

TEMA 3. ALMACENAMIENTO DE LA INFORMACIÓN

- 3.1. [Introducción. Jerarquía de memoria y memoria interna del computador](#)
- 3.2. [Memoria interna del computador](#)
 - 3.2.1. [Mapa de memoria de un computador](#)
 - 3.2.2. [Configuración de la memoria principal](#)
- 3.3. [Recursos para mejorar las prestaciones de la memoria principal](#)
 - 3.3.1. [Introducción. Aumento de la capacidad real y aumento de la velocidad de acceso](#)
 - 3.3.2. [Memoria virtual](#)
 - 3.3.3. [Memoria cache o inmediata](#)
 - 3.3.4. [Protección de la memoria principal](#)
- 3.4. [Dispositivos de almacenamiento](#)

El almacenamiento de la información es una operación imprescindible en todo computador tipo Von Neumann (Máquina de Programa Almacenado).

La información a procesar se encuentra almacenada en una configuración de memoria con distintos niveles de jerarquía; cuanto más inmediato sea el procesamiento de la información más cerca del procesador se encontrará ésta. Cada nivel de jerarquía de memoria se caracteriza por unos tiempos de acceso y por unas capacidades de almacenamiento; los niveles de jerarquía superior (los más cercanos al procesador) tienen los menores tiempos de acceso y las menores capacidades de almacenamiento.

Los niveles de jerarquía básicos existentes en todo sistema computador son: el de registro, el de memoria principal y el de memoria secundaria. El nivel de memoria principal representa el espacio de trabajo del procesador; a este nivel se dirige el procesador para buscar tanto las instrucciones, que definen qué tratamiento hay que aplicar, como los datos a los que hay que aplicar dicho tratamiento. Por lo tanto, para este nivel se demandan una gran capacidad de almacenamiento y un tiempo de acceso pequeño. Para que la configuración e implementación de este nivel impliquen un coste razonable, aparecen los niveles de memoria caché y de memoria virtual; el de memoria caché para acelerar el acceso a la información y, el de memoria virtual para aumentar el área de trabajo del computador (aumentar la capacidad).

Este tema trata los mecanismos de memoria caché y el de memoria virtual. El primero se resuelve exclusivamente por hardware y la memoria virtual implica tanto elementos de tipo hardware como de tipo software. Además, al final del tema se presentan los distintos dispositivos de lectura/escritura comerciales para configurar el nivel de jerarquía de memoria principal, así como los existentes para el nivel de memoria secundaria.

3.1.INTRODUCCIÓN. JERARQUÍA DE MEMORIA Y MEMORIA INTERNA DEL COMPUTADOR

Existen muchas tecnologías para fabricar las memorias. Desde el punto de vista de su utilización, éstas se caracterizan por tres propiedades fundamentales:

- Coste por bit.
- Tiempo que se tarda en acceder a la información.
- Capacidad de almacenamiento o tamaño típico.

El coste por bit decrece muy rápidamente con la velocidad y a menor velocidad mayor capacidad. Ahora bien, para un computador se desea:

- que tenga mucha memoria (una gran capacidad de almacenamiento),
- que ésta sea muy rápida. Ya que la memoria abastece a la CPU de datos e instrucciones, y esta CPU se constituye con circuitos integrados rapidísimos, se obliga a no gastar tiempo en las lecturas y escrituras con el fin de que sea lo más rápido posible el procesado. Por otro lado, se desea también
- que el precio total del sistema de memoria sea reducido.

Por lo tanto, la memoria ideal debe poseer una capacidad elevada junto a un tiempo de acceso muy pequeño y un precio reducido. Dichas características, por ahora, son tecnológica y económicamente irreconciliables.

Por todo ello, la memoria de los computadores se suele estructurar en varios niveles. Existirá un nivel rápido de pequeña capacidad y niveles sucesivos de menor velocidad, pero mayor capacidad. La información se ubicará en uno de los niveles de acuerdo con su probabilidad de uso; así, un programa o unos datos poco empleados estarán almacenados en el nivel inferior, más lento y de mayor capacidad; y, si en un momento determinado se necesitan, se pasarán al nivel superior más rápido para ser utilizados.

En lo referente al sistema de memoria de un sistema computador, se pueden distinguir los siguientes niveles jerárquicos:

- a) REGISTROS. En ellos estará la información inmediata a procesar. Son memorias de tipo RAM con capacidad del orden de bytes. Tecnología de semiconductores.
- b) MEMORIAS CACHE. Son memorias de tipo RAM de baja capacidad (Kbytes) y muy rápidas. Tecnología de semiconductores.
- c) MEMORIA PRINCIPAL. La mayor parte de los dispositivos que configuran esta memoria son de tipo RAM de gran capacidad (Mbytes) y tiempo de acceso rápido. Tecnología de semiconductores.

