fig. 4.14. The black line in the figure is when pre-compare bits are match. It is third kind of condition that is discussed in the above text. The access time of different pre-comparison bit numbers, from one bit to five bits, are longer than the conventional 10T CAM from 9.5% to 20.2%. On the other hand, the gray line in the figure is when the pre-comparison bits are mismatch. This is the first kind of condition, the access time can be reduced and smaller than the conventional 10T CAM from 35.1% to 39.9%. From the Fig. 4.13, we can also find out that the pre-comparison bit number from two bits to five bits. As its bit number increases, the access time for two kinds of condition are increasing step by step. Only the access time of one bit is some different. The other bits have this regular condition. So, if the pre-comparison bit increases, the access time penalty would increase to lower the performance for comparison.

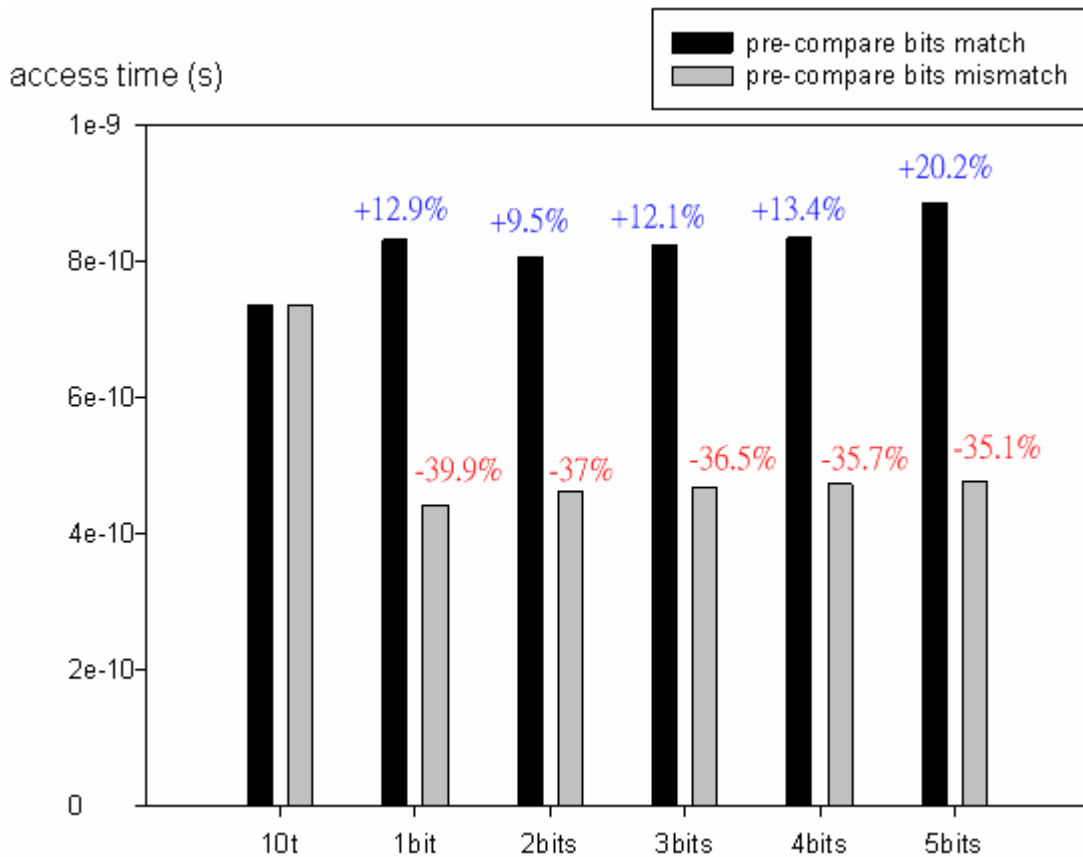


Fig. 4.14 Access time for conventional 10T CAM and enhanced pre-comparison CAM



4.3.3 Power Delay Product

Power delay product is a pointer to see if the design is better. So, the power delay is the product by Fig. 4.11 and Fig. 4.14. Using the power delay product value of conventional 10T CAM can be basis. Availing the value compares the other improved pre-comparison CAM to see the result. Fig. 4.15 is the result of power delay product. Because the improved pre-comparison CAM has two kinds of access time, this figure also has two different data. For the longer access time, its power delay product is also under the value of the conventional 10T CAM, and only the five bits improved pre-comparison CAM is greater than 10T CAM. For the shorter access time, its power delay is actually shorter than the conventional 10T CAM. So, except for five bits and more than five bits improved pre-comparison CAM, the power delay product would be smaller than the conventional 10T CAM. So, my low power design is suitable for one to four bits from the viewpoint of power delay product.

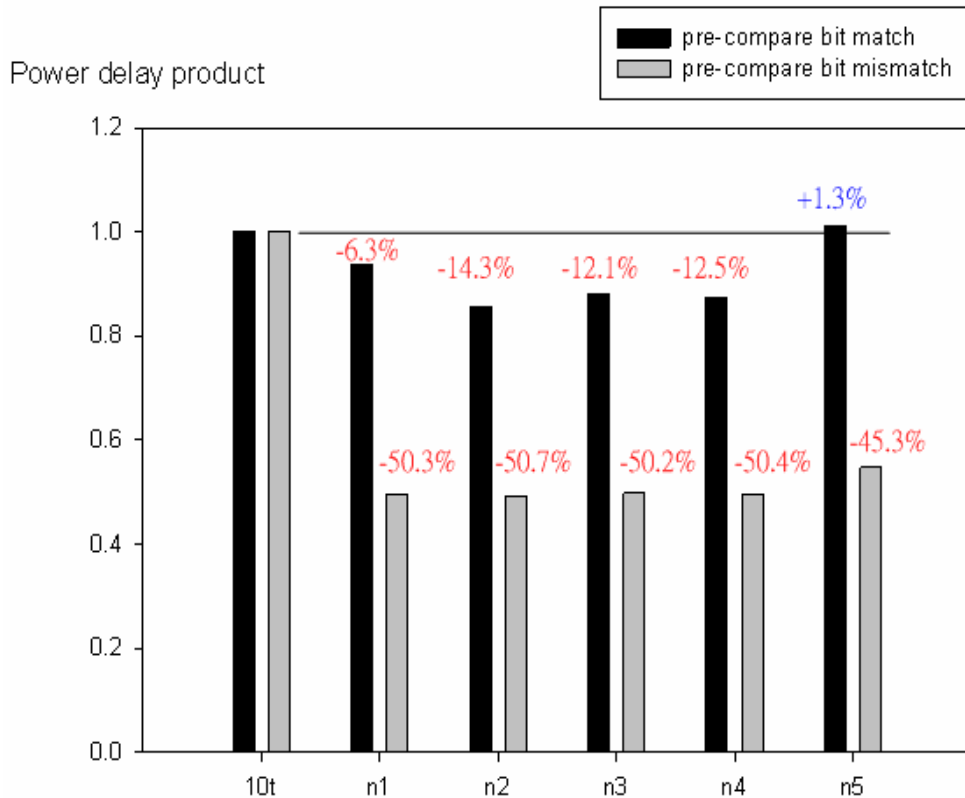


Fig. 4.15 Power delay product for conventional and enhancement of pre-comparison CAM



4.3.4 Transistor Penalty for Pre-comparison Circuit

As the bit of improved pre-comparison CAM is more, the match line capacitance is smaller because of the reduction of XOR logic gate. But there is another problem for more pre-comparison bit. That is the transistor penalty and area penalty for the pre-comparison circuit. If the pre-comparison bit increases, the logic gate would increase to do the pre-comparison. We will do some calculation for the relation of pre-comparison bit and additional transistors. One bit improved pre-comparison CAM needs additional transistors is as follow:

1bit:

$$(4 \cdot 3) + 2 + 2 + 4 - 4 = 16 \text{ transistors} \quad (4.7)$$

2bits:

$$(4 \cdot 3 \cdot 2 + 4) + 2 + 2 + 4 - 4 \cdot 2 = 28 \text{ transistors} \quad (4.8)$$

nbits:

$$\begin{aligned} & 3 \cdot (\text{two input NAND gate}) \cdot n + (n \text{ input NAND gate}) + (\text{one inverter}) + (\text{transmission gate}) \\ & + (\text{one two input NOR gate for output logic}) - (\text{reduced XOR gate}) \cdot n \\ & = 3 \cdot 4 \cdot n + 2 \cdot n + 2 + 2 + 4 - 4 \cdot n \\ & = 10n + 8 \end{aligned} \quad (4.9) \quad (\text{for } n > 1)$$

From the above count of additional transistor, we can find out some regulation. As the bit number increases one bit, it would needs more ten transistors. Through the regularion, a formula 4.9 is produced. Fig 4.16 shows the corresponding additional transistors for different bit number of improved pre-comparison CAM.

numbers of additional transistor

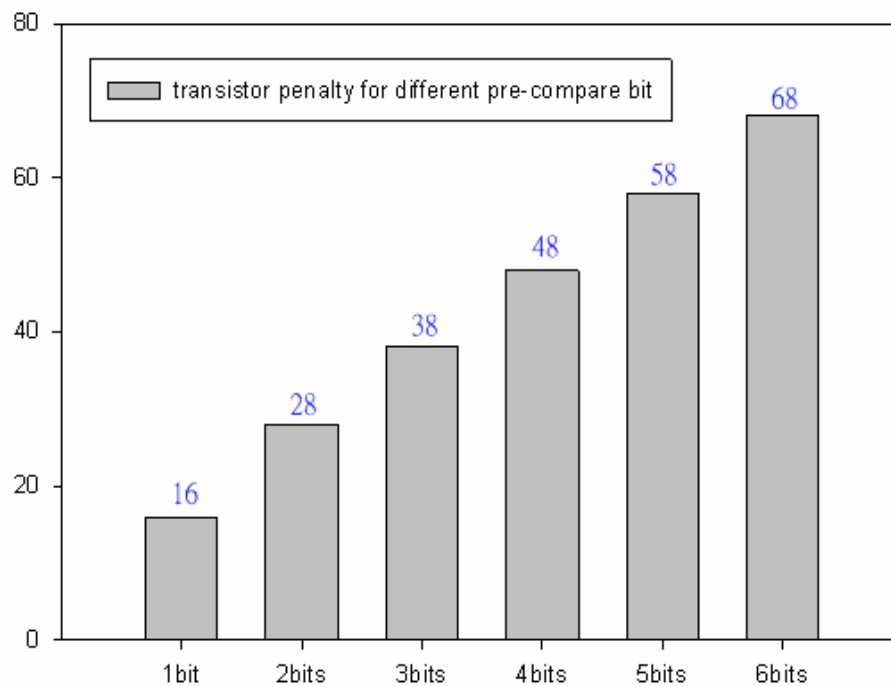


Fig. 4.16 The transistor penalty for different bit number

4.3.5 The Optimal Bits for Pre-comparison Circuit

From the above text, we make some discussion about improved pre-comparison CAM. The consideration has many aspects, such as match line capacitance, prediction precision, power consumption, access time, additional transistors penalty, power delay product, and so forth. For match line capacitance, the bit is greater, the capacitance is smaller. Prediction precision is the same. But its bit number goes to three, the prediction precision has reach almost ninety percentage. As the pre-comparison bit increases, the precision increases more and more slow.

For the viewpoint of power consumption, the bit number that is two or three or four bits has lower power dissipation. The access time would increase when the pre-comparison bit increases. The additional transistor penalty is less for less pre-comparison bit. For power delay product, the three bits and four bits are both good. From the above description, the optimal bit numbers for improved pre-comparison CAM is between three bits and four bits. The power consumption and power delay product are the main consideration. For three and four pre-comparison bits, their power consumption and power delay product are almost the same and the minimum value, so we get the conclusion that three bits to four bits is the optimal bit numbers for enhanced pre-comparison CAM array.

4.4 Layout

Fig. 4.17 is the layout for one bit CAM cell. Fig. 4.18 has two kinds of CAM for one word line and thirty-two bits. First one is conventional NOR type 10T CAM, and the other one is four bits pre-comparison NOR type 10T CAM. The pre-comparison CAM has the area overhead about the pre-comparison circuit. Fig. 4.19 is a four bits CAM array. The array size has thirty-two words and thirty-two bits. Table 4.2 shows the area penalty for different bits of pre-comparison CAM. For four bits pre-comparison CAM, its area penalty is 28.1% than than the conventional 10T CAM.

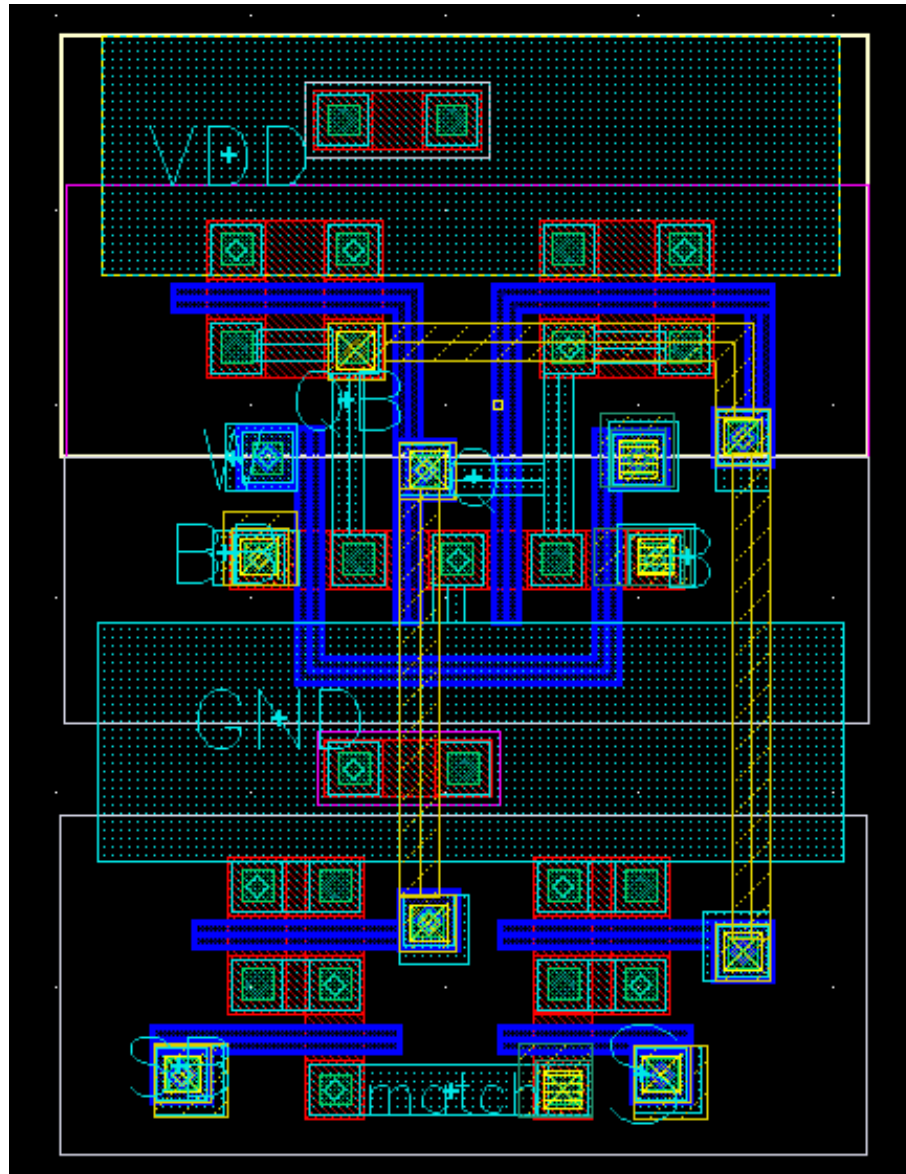


Fig. 4.17 Layout for 1bit 10T CAM cell

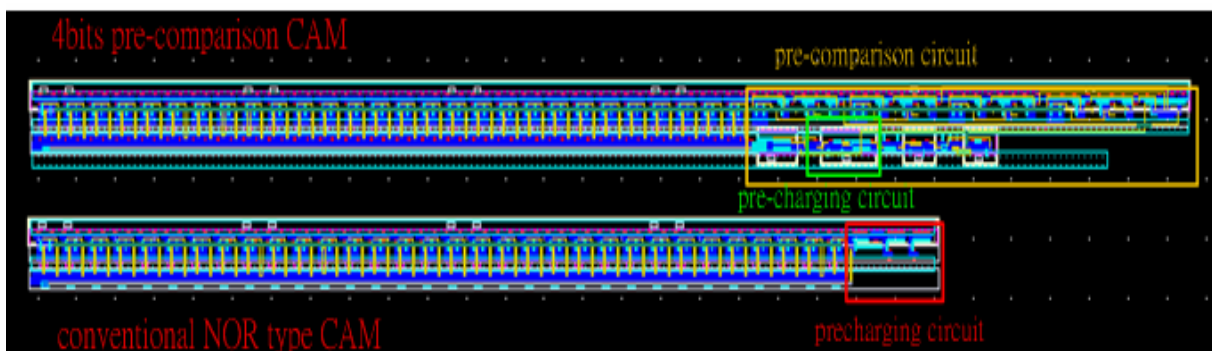


Fig. 4.18 Layout for 32 bits conventional 10T CAM and 4 bits pre-comparison CAM

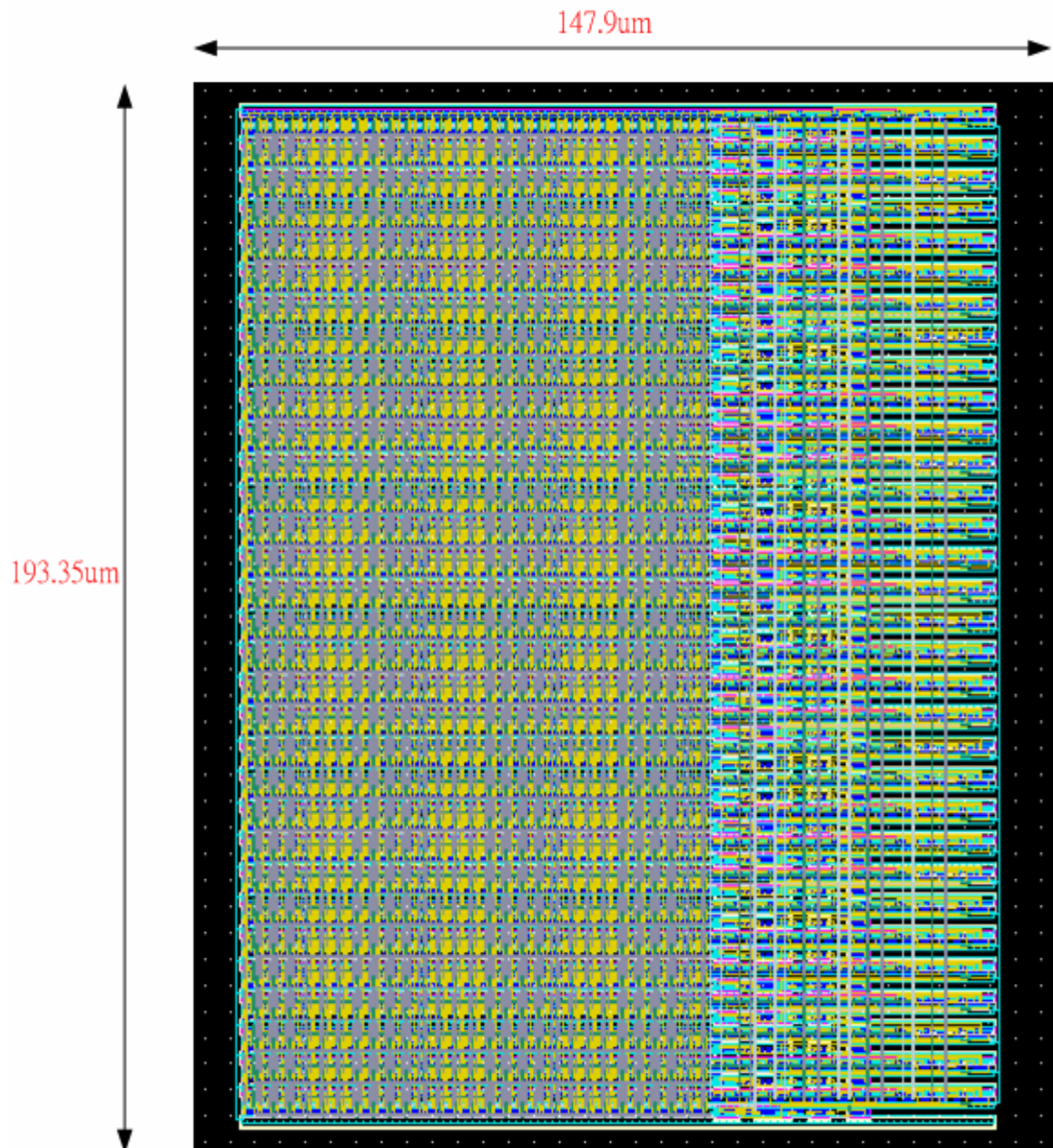


Fig. 4.19 Layout for four bits CAM array with 32words x 32bits

Table 4.2 Area penalty for pre-comparison CAM

	conventional	2bits	3bits	4bits	5bits
Area(μm^2)	22370.595	24437.5065	26633.9625	28656.4035	30678.8445
overhead	0%	9.2%	19.1%	28.1%	37.1%

4.5 Conclusions

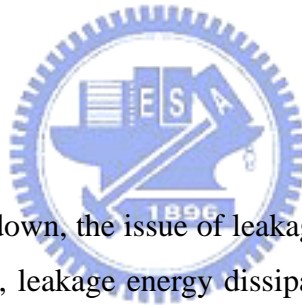
Pre-comparison CAM has improved in this chapter. The improved pre-comparison CAM has lower match line capacitance and lower power consumption. It has two kinds of access time, one is shorter and the other is longer. Although, the longer one is larger than the access time of conventional CAM, its power delay product is still smaller than the conventional one for one pre-comparison bit to four pre-comparison bits. For four bits improved pre-comparison CAM, it reduces 22.8% power consumption, increases 13.4% access time, lower 35.7% access time, and decrease 12.5% and 50.4% power delay product. The area overhead of four bits pre-comparison CAM is 28.1% for additional pre-comparison circuits.

The pipeline technique may be used in the pre-comparison CAM. If the pre-comparison CAM has one register, the pre-comparison result can be saved. It will have two kinds of condition. When the pre-comparison is mismatch, the register is transparent. It means that the cycle time is one. If the pre-comparison is match, the register will store the pre-comparison result. It needs two cycles. The pre-comparison result is used to decide one cycle or two. For a CAM array, it may have many word lines. Only when all the pre-comparison results of all words are mismatch, the cycle of the system would be one. Otherwise, the cycle time is two. With pipeline, the average cycle time is between one and two. Because the percentage of mismatch for all word lines is small, the cycle time will be close to two. Its efficiency is worse. Therefore, the pipeline technique is not appropriate to the pre-comparison CAM.

Chapter 5 Application of Dual V_{dd} and Power Gating for CAM and TLB

For deep submicron, the leakage current is a serious problem. Leakage current occurs in sleep mode or standby mode for logic circuit. Sleep mode or standby means that the logic gate is not in action state and has no data variation. So, memory devices, such as CAM, SRAM, and TLB, are affected by leakage current obviously. This is because memory devices are stored elements and most of time for memory devices do not work and only hold data. In the hold data state, the leakage current would appear to cause additional power dissipation. So leakage current must be put attention and can not be ignored. The following section would discuss the leakage problem.

5.1 Leakage current



As the technology scales down, the issue of leakage current is more and more serious. In deep submicron CMOS design, leakage energy dissipation occupies a large portion for total power dissipation. Chip designers scale down the supply voltage to reduce the dynamic energy dissipation, but it also scales down the threshold voltage that would make the leakage current rise, even the transistors do not activate. Fig. 5.1 illustrates the magnitude of the problem with data from existing technologies and projections based on the international technology roadmap for semiconductor (ITRS). As it can be seen, Fig. 5.1, even in current-generation technology, subthreshold leakage power dissipation is comparable to the dynamic power dissipation. The fraction of the leakage power will increase significantly in the near future. In fact, the off-state subthreshold leakage component of the total power in a microprocessor may exceed active power as the technology decreases [38]-[40].

Subthreshold leakage is a problem for all transistors logic circuit. For on-chip caches, they are growing fraction of the total number of microprocessor devices. Furthermore, the leakage power is becoming the dominant fraction of total power consumption because of many memory devices. Most of data in Caches is accessed relatively infrequently due to

either temporal or spatial locality, thus, they will increase the leakage power. The leakage power dissipation for memory has very large portion. How to reduce the leakage current for deep submicron is a very important issue.

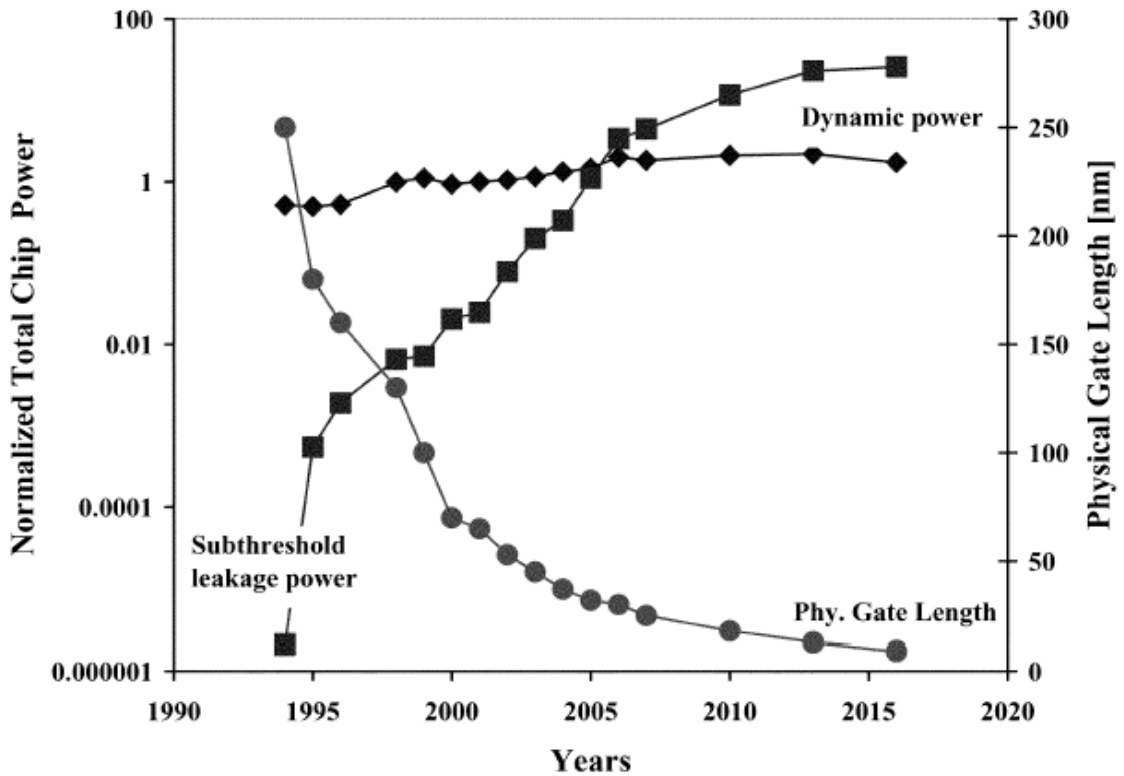


Fig. 5.1 Normalized dynamic and static power dissipation

5.2 Low leakage design

There are many methods to reduce the leakage power, such as gating vdd or called power gating, multi-threshold voltage, dynamic voltage scaling, and so on. The following section would introduce those methods.

5.2.1 MTCMOS

Multi-threshold CMOS (MTCMOS) circuit is a technical for low leakage design. It divides several different modes for logic circuit. In general, it has action mode and sleep mode. MTCMOS logic circuit is as Fig. 5.2. Action mode means that the logic circuit is working.

When MTCMOS logic circuit is activated, the nMOS that is connected between ground and virtual ground and the pMOS that is connected between vdd and virtual vdd would be turned on. The action mode would let the virtual ground pull down to ground and the virtual vdd pull up to vdd. On the other hand, sleep mode means that the logic circuit is not to work. So it is called sleep mode. When the logic circuit is in sleep mode, the logic circuit would not work. So, the gated nMOS and pMOS won't be turned on. Because the gated MOS turns off, the leakage current would be reduced. The leakage current can be reduced [41]-[42], [44], [47].

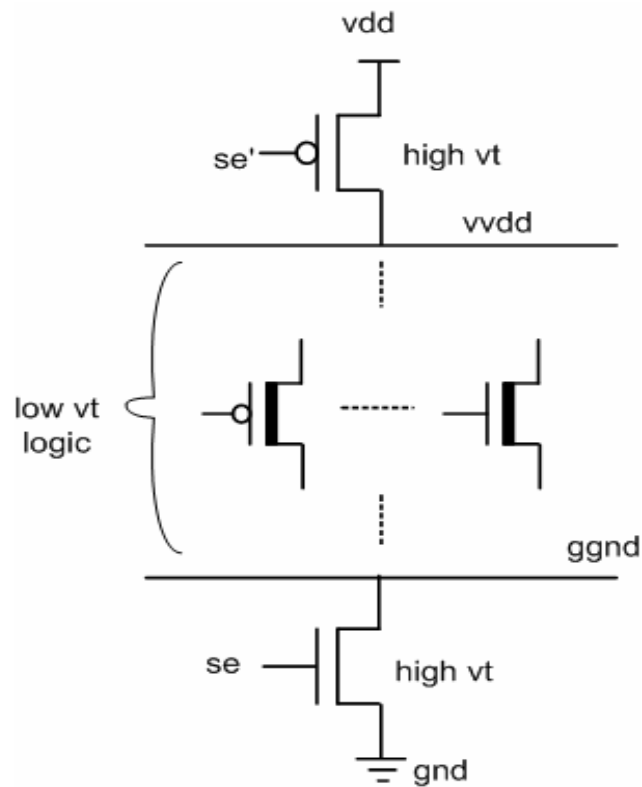


Fig. 5.2 Conventional MTCMOS

From Fig. 5.2, we see two different threshold voltages. One is high threshold voltage that is abbreviated high V_t and the other is low threshold voltage that is called low V_t . The characteristic for high V_t and low V_t is the speed and leakage current. To increase operating speed, low V_t MOS is used for logic gates. During the long standby time, the power supply is disconnected with high V_t MOS. The high V_t has lower leakage current, but its speed is slower and worse performance. On the contrary, the low V_t has faster speed, but its leakage current is very seriously. So the gated MOS, nMOS and pMOS, are used high V_t transistors and the logic gate circuit are used low V_t transistors, such as Fig. 5.2. This method is called

multi-threshold CMOS circuit and is abbreviated MTCMOS.