- d) **MEMORIAS INTERMEDIAS.** Memorias de acceso secuencial, aleatorio o híbridas. Con tiempos de acceso del orden de milisegundos y con una capacidad elevada (Gbytes). Realmente este nivel no se ha extendido mucho y suele ser sustituido por un nivel de MEMORIA SECUNDARIA referente al disco de cabeza fija y de uso exclusivo del sistema operativo.
- e) **MEMORIAS AUXILIARES.** Son lentas y de gran capacidad. Como ejemplos representativos se pueden dar los discos de cabezas móviles y las cintas magnéticas.
- En la [Tabla 3.1](#) se muestran las características de los distintos niveles de jerarquía de la memoria.

Jerarquía	Capacidad (en Bytes)	Tiempo de acceso	Ancho de banda	Tipo	Acceso elemental
REGISTROS	6-200	0.25-0.5 ns	20-100 GB/s	Biestables	Palabra
M. CACHE	8 K-8 M	0.5-25 ns	5-10 GB/s	SRAM	Palabra
M. PRINCIPAL	1 M-20 G	80-300 ns	1-5 GB/s	DRAM	Palabra
M. SECUNDARIA	5 G-100 G	5-20 ms	20-150 MB/s	A. Directo	Sector
M. AUXILIAR	300 K-800 G	minutos	1-5 MB/s	A. Secuencial	Registro

Tabla 3.1. Características de los distintos niveles de jerarquía de la memoria.

Dado que el interés de la información varía según las necesidades de los usuarios, se debe producir un movimiento continuo ascendente y descendente de la información en la jerarquía de memoria, de forma que, aquella que debe ser utilizada por la unidad central de proceso en un momento determinado, se encuentre en los niveles superiores.

Se puede decir que, aquellos módulos de información necesarios en un momento determinado, por ejemplo porque un usuario llama a su fichero para editarlo desde un terminal, son como burbujas de información que deben subir por la jerarquía de memoria, para estar disponibles en el nivel requerido. Si esta información es modificada por el procesador y debe conservarse una vez terminado su proceso, deberá bajar otra vez a su nivel correspondiente, si, por el contrario, la información no ha sido modificada, por ejemplo en un programa que ha sido ejecutado, se destruye (la burbuja explota).

Por *memoria interna del computador* se entiende la definida por los niveles de jerarquía: Registros, Memoria Caché (si la hay) y Memoria Principal.

La Memoria Principal es el espacio de trabajo de la Unidad Central de Proceso; la definición de una dirección por parte de la Unidad de Control siempre se corresponde con una dirección de este espacio, independientemente de que exista o no Memoria Caché. La Memoria Caché se presenta de forma transparente a la Unidad de Control; es otro dispositivo hardware el encargado de la gestión de la información de la Memoria Caché.

La memoria Caché contiene la misma información que determinadas posiciones de la Memoria Principal; si la Unidad de Control direcciona alguna posición de memoria con contenido copiado en la Memoria Caché, ésta se la suministra rápidamente (antes de lo que lo haría la Memoria Principal).

3.2.MEMORIA INTERNA DE UN COMPUTADOR

La memoria interna del computador la forman los tres primeros niveles de jerarquía presentados (Registros, memoria cache y memoria principal). Los registros, entre los que cabe destacar los registros de uso general incluidos en la unidad aritmética y lógica, forman el nivel más rápido y de menos capacidad. La memoria principal, que es la memoria donde tienen que residir los programas y sus datos para poder ser ejecutados por el computador, forman el tercer o segundo nivel, según que el computador disponga o no de memoria cache. Por su lado, esta memoria caché es una memoria auxiliar que se emplea para acelerar los accesos a la memoria principal, pero de forma transparente al usuario, esto es, funcionalmente el computador se comporta como si esta memoria no existiese, pero trabaja más rápido que sin ella. La [Figura 3.1](#) muestra estos tres niveles de jerarquía, así como la relación de la CPU con el resto de los dispositivos de memoria.

Hoy en día la memoria interna de los computadores es de semiconductores. Los registros se construyen con la misma lógica que el resto del computador y la memoria principal con las pastillas de memoria vistas anteriormente. Sin embargo, las memorias auxiliares son en su mayoría del tipo magnético.

En esta sección se va a ver la forma de configurar la memoria principal de un computador.