5.2.1.1 Selective MTCMOS

Section 5.1.1 has introduced the basic MTCMOS design. Here, we would introduce Selective MTCMOS. Selective MTCMOS is in Fig. 5.3. It is NAND cell and NOR cell. The main part is composed of low V_t transistor and the gated MOS is composed of high V_t transistor. The signal, MTE, is used to control the gated MOS to turn on or turn off. When MTE is logic one, it is in the active mode and results in performing fast logic operation with low V_t transistors. If it is in sleep mode or called standby mode, the MTE would be set logic zero to cut off the subthreshold leakage path from vdd to ground [43].

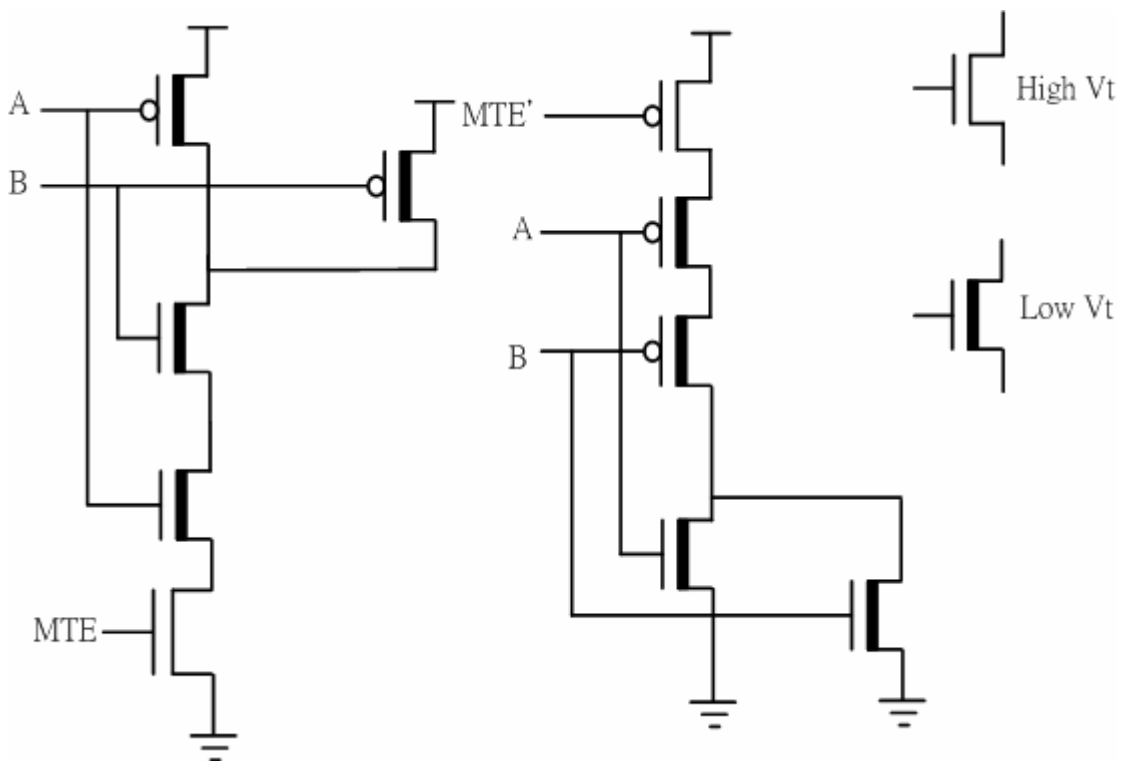


Fig. 5.3 Selective MTCOM

In the logic circuit design, there are some critical path and some non-critical path. Critical path means that it is longest sensitized path. The delay of the logic circuit is equal to the delay of the critical path. So, critical path dominates the performance of the logic circuit.

How to reduce the delay of critical path is an important issue. Fig. 5.4 is an example for MTCMOS design. The circuit is composed of many logic gates, such as and gate, or gate, and etc. The critical path of the circuit is the blue color in the figure. The critical path of the circuit is composed of Selective MTCMOS. The non-critical path of the circuit is composed of high V_t , because the performance of non-critical path is not the most important thing and using high V_t can reduce the leakage current and leakage power consumption.

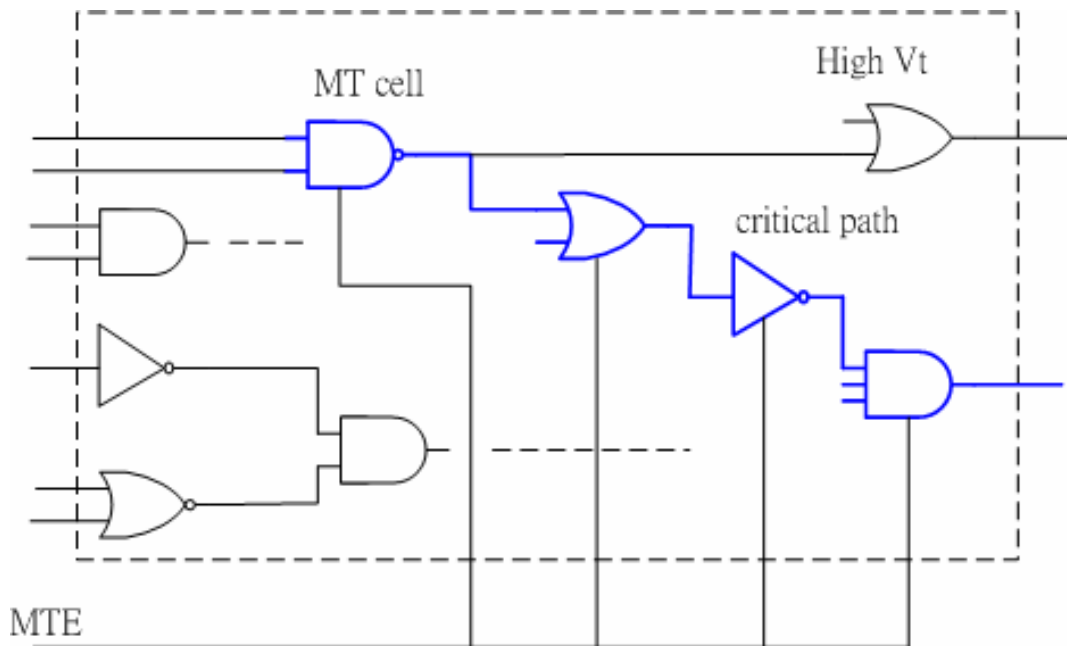


Fig. 5.4 Selective-MT circuit

5.2.1.2 Improved Selective MTCMOS

We have introduced Selective-MT circuit. But it has one drawback is its area. The MT logic gate must add additional gated MOS to control the leakage path. So there is some area overhead. In order to reducing the area overhead, there is another paper to provide new thinking, in Fig. 5.5 [45].

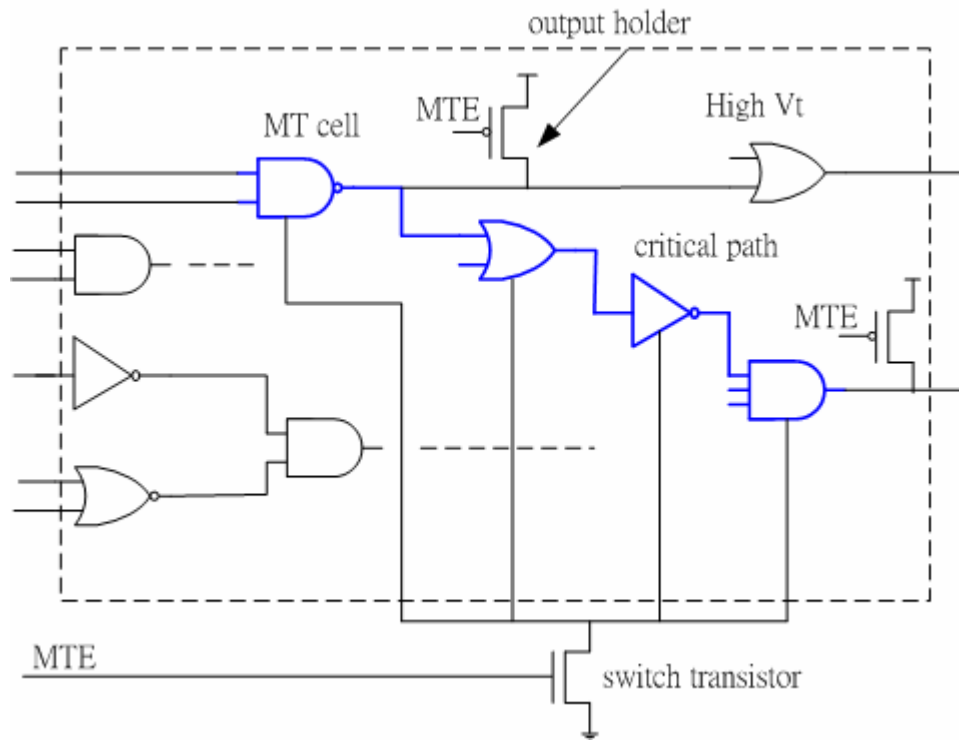


Fig. 5.5 Improved selective-MT circuit



5.2.3 Gated Vdd

Before we introduce the technique of gated vdd, stacking effect would be introduced firstly. Fig. 5.6 can explain this phenomenon, stacking effect. On the left of figure is an off nMOS transistor and it has a leakage current in the steady state. On the right of figure have two nMOS transistors that are stacked and off.

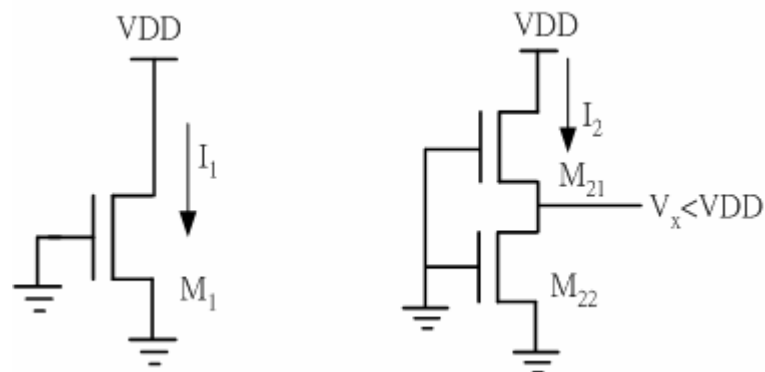


Fig. 5.6 Stacking effect due to self-reverse biasing

The node of V_x is slightly higher than ground in the steady state, so the transistor M_{21} has a negative V_{gs} to make the pn junction reversely biased. This phenomenon would make the leakage current I_2 is smaller than I_1 due to the reversed bias of transistor M_{21} . By stacking effect, the leakage current can be reduced [46].

Although stacking effect can reduce the leakage current, it also has some tradeoff. The tradeoff is the time delay. We can see Fig.5.7. It is a comparison between conventional MOS and stacked MOS for time delay and leakage current. We can find out that the difference between conventional MOS and stacked MOS. The time delay of stacked MOS is longer than conventional MOS by three times. On the other hand, the leakage current of stacked MOS can save almost ten times leakage. We only talk about one additional MOS for stacking effect. If the stacked MOS is more, the reduced leakage would be more, but the delay would be worse. So, the gated vdd is used with one gated nMOS in general.

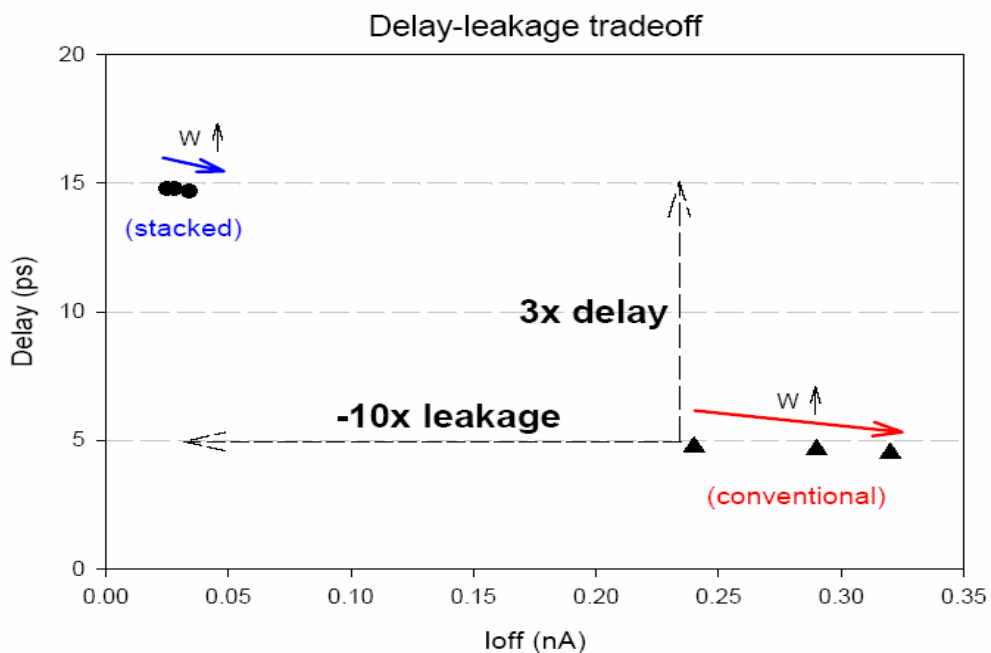


Fig. 5.7 Delay-leakage tradeoff of stacking effect

Fig. 5.8 is used the technique of gated vdd to SRAM. Figure (a) is without diode. Figure (b) is with diode. We know that the gated vdd has additional nMOS to make the original ground to virtual ground. In the active mode, the nMOS is turned on to activate. In the standby mode, the nMOS is turned off to reduce the leakage current. A control signal controls nMOS to turn on or turn off.

For figure (a), there would be a problem in the virtual ground, $vss0$. If the standby mode lasts some long time, the small leakage would be continue to charge the node of $vss0$. This condition would make the node of $vss0$ go to higher. If the voltage of $vss0$ is charged to half of vdd , then the original state that is stored in the memory device would be destroyed. Additional diode, figure (b), is added in the virtual ground to make the virtual ground, $vss1$, keep a steady voltage value. The original state in the memory device would not be destroyed and has the same function to reduce the leakage current.

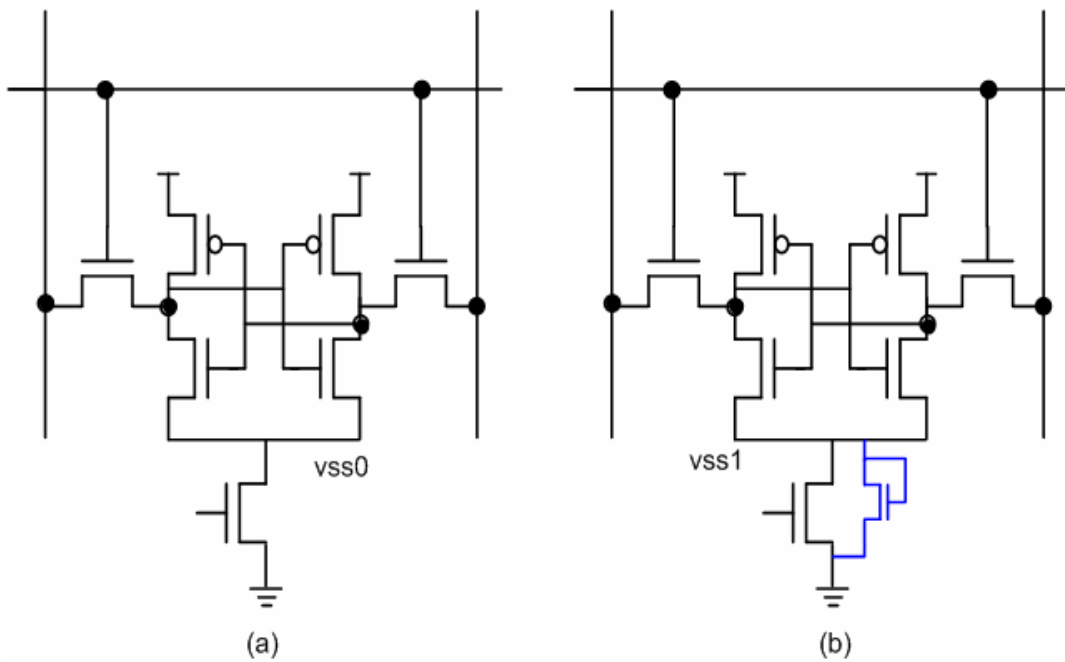


Fig. 5.8 Gated vdd SRAM (a) without diode (b) with diode

Fig. 5.9 shows the two kinds of gated vdd comparison for Fig. 5.8 (a) and (b). The virtual ground without diode, $vss0$ in Fig. 5.9, is obviously to be charged beyond the half of vdd . It will destroy the stored data. The virtual ground with diode, $vss1$ in Fig. 5.9, is also to keep the stable value, 100mv, and would not affect the value that is stored in memory device. So the gated vdd with diode is feasible to reduce leakage current.

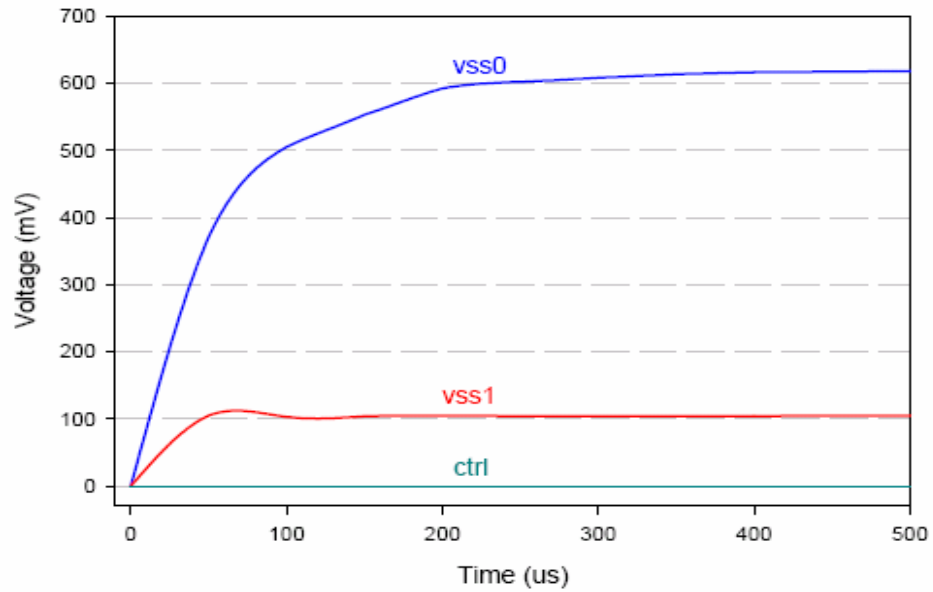


Fig. 5.9 Voltage of virtual GND increases after turning off gating device

For memory, the static noise margin (SNM) is important. The value of static noise margin can show how sensitive the memory device is. If the value of static noise margin is bigger, it can endure more noise disturbances. Fig. 5.10 (a) shows a latch that comprises two inverters and two static noise sources, and Fig. 5.10 (b) shows the graphical view of SNM. Fig. 5.11 (a) shows a latch that comprises two inverters with gated vdd and two static noise sources, and Fig. 5.11 (b) shows scale factor, n , for all the MOS size.

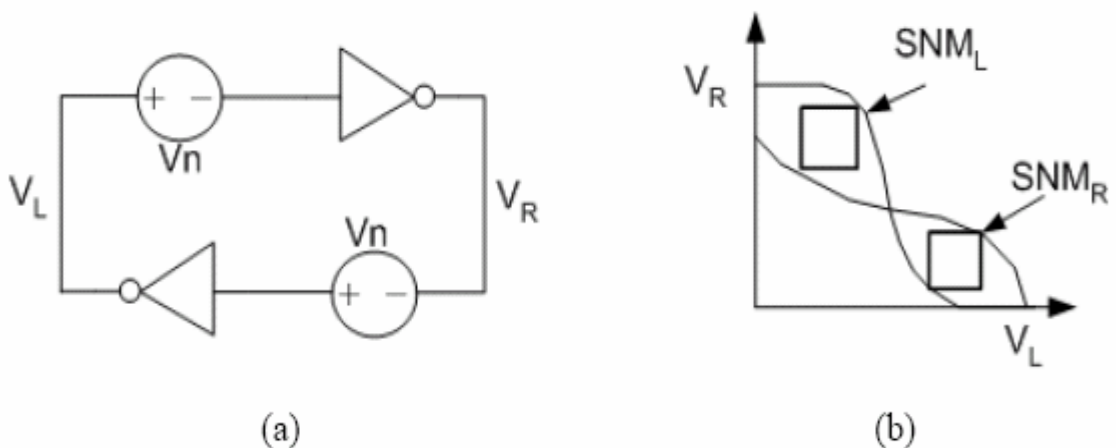


Fig. 5.10 (a) A latch with static noise sources and (b) Static noise margin

In Fig. 5.11 (b), the nMOS and pMOS of SRAM have the same size and are n times for the MOS of gated vdd and diode. Fig. 5.12 is the simulation result about the relation of SNM, power consumption, and the scale factor n . As the scale factor increases, the SNM and power consumption would both increase.

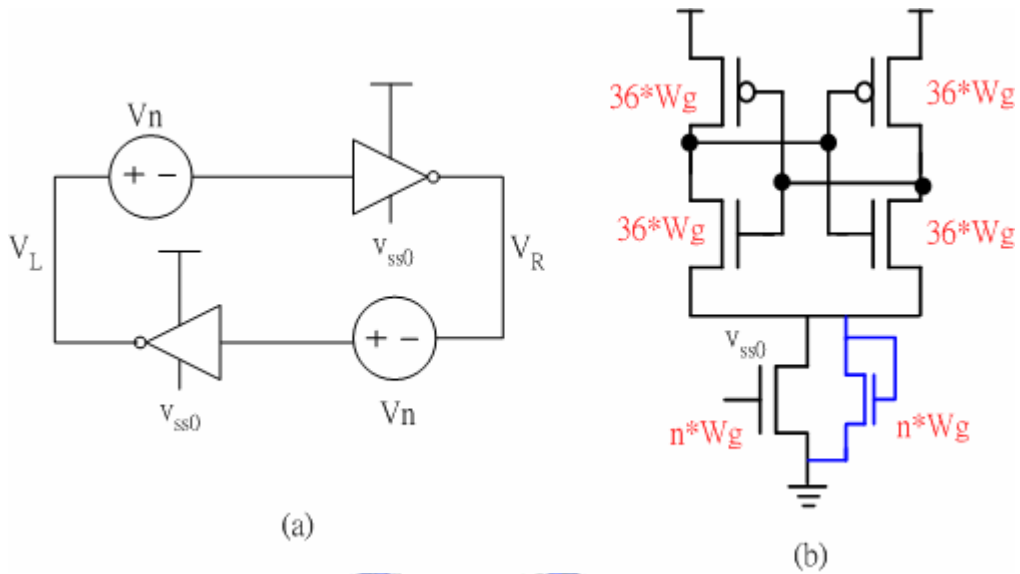


Fig. 5.11 (a) Two gated vdd inverters and two static noise sources (b) Scale factor of gated vdd for SRAM

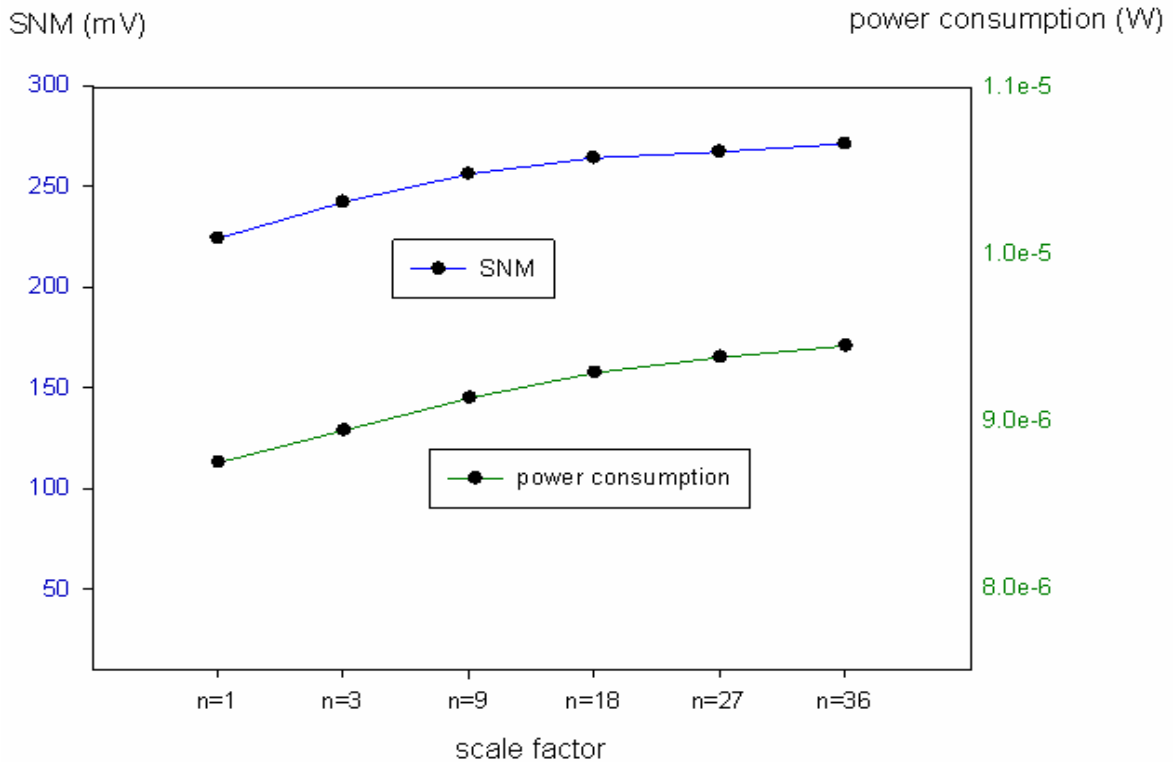


Fig. 5.12 SNM and leakage power versus scale factor n

5.2.3 Dual vdd

The technique of dual vdd is like MTCMOS. Dual vdd also has two kinds of supply voltage, V_{DDH} and V_{DDL} . V_{DDH} is the original supply voltage, and V_{DDL} is the lower supply voltage. The technique of dual vdd can reduce much dynamic power dissipation because dynamic power dissipation is proportional to the square of supply voltage. Its thinking is similar with MTCMOS. If the circuit is non-critical path, it can use the lower supply voltage. Of course, the lower supply voltage would degenerate the performance. If the circuit is critical path, it must use the original supply voltage [48].

For memory, the dual vdd is not used to reduce dynamic power dissipation. It is used to reduce the leakage power. The principle is similar with MTCMOS. On one hand, performance demands require the use of fast high-leakage transistors. On the other hand, energy efficiency requires low-leakage transistors. If the memory cells are not intended to be accessed for a time period, they would be placed in a sleep mode or standby mode by supplying a standby voltage to the memory cells. The leakage power is significantly reduced due to the decreases in both leakage current and supply voltage. Supply voltage reduction is especially effective for leakage power reduce because of short-channel effects, such as drain-induced barrier lowering (DIBL). DIBL results in a superlinear dependence of leakage current on the supply voltage.

5.3 Low leakage SRAM

SRAM is used very widespread and indispensable for storing data. However, memory device element has the common problem. That problem is the leakage power. In the above text, the stacking effect can reduce the leakage current. Some papers use this concept to make the SRAM have virtual ground. Using a signal controls it to turn on or turn off in active mode or sleep mode [46]-[47], [49]-[50].

5.3.1 Power gating SRAM

Fig. 5.13 is the method of gated vdd, or called power gating, used in a word line of SRAM. The power gating transistors are composed of two nMOS. One nMOS is used to be a diode. A diode is composed of nMOS, and its gate node and drain node are connected together to the virtual ground. The other nMOS is control by a signal. The signal can be the same of word line signal. Because only the word line is turned on, the SRAM would be read out or written in data, and the virtual ground must go to the ground. On the other time, the SRAM is on standby mode and the controlled nMOS would be turn off to reduce the leakage current.

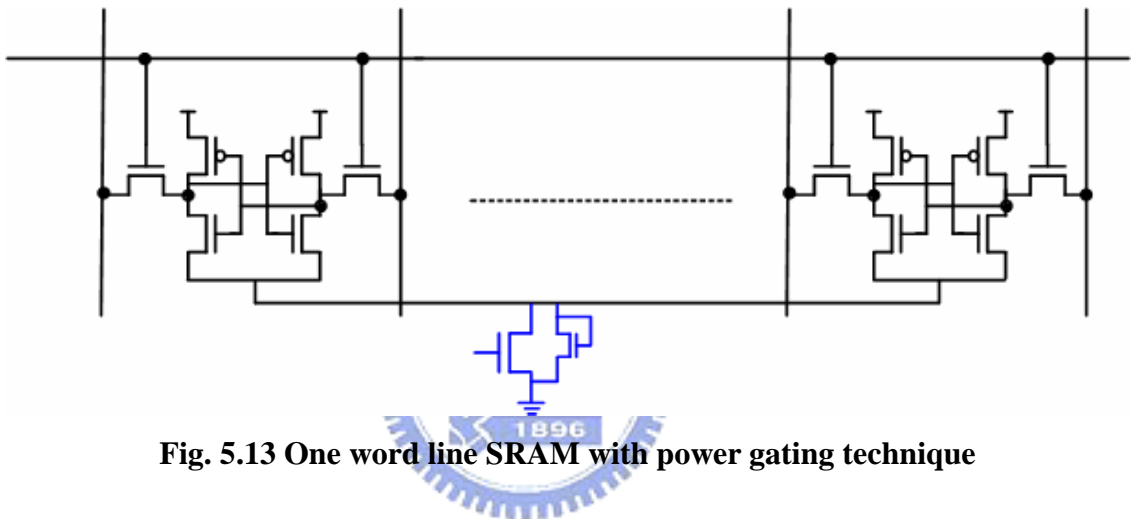


Fig. 5.13 One word line SRAM with power gating technique

5.3.2 Simulation

We simulate some cases about power gating SRAM to compare its characteristic. The technology use TSMC 100nm CMOS technology. Use one word line SRAM and have thirty-six bits to compare the scale factor and standby power. We use the conventional SRAM and the different scale factor of power gating SRAM. Fig. 5.14 is the power consumption comparison and the access time comparison.

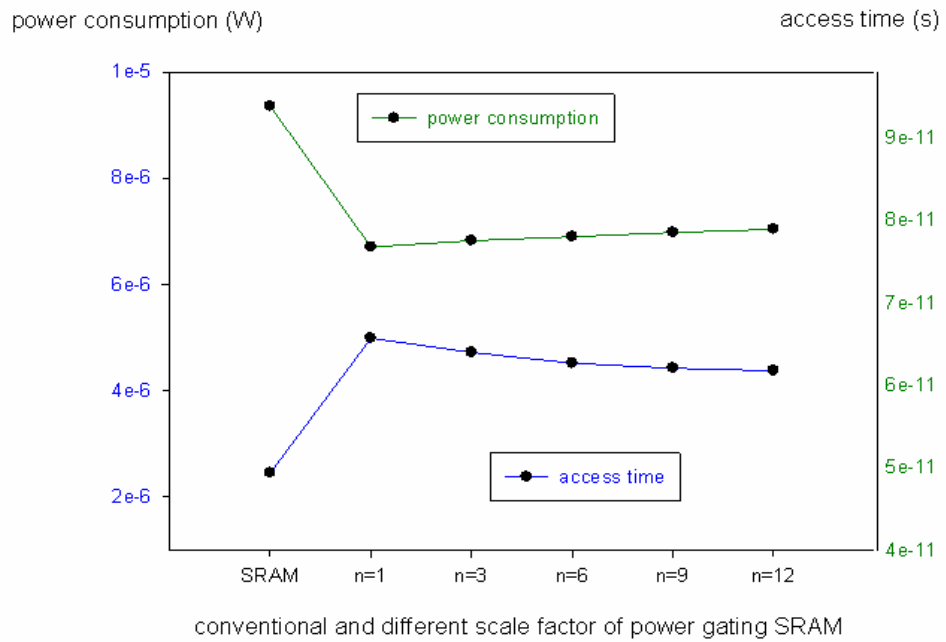


Fig. 5.14 Power and timing for scale factor of power gating SRAM

As the scale factor is small, the reduced leakage power is more. The scale factor increases bigger, the percentage of reduced leakage would be decreased. But there is a tradeoff for power gating SRAM. If the scale factor is smaller, it can save more power. The cost is its access time because the smaller scale factor makes the access time longer. Power delay product is also compared. For different scale factors, their power delay products are almost the same in Fig. 5.15.

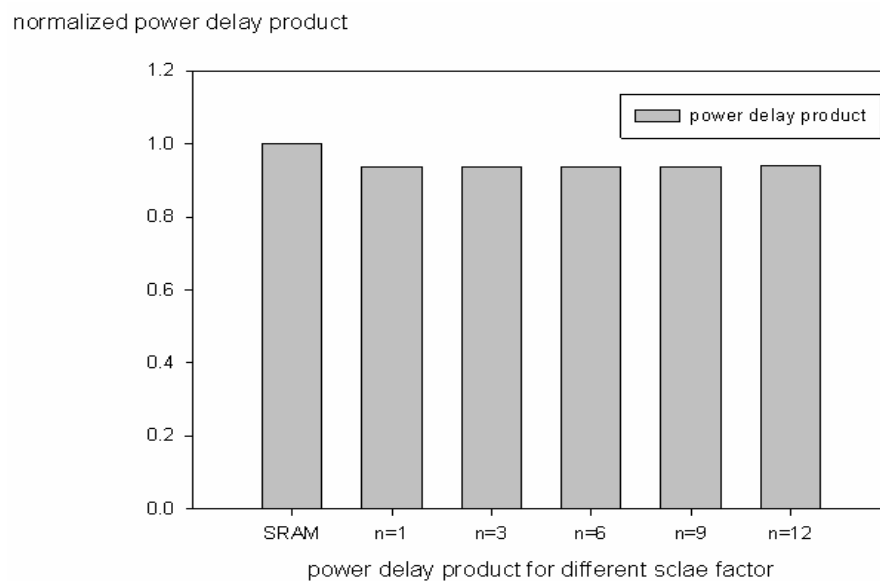


Fig. 5.15 Power delay product for different scale factor of power gating SRAM

Fig. 5.16 is the simulation for different dual voltage value. The lower dual voltage can have the lower leakage power, but the access time for different dual vdd is the contrary to the leakage power. The lower dual vdd makes the access time longer. Fig. 5.17 also compares the power delay product. We can observe that the power delay product for different dual vdd. The lowest power delay product is the value of 0.5V or 0.6V. If we choose the value of 0.5V or 0.6V, it would have some problems. The value of 0.5V or 0.6V is only half of original supply voltage. If any noise comes, the original logic one, would be smaller than half of supply voltage in standby mode. When it comes back in active mode, the original logic one would be viewed as logic zero because of lowering half of supply voltage. So we take the value of 0.8V for dual vdd to prevent this condition.

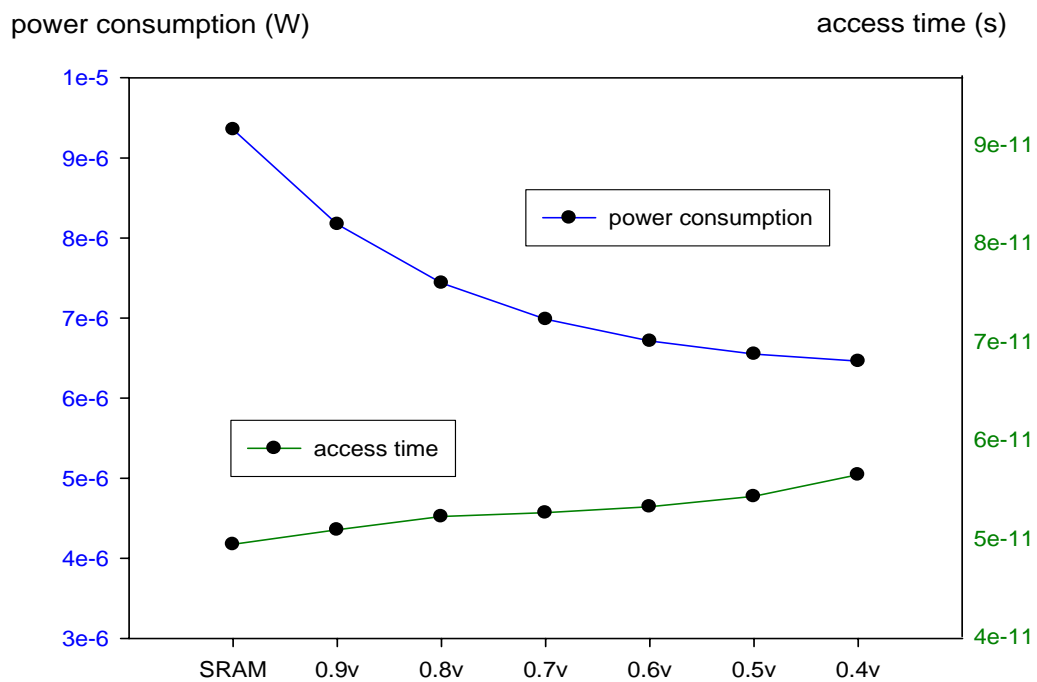


Fig. 5.16 Comparison of power consumption and access time for dual vdd SRAM

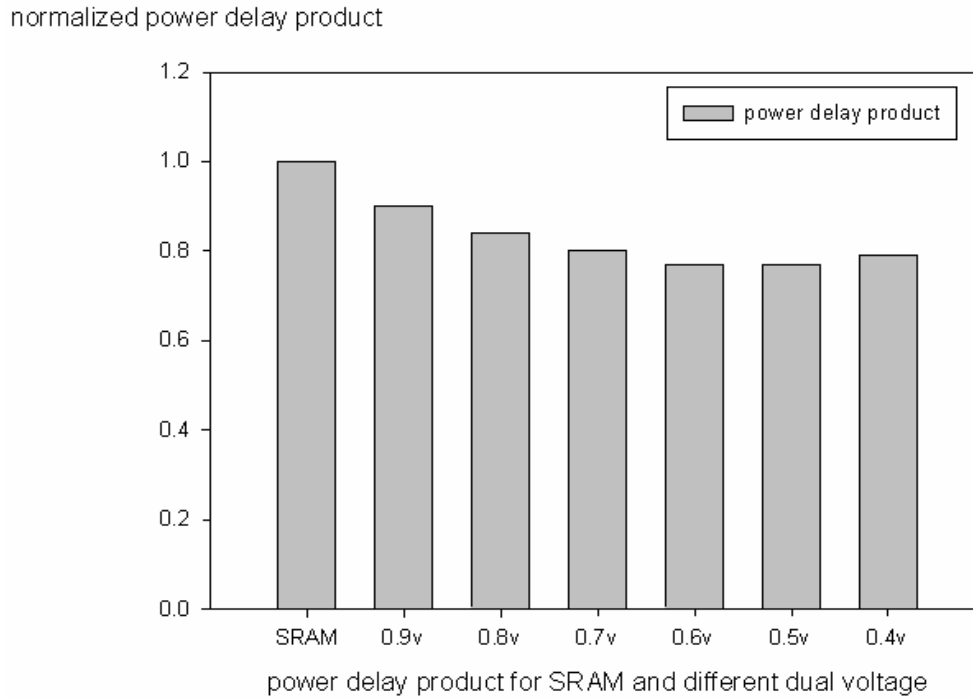


Fig. 5.17 Power delay product for different scale factor of dual vdd SRAM

5.3.3 SRAM Combines Power gating and dual vdd

In this section, we combine the technique of power gating and dual vdd technique to SRAM, such as Fig. 5.18. We also simulate it to observe the power consumption, access time, power delay product, and compare these conditions for conventional SRAM. Fig. 5.19 is the comparison result for leakage power and access time for conventional SRAM and power gating and dual vdd SRAM. We use one word line and thirty-six bits to simulate. We can find out that the lower vdd can have a lot of power saving. But the access time will be greater than the conventional SRAM. For fixed size power gating and six kinds of dual vdd SRAM, the 0.6v one has smaller access time. The simulation is based on the smallest size of power gating device, and we also do the comparison for power delay product in Fig. 5.20.

We also let the size of power gating be variation and fixed the dual vdd is 1v and 0.8v. We also do three kinds of simulation, leakage power consumption, access time, as well as power delay product. The simulation result is in Fig. 5.21 and Fig. 5.22.

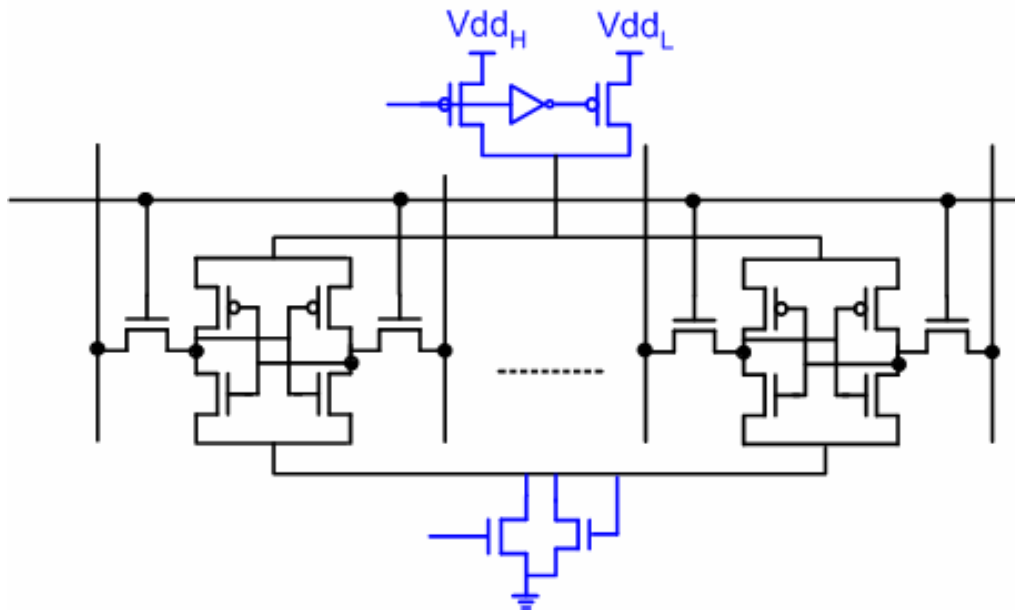
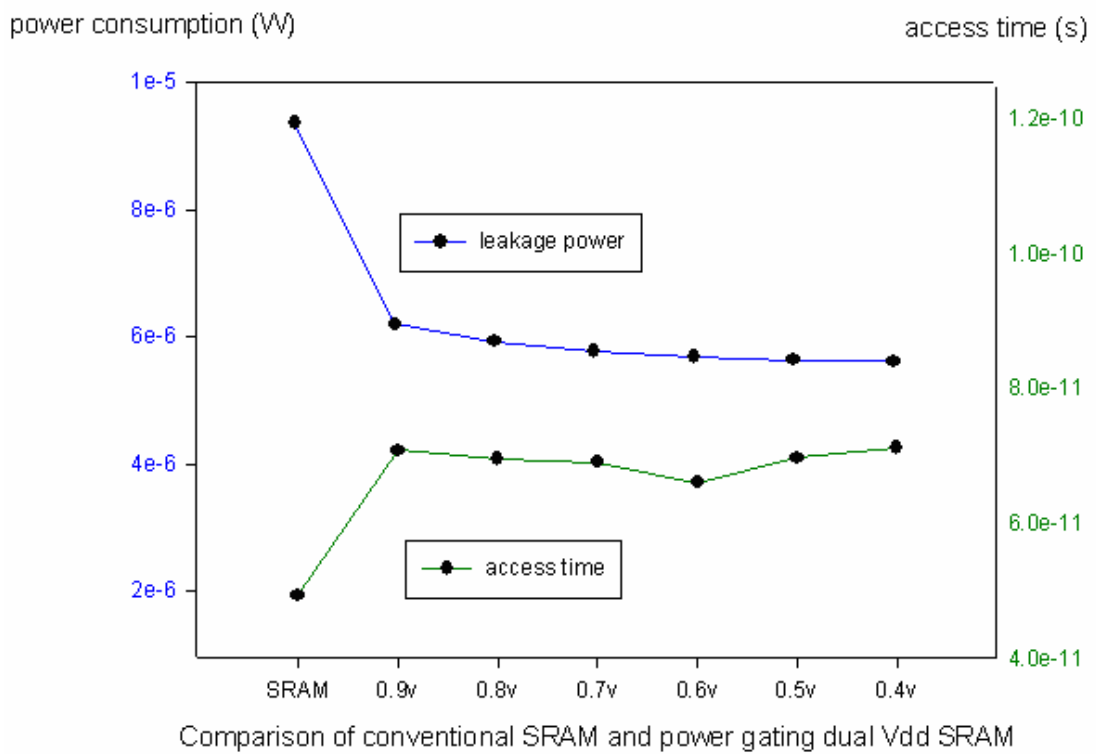


Fig. 5.18 Combine power gating and dual vdd SRAM



Comparison of conventional SRAM and power gating dual Vdd SRAM

Fig. 5.19 Power consumption and access time for power gating dual vdd SRAM

normalized power delay product

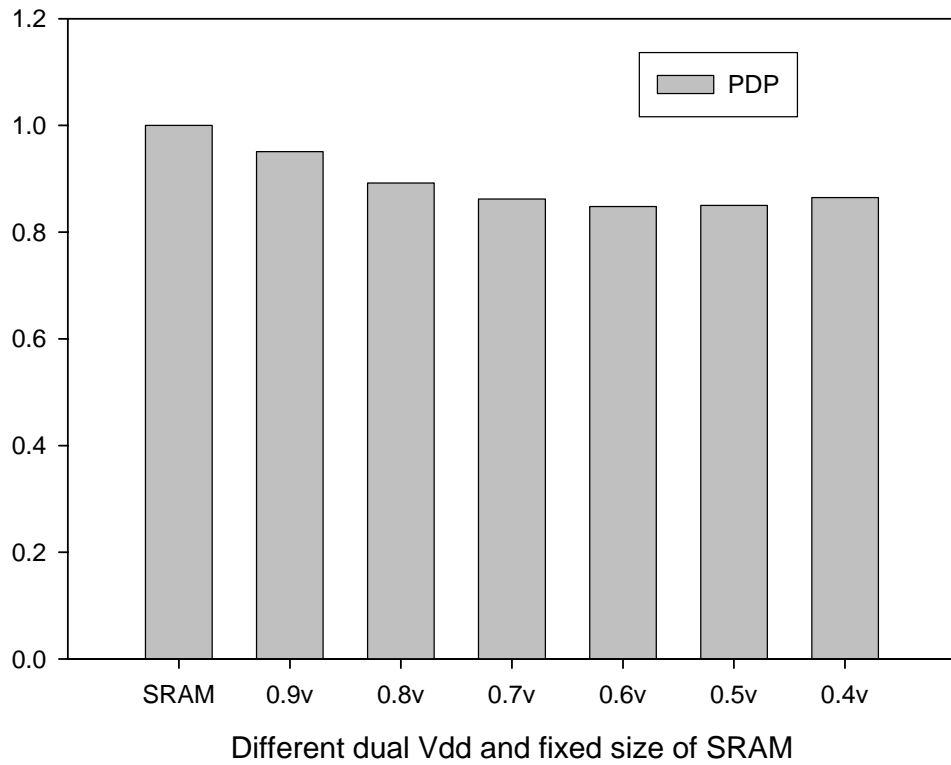
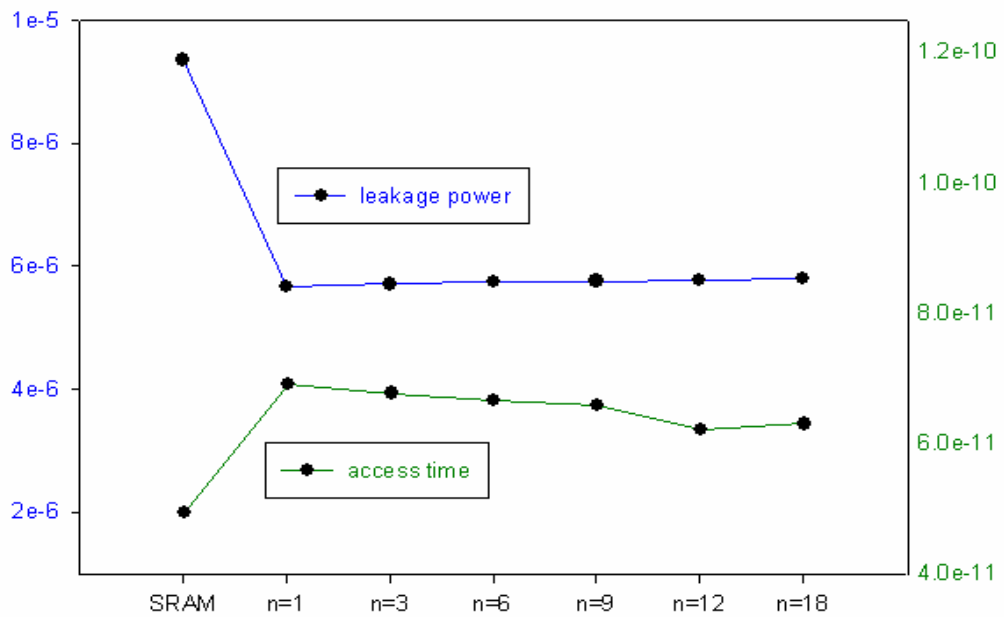


Fig. 5.20 Power delay product for combining dual vdd and fixed size of power gating SRAM



power consumption (W)

access time (s)



Fixed dual Vdd, 1v and 0.6v, and change the scale factor for power gating SRAM

Fig. 5.21 Power consumption and access time for power gating SRAM

normalized power delay product

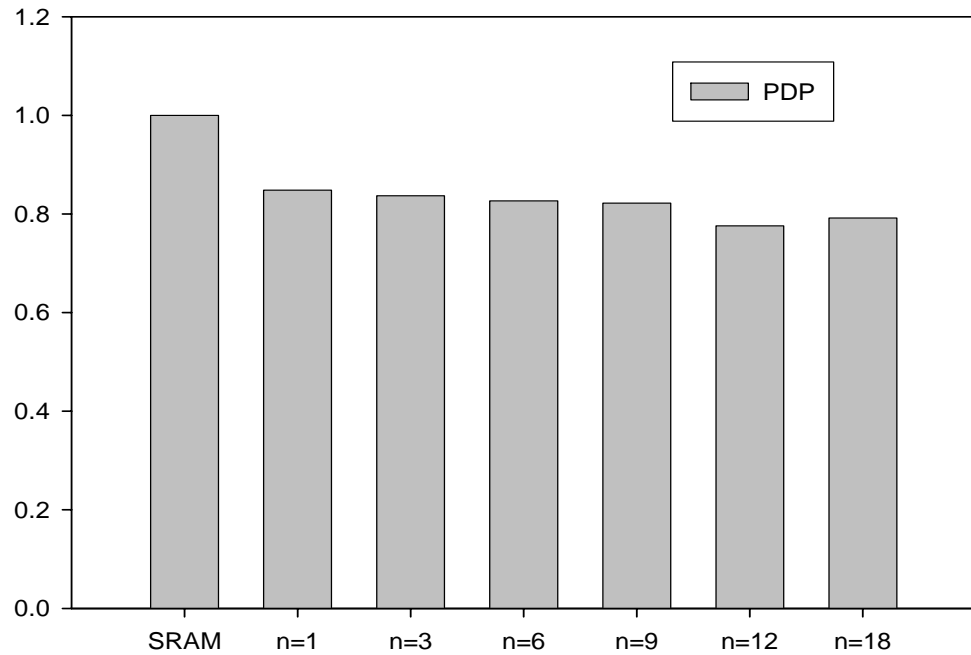


Fig. 5.22 Power delay product for power gating SRAM



5.4 Simulation for Pre-comparison CAM

5.4.1 Simulation of Dynamic Power for Pre-comparison CAM

Fig. 5.23 is the simulation result for dynamic power consumption of CAM array. we use the same size 36x36 and compare three kinds of CAM, conventional 10T, 3bits pre-comparison, and 4bits pre-comparison CAM. The CMOS technology is 100nm. From the simulation result, the reduced power consumption is not as much as the 0.13um. For 4bits pre-comparison CAM, it can reduce 22.8% in 0.13um, but it only reduces 9.6% in 100nm. This is because that the technology scales down and the effect of dynamic power is not so important. On the contrary, the leakage is more important than the dynamic for deep submicron design. In the following text, it will focus on the leakage reduction for my pre-comparison CAM.

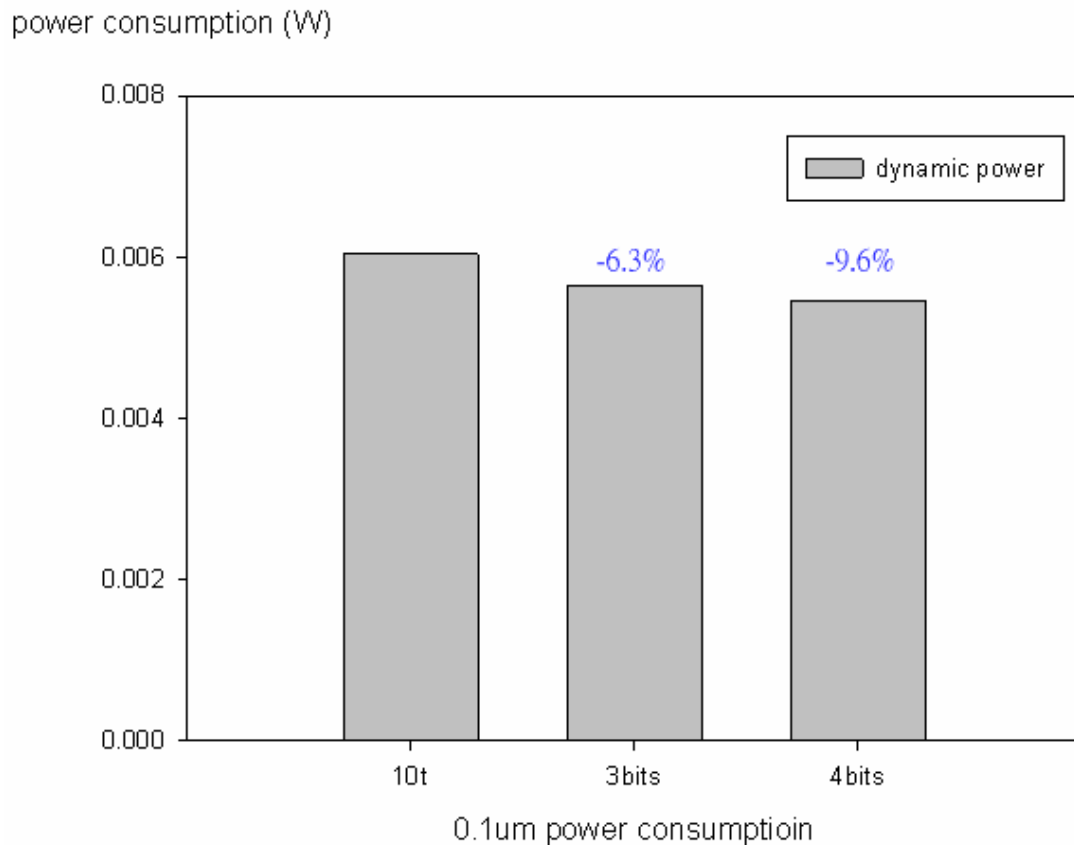


Fig. 5.23 Dynamic power consumption for 100nm CAM

5.4.2 Simulation of Leakage Power for Pre-comparison CAM

Pre-comparison CAM combines the techniques of power gating and dual vdd in Fig. 5.24. The leakage power is compared in Fig. 5.25. There are four types of pre-comparison CAM, original, power gating, dual vdd, and power gating as well as dual vdd pre-comparison CAM. The power gating pre-comparison CAM has 22% leakage reduction. The dual vdd has 26.3% power saving. Combining power gating and dual vdd can reduce 31.1% power dissipation. So combining these two techniques actually reduces the leakage power consumption.

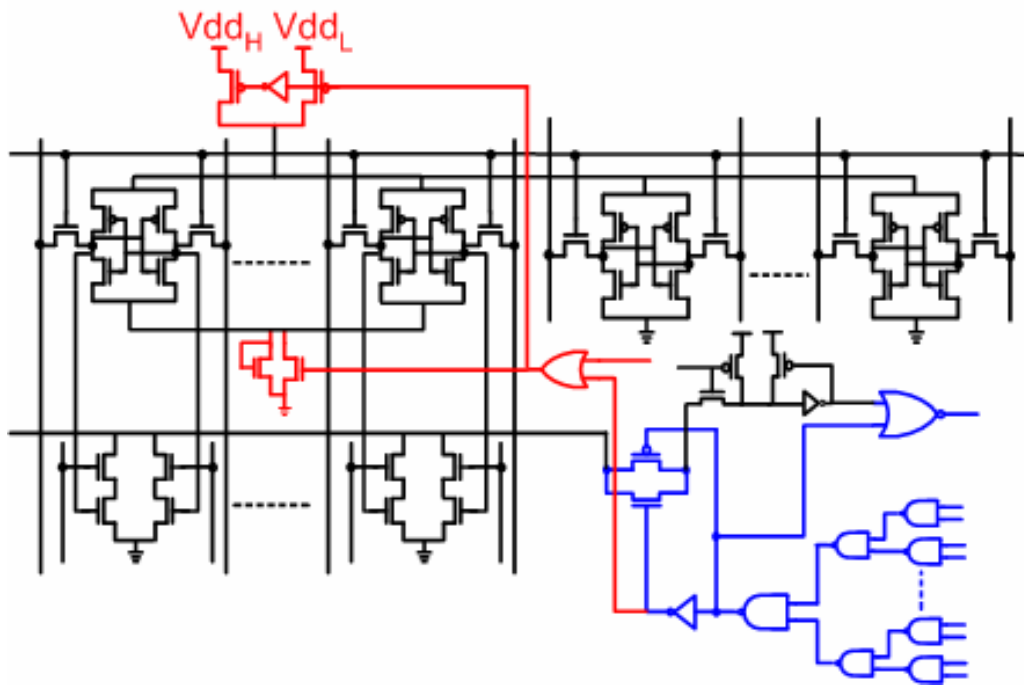


Fig. 5.24 Pre-comparison CAM combines power gating and dual vdd

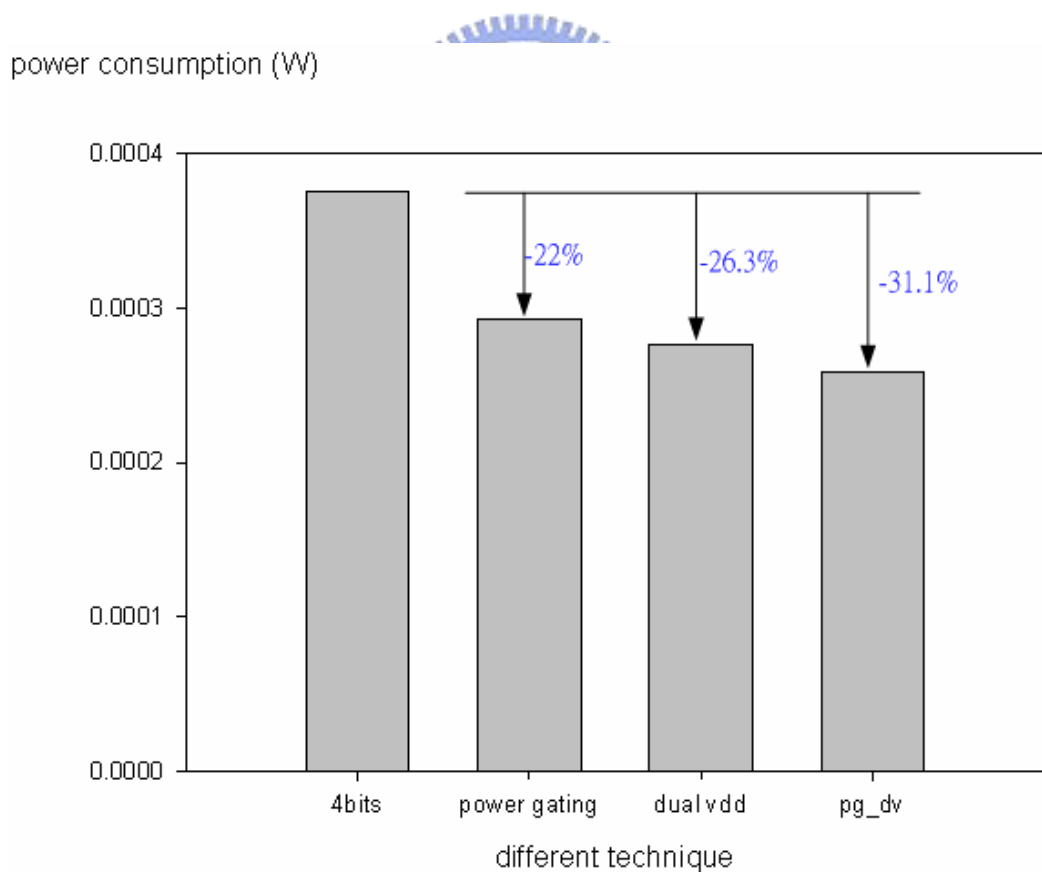


Fig. 5.25 Power consumption for power gating and dual vdd pre-comparison CAM

5.5 Component of TLB

In chapter 2, we have introduced what is TLB. We would introduce components of TLB in this section. From Fig. 5.26, we can know that TLB is mainly composed of CAM array and SRAM array. There are some peripheral architecture circuits in CAM array and SRAM array, such as bits line sense amplifiers, bit line precharge circuits, word line address decoder, match line precharge circuit, and etc. When a virtual address is sent into TLB, it would be compared through the CAM array in parallel. If it is match, the match signal would be transmitted to SRAM. The chosen word line of SRAM would send the data, real address, out and return a hit signal. If it is mismatch, the match signal would be logic low and do not turn the word line of SRAM on. The SRAM of TLB is some different for the conventional SRAM. For the comparison part, CAM, its SRAM only has row decoder and does not have column decoder. For the stored real address part, its SRAM also has no column decoder. It has two conditions to access the word line. Firstly, the real address is written into SRAM. Secondly, the match of CAM is match and sends the match signal to pull the corresponding word line up to read out the real address, Fig. 5.27. The other peripheral circuits are the same as the traditional SRAM, such as, sense amplifier, equalizer, pass transistors, and so on.

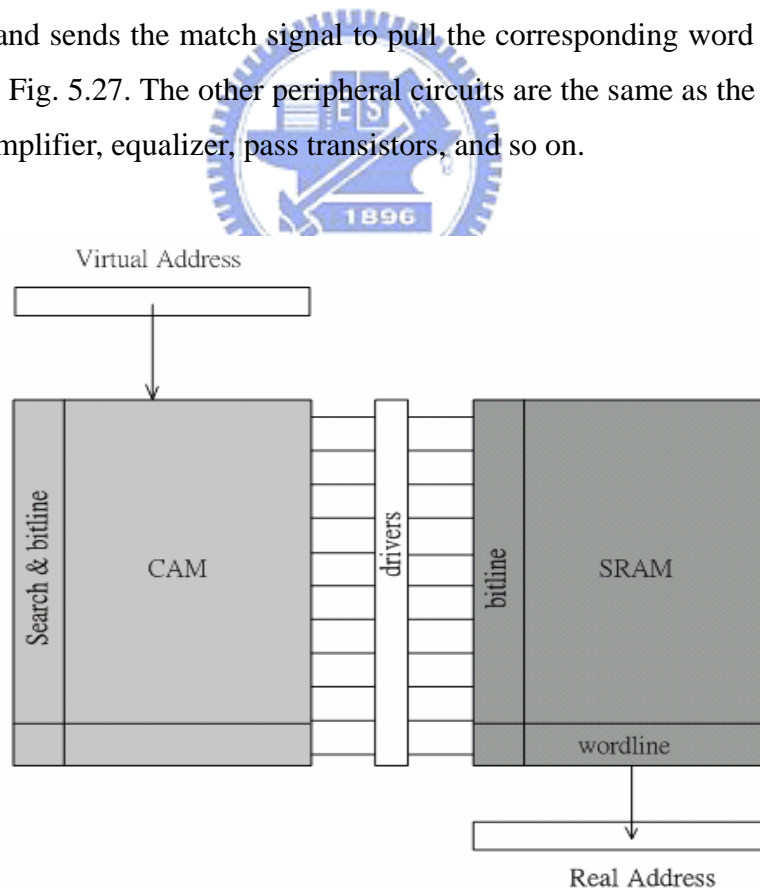


Fig. 5.26 Component of TLB

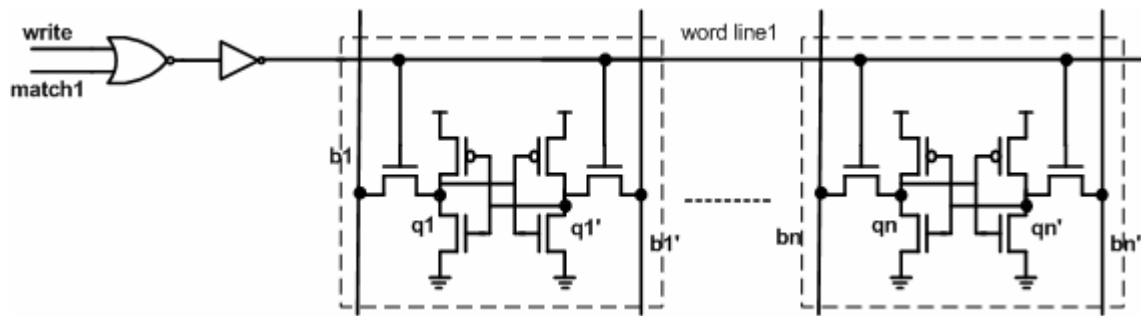
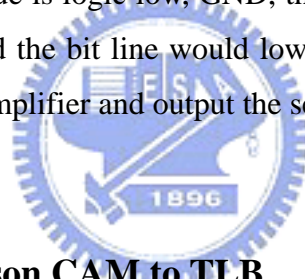


Fig. 5.27 The SRAM part of TLB

Fig is one bit organization for SRAM. Its main part is the six transistors that compose the two latches and two pass transistor gate. The peripheral circuits have equalizer and sense amplifier. Equalizer is used to precharge bit line and bit line bar to vdd in read operation. The wanted word line would turn on after the equalizer precharges the bit line and bit line bar. If the stored value is logic high, VDD, the bit line and the stored value is the same and no current flows. If the stored value is logic low, GND, the bit line would have current flow that flows into the stored node, and the bit line would lower its voltage. The bit line and bit line bar are sensed through sense amplifier and output the sense value.



5.5.1 Apply Pre-comparison CAM to TLB

In chapter3, we have introduced my low power pre-comparison CAM cell. Here, apply pre-comparison CAM into TLB. TLB is mainly composed of SRAM and CAM. Fig is the TLB that is used the application of pre-comparison CAM for one word line. The pre-comparison CAM has four bits pre-comparison circuit. The chosen bit number is discussed in chapter four. We use the result of optimal pre-comparison bit number in chapter four. The left part is the CAM cell, and it has twenty-two bits. The right side is the SRAM, and it has thirty-six bits. The pre-comparison circuit is on the center of the Fig. 5.28 that is marked by blue color. For each access, the CAM would compare to see if it is match. If the comparison is match, the match signal would be sent into the word line of SRAM. It would turn on the word line of SRAM to read the data that is stored in SRAM. There are two signals to control the stored data SRAM. One is match line signal and the other is the original SRAM word line signal. These two signals are combined with OR gate. Word line signal of SRAM is

responsible for writing in data. Match signal is responsible for reading out data that is stored in SRAM. Fig is one entry TLB. Fig is complete TLB. In the following discussion about TLB, its size of array is thirty-two entries. Each word line has twenty-eight bits CAM cell, four bits pre-comparison circuit, and thirty-six bits SRAM cell.

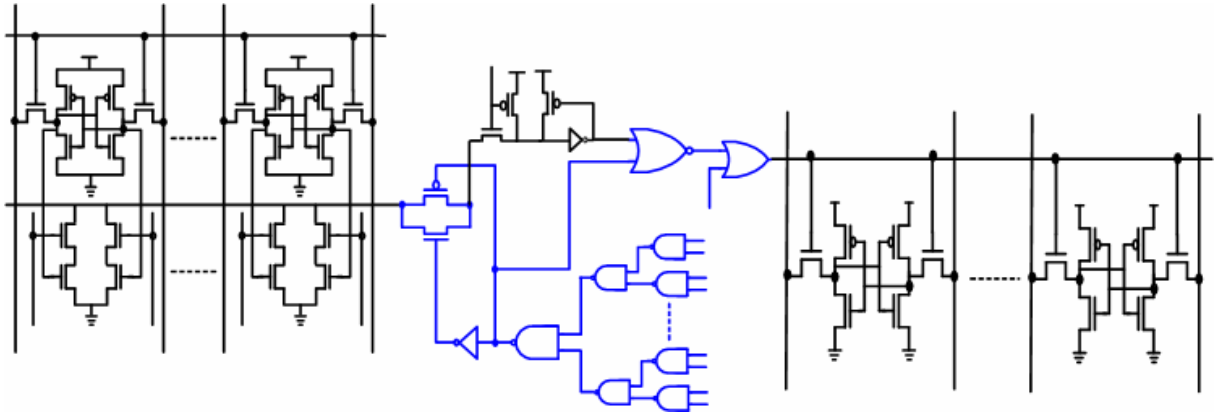


Fig. 5.28 Pre-comparison TLB

5.5.2 Apply Power Gating and Dual Vdd to Pre-comparison TLB

Fig. 5.29 shows the combination of power gating and dual vdd in pre-comparison TLB for one word line. The dual vdd is used 1V and 0.8V. The simulation result is in Fig. 5.30. The size of TLB is 32 entries. The array size of CAM has 32x32, and the array size of SRAM has 32x36.

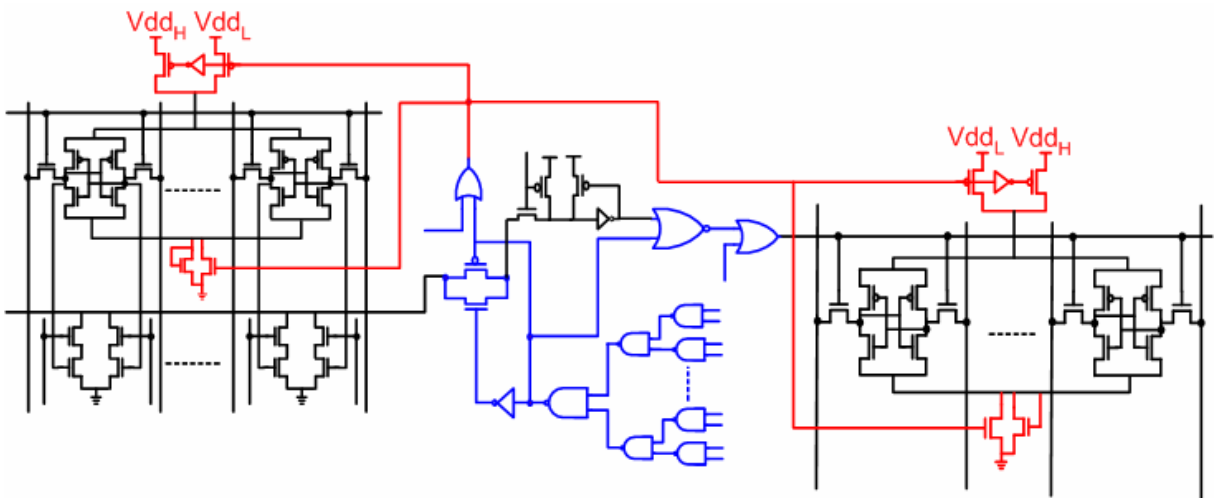


Fig. 5.29 Pre-comparison TLB combines power gating and dual vdd

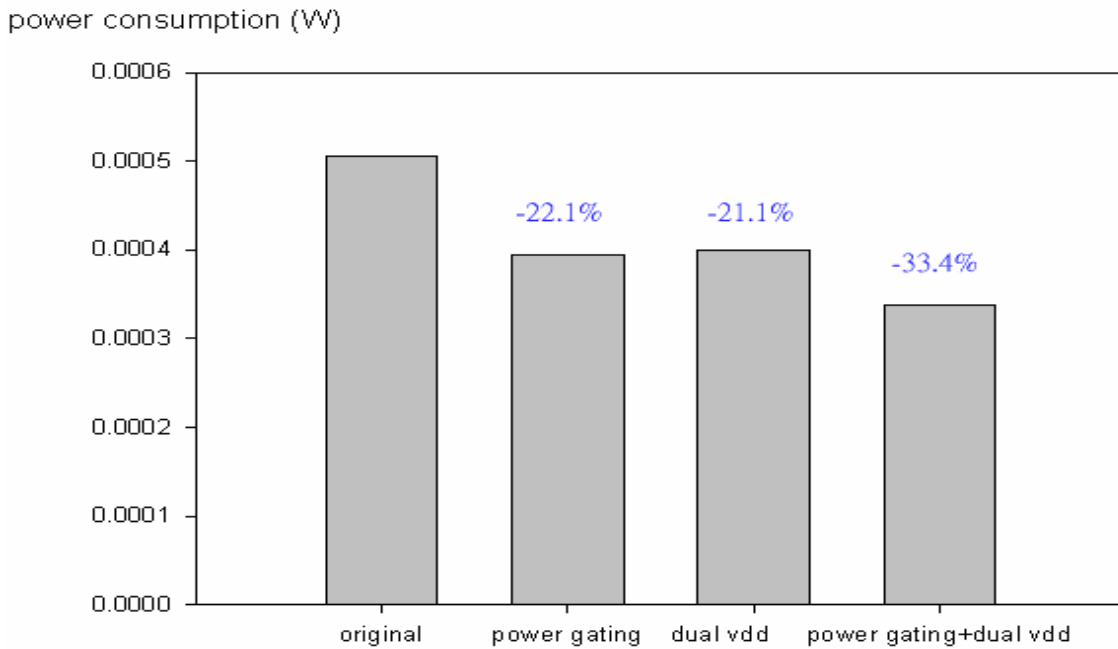


Fig. 5.30 Leakage power for four types of TLB



5.6 Conclusions

In this chapter, the technique of power gating and dual vdd are combined and applied in SRAM, pre-comparison CAM, and pre-comparison TLB. The dual vdd is used 1V and 0.8V. From the simulation result, the leakage power of pre-comparison CAM can reduce 22% for only power gating, 26.3% for only dual vdd, and 31.1% for power gating and dual vdd. The leakage power of pre-comparison TLB can reduce 22.1% for only power gating, 21.1% for only dual vdd, and 33.4% for combing these two techniques.

Chapter 6 Conclusions and Future Work

6.1 Conclusions

Memory hierarchy is introduced in chapter2. The component of memory hierarchy has Cache, translation lookaside buffer (TLB), and etc. Some low power design for Cache and TLB are also described in this chapter. Content-addressable memory (CAM) is widely applied to Cache and TLB. The NOR type and NAND type CAM are also introduced in this section.

Dynamic power is always the dominant power for circuit design. A pre-comparison CAM is proposed in chapter3. It avails a pre-comparison circuit to make the pre-comparison. Through the result of pre-comparison, it can efficiently reduce the times of discharging and find out some mismatch word lines for some pre-comparison bits in advance. The simulation is based on TSMC 0.13um CMOS technology. The result shows that the one pre-comparison bit of NOR type 10T CAM can save 6%, and two bits can save 19.8% for 32 word lines and 32 bits lines of CAM array.

Pre-comparison CAM has improved in chapter4. The enhancement of pre-comparison CAM can have 3% reduction of match line capacitance as the pre-comparison bits increase. The access time for the improved pre-comparison CAM reduces 6.4% and 10.4% than the original pre-comparison CAM for one pre-comparison bit. For two pre-comparison bits, it can reduce 37.6% and 9.7%. Take four pre-comparison bits for index. The following comparison is about conventional NOR type 10T CAM and four bits pre-comparison NOR type 10T CAM. The low power one increases 13.4% and decreases 35.7% access time than traditional one. The power consumption can reduce 22.8%. The power delay product has 12.5% and 50.4% reduction. The penalty is additional forty-eight transistors for comparison circuit.

Leakage current is described in chapter5. As the technology scales down, the leakage current is not neglected. The techniques of power gating and dual vdd are applied in this section. The simulation is based on the TSMC 100nm CMOS technology. The pre-comparison CAM for four bits combining these two skills is actually reduced 31.1% leakage power than the pre-comparison CAM for four bits. Apply the four bits of pre-comparison CAM into TLB.

The four bits of pre-comparison TLB can save 33.4% leakage power than the original four bits of pre-comparison TLB.

6.2 Future Work

The dual vdd in this thesis uses two independent supply voltages. It will increase the difficulty of layout because of two power line and need more area overhead. Therefore, how to use only one supply voltage is a way. Using a DC/DC converter can reach this goal. The concept shows in Fig. 6.1. When supply voltage comes in the system, it will go to the DC/DC converter firstly. The DC/DC converter can send out the only voltage which is VDD high or VDD low to the memory circuit. A signal comes from the memory circuit. It can control the DC/DC converter to decide VDD high or VDD low. This method can reduce the complexity of layout for the original dual vdd architecture.

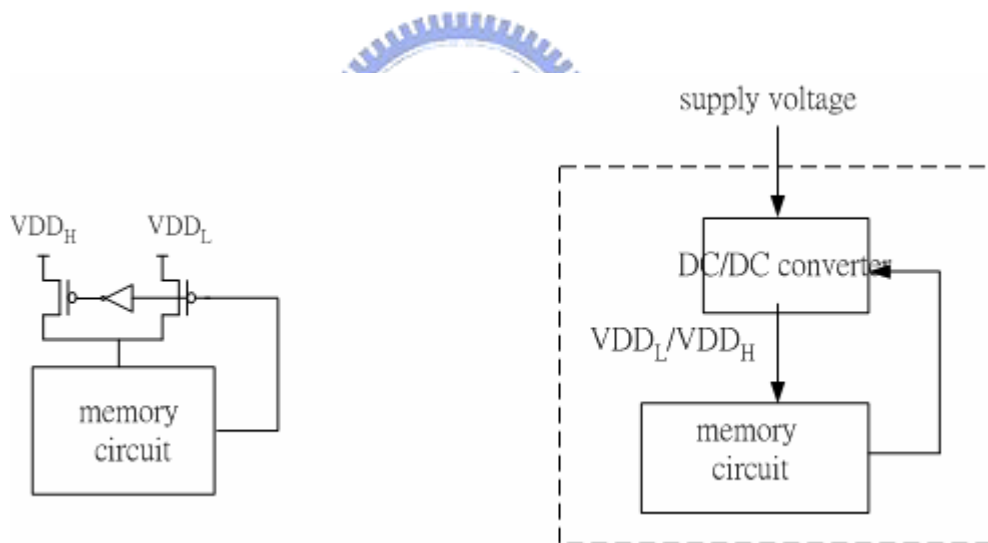


Fig. 6.1 Use a DC/DC to produce high or low voltage

References

- [1] John. L. Hennessy, and David A.Patterson,“ **Computer Architecture - A Quantitative Approach,**” Morgan Kaufmann, 3nd edition.
- [2] John. L. Hennessy, and David A.Patterson,“ **Computer Organization & Design – The Hardware / Software Interface,**” Morgan Kaufmann, 3nd edition.
- [3] M. Yoshimoto, K. Anami, H. Shinohara, T. Yoshihara, H. Takagi, S. Nagao, S. Kayano, T. Nakano,” **A divided word-line structure in the static RAM and its application to a 64K full CMOS RAM,**” IEEE Journal of Solid-State Circuits, Oct 1983, pp. 479-485.
- [4] J. M. Rabaey, A.Chandrakasan, and B.Nikolic,“ **Digital Integrated Circuits,**” Prentice Hall, 2nd edition.
- [5] Zhao Xue-Mei, Ye Yi-zheng, Yu Ming-yan, Li Xiao-ming,” **A Fast Low Power Embedded Cache Memory Design,**” ASIC, Oct. 2001, pp. 566-569.
- [6] M. Sinha, S. Hsu, A. Alyandpour, W. Burlison, R. Krishnamurthy, S. Borkar,” **Low voltage sensing techniques and secondary design issues for sub-90nm caches,**” proc. of ESSCIRC, Sept. 2003, pp. 413-416.
- [7] M. Sinha, S. Hsu, A. Alyandpour, W. Burlison, R. Krishnamurthy, S. Borkar,” **High-performance and low-voltage sense-amplifier techniques for sub-90nm SRAM,**” IEEE International SOC Conference, Sept. 2003, pp. 113-116.
- [8] J. R. Haigh, M. W. Wilerson, J. B. Miller, T. S. Beatty, S. J. Strazdus, L. T. Clark,“**A low-power 2.5-GHz 90-nm level 1 cache and memory management unit,**” IEEE Journal of Solid-State Circuits, May. 2005, pp. 1190-1199.
- [9] K. Ghose, M. B. Kamble,” **Reducing power in superscalar processor caches using subbanking, multiple line buffers and bit-line segmentation,**” In Proceedings ISLPED, 1999, pp. 70-75.

[10] A. Agarwal, Li Hai, K. Roy," **DRG-cache: a data retention gated-ground cache for low power,**" Design Automation Conference, June 2002, pp. 473-478.

[11] A. Agarwal, Li Hai, K. Roy," **A single-V/sub t/ low-leakage gated-ground cache for deep submicron,**" IEEE Journal of Solid-State Circuits, Feb. 2003, pp. 319-328.

[12] Norman W. Petty, Boulder, and Colo," **Content-addressable memory implemented with a memory management unit,**" US patent 6,026,467, Feb. 15, 2000.

[13] Nik Shaylor, et al," **Method and apparatus for a high-performance embedded memory management unit,**" US patent 6,233,667, May 15,2001.

[14] I. Kadayif, A. Sivasubramaniam, M. Kandemir, G. Kandiraju, and G. Chen," **Generating physical addresses directly for saving instruction TLB energy,**" Proc. MICRO, Nov. 2002, pp. 185-196.

[15] P. Petroy, A.Orailoglu," **Virtual page tag reduction for low-power TLBs,**" ICCD, Oct. 2003, pp. 371-374.



[16] B. Jacob, T. Mudge," **Uniprocessor virtual memory without TLBs,**" IEEE Transactions on Computers, Vol. 50, No. 5, May. 2001, pp. 482-499.

[17] Jin-Hyuck Choi, Jung-Hoon Lee, Gi-Ho Park, and Shin-Dug Kim," **An advanced filtering TLB for low power consumption,**" SCAB-PAD, Oct. 2002, pp. 93 – 99.

[18] Jung-Hi Min, Jung-Hoon Lee, Seh-Woong Jeong, and Shin-Dug Kim," **A selectively accessing TLB for high performance and lower power consumption,**" ASIC, Aug. 2002, pp. 45-48.

[19] Jung-Hoon Lee, Gi-ho Park, Sung-Bae Park, and Shin-Dug Kim," **A selective filter-bank TLB system [embedded processor MMU for low power],**" ISLPED , Aug. 2003, pp. 312 – 317.

[20] V. Delaluz, M. Kandemir, A. Sivasubramaniam, M. J. Irwin, N. Vijaykrishnan," **Reducing dTLB energy through dynamic resizing,**" Proceedings 21st International Conference on Computer Design, Oct. 2003, pp. 358-363.

[21] David Channon, David Koch," **Performance Analysis of Re-configurable Partitioned TLBs,**" Proceedings of the Thirtieth Hawaii International Conference on System Sciences, Vol. 5, Jan. 1997, pp. 168-177.

[22] F. Shafai, K.J. Schultz, G.F.R. Gibson, A. G. Bluschke, D.E. Somppi, "**Fully Parallel 30-MHz 2.5-Mb CAM,**" IEEE Journal of Solid-State Circuits, Vol. 33, No. 11, Nov. 1998, pp. 1690-1696.

[23] H. Miyatake, M. Tanaka, and Y. Mori, "**A Design for High-Speed Low-Power CMOS Fully Parallel Content-Addressable Memory Macros,**" IEEE Journal of Solid-State Circuits, Vol. 36, No.6, pp. 956-968, June 2001.

[24] I. Arsovski, T. Chandler, and A. Sheikholeslami, "**A Ternary Content-Addressable Memory (TCAM) Based on 4T Static Storage and Including a Current-Race Sensing Scheme,**" IEEE Journal of Solid-State Circuits, Vol. 38, No.1, pp. 155-158, Jan 2003.

[25] A. Roth, D. Foss, R. Mckenzie, and D. Perry, "**Advanced Ternary CAM Circuit on 0.13um Logic Process Technology,**" Proceedings of the IEEE 2004 Custom Integrated Circuits Conference, pp. 465-468, Oct. 2004.

[26] Gandhi Thirugnanam, N. Vijaykrishnan, and Mary Jane Irwin," **A Novel Low Power CAM Design,**" in Proc. IEEE ASIC/SOC Conf., September 2001, pp. 198-202.

[27] C. Zukowski and S. Wang," **Use of Selective Precharge for Low-Power Content Addressable Memories,**" IEEE International Symposium on Circuits and Systems, June 9-12, 1997, pp. 1788-1791.

[28] T.Juan, T. Lang, J. J. Navarro," **Reducing TLB power requirements,**" in Proc. 1997 Int. Symp.on Low Power Electronics Design, 1997, pp. 196-201.

- [29] T. Jamil," **RAM versus CAM**," IEEE Potentials, April/May 1997, pp. 26-29.
- [30] Chi-Sheng Lin, Jui-Chuan Chang, and Bin-Da Liu," **A Low-Power Precomputation-Based Fully Parallel Content-Addressable Memory**," IEEE J.Solid-State Circuit, Vol. 38, No. 4, April 2003, pp. 654-662.
- [31] Chi-Sheng Lin, Jui-Chuan Chang, and Bin-Da Liu," **Low-Power And Low-Voltage Fully Parallel Content-Addressable Memory**," ISCAS, May 2003, pp.25-28.
- [32] C. S Lin, J. C. Chang, and B. D. Liu," **Design For Low-Power, Low-cost, and HighReliability Precomputation-Based Content-Addressable Memory**," APCCAS, Oct. 2002, pp. 319-324.
- [33] Aristides Efthymiou, and Jim D. Garside," **A CAM With Mixed Serial-Parallel Comparison for Use in Low Energy Caches**," IEEE transactions on VLSI systems, March 2004, pp. 325-329.
- [34] Aristides Efthymiou, and Jim D. Garside," **An Adaptive Serial-Parallel CAM Architecture for Low-Power Cache Blocks**," ISLPED, August 2002, pp. 136-141.
- [35] Kuo-Hsing Cheng, Chia-Hung Wei, Shu-Yu Jiang," **Static Divided Word Matching Line For Low-Power Content Addressable Memory Design**," ISCAS, May 2004, pp. 629-632.
- [36] Yi-Liang Hsiao, Ding-Hao Wang and Chein-Wei Jen," **Power Modeling and Low-Power Design of Content Addressable Memories**," ISCAS, May 2001, pp. 926-929.
- [37] H. Qin, Y. Cao, D. Markovic, A. Vladimirescu, and J. Rabaey," **SRAM Leakage Suppression by Minimizing Standby Supply Voltage**," Proceedings of 5th International Symposium on Quality Electronic Design, pp. 55-60, 2004.
- [38] [Http://www-device.eecs.berkeley.edu](http://www-device.eecs.berkeley.edu): BSIM 100nm and 70nm predictive technology process files.

[39] International Technology Roadmap for Semiconductors 2001 edition, Semiconductor Industry Association, <http://public.itrs.net>.

[40] Nam Sung Kim, Krisztian Flautner, David Blaauw, and Trevor Mudge," **Circuit and Microarchitectural Techniques for Reducing Cache Leakage Power**," IEEE Transactions on Very Large Scale Integration System, Vol. 12, No. 2, Feb. 2004, pp. 167-184.

[41] J. Kao and A. Chandrakasan," **Dual-Threshold Voltage Techniques for Low-Power Digital Circuits**," IEEE Journal of Solid-State Circuits, July 2000, Vol. 35, No. 7, pp. 1009-1018.

[42] Takeshi Kitahara, Naoyuki Kawabe, Fimihiro Minami, Katsuhiro Seta, and Toshiyuki Furusawa," **Area-efficient Selective Multi-Threshold CMOS Design Methodology for Standby Leakage Power Reduction**," Proceedings of the Design, Automation and Test in Europe Conference and Exhibition, 2005, pp. 646-647.

[43] Kimiyoshi Usami, Naoyuki Kawabe, Masayuki Koizumi, Katsuhiro Seta, and Toshiyuki Furusawa," **Automated Selective Multi-Threshold Design for Ultra-Low Standby Applications**," ISLPED, 2002, pp. 202-206.

[44] Benton H. Calhoun, Frank A. Honore, and Anantha P. Chandrakasan," **A Leakage Reduction Methodology for Distributed MTCMOS**," IEEE Journal of Solid-State Circuits, May 2004, Vol. 39, No. 5, pp. 319-328.

[45] A. Agarwal et al," **A Single-Vt Low -Leakage gated-ground Cache for Deep Submicron**" IEEE Journal of Solid-State Circuits, Feb. 2003, pp 319-328.

[47] Yibin Ye, Shekhar Borkar , and Vivek De," **A New Technique for Standby Leakage Reduction in High-Performance Circuits**," in Symp. VLSI Circuits, June 1998, pp. 40-41.

[47] K. Nii, H. Makino, Y. Tujihashi, C. Morishima, Y. Hayakawa, H. Nunogami, T. Arakawa, and Hisanori Hamano," **A Low Power SRAM using Auto-Backgate-Controlled MT-CMOS**," International Symposium on Low Power Electronics and Design, August 1998, pp. 293-298.

[48] Dinesh Somasekhar, Yibin Ye, and Kaushik Roy,” **An Energy Recovery Static RAM Memory Core,**’ In IEEE Sym. On Electronics, 1995, pp. 62-63.

[49] K. Zhang, U. Bhattach, Z. Chen, F. Hamzaoglu, D. Murray, N. Vallepalli, Y. Wang, B. Zheng, and M. Bohr,” **SRAM Design on 65-nm CMOS Technology With Dynamic Sleep Transistor for Leakage Reduction,**“ IEEE Journal of Solid-state Circuits, Vol. 40, No. 4, April 2005, pp 895-901.

[50] M. Yamaoka, et al,” **A 300-MHz 25uA/Mb-Leakage On-Chip SRAM Module Featuring Process-Variation Immunity and Low-Leakage-Active Mode for Mobile-Phone Application Processor,**” IEEE Journal of Solid-State, Vol. 40, No. 1, Jan. 2005, pp. 186-194.



Vita

PERSONAL INFORMATION

Chinese Name: 張維耿

English Name: Wei-Keng Chang

Birth Date: Dec 9, 1980

Birth Place: Taipei, Taiwan, R.O.C.

Address: Department of Electronics Engineering
National Chiao Tung University
1001 Ta-Hsueh Road
Hsin-chu, Taiwan 30050, R.O.C.

E-Mail Address: wagon.ee92g@nctu.edu.tw



EDUCATION

B.S. [2003] Department of Electronic Engineering, National Cheng-Kung University.

M.A. [2005] Institute of Electronics, National Chiao-Tung University.