3.2.1. Mapa de memoria de un computador

Como hemos visto anteriormente, para que un computador pueda ejecutar un programa necesita que éste, así como sus datos, se encuentre en memoria principal. Para referenciar o direccionar estos datos e instrucciones, el computador genera y manipula direcciones de memoria principal. Por la propia construcción del computador, estas direcciones se ven limitadas a un cierto tamaño, correspondiente al número de bits que es capaz de manejar ese computador en las operaciones de direccionamiento. Por ejemplo, si el contador de programa es un registro que tiene m bits, las direcciones de las instrucciones que puede direccionar ese computador se limita al rango $[0, 2^m-1]$. Por su lado, las direcciones de los operandos vendrán también limitadas a m bits. Con frecuencia este número m de bits, que determina el ancho de las direcciones, coincide con el ancho de palabra del computador, pero también puede que sea distinto. Por

ejemplo, existen computadores de 32 bits con un ancho de direcciones de 24 bits, mientras que los microprocesadores de 8 bits suelen emplear direcciones de 16 bits.

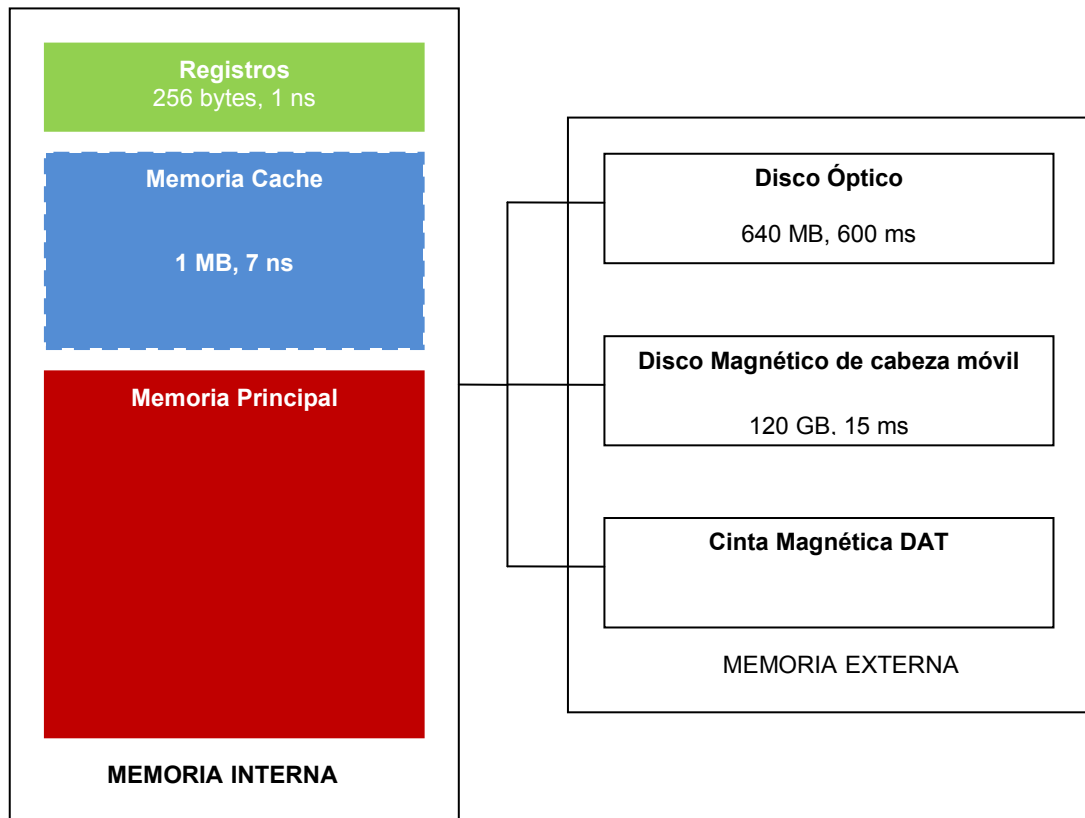


Figura 3.1. Memoria de un computador.

Se llama **MAPA de MEMORIA** a todo el espacio direccionable por un computador. Este espacio viene determinado por el ancho de las direcciones, puesto que, decir que se dispone de m bits de dirección, equivale a decir que se tiene un mapa de memoria de 2^m **posiciones**. Por otro lado las m líneas de dirección forman lo que se llama el BUS DE DIRECCIONES. Por extensión, muchas veces se dice que un computador tiene un **mapa de memoria de m bits**, refiriéndose al ancho de las direcciones que maneja.

Generalmente el computador no se equipa con toda la memoria necesaria para llenar su mapa de memoria. Por el contrario, como se muestra en la [Figura 3.2](#), se suele dotar de una parte de la misma, quedando espacio libre para posibles ampliaciones.

El rápido abaratamiento de las memorias de semiconductores ha hecho que los computadores se construyan con memorias principales cada vez mayores.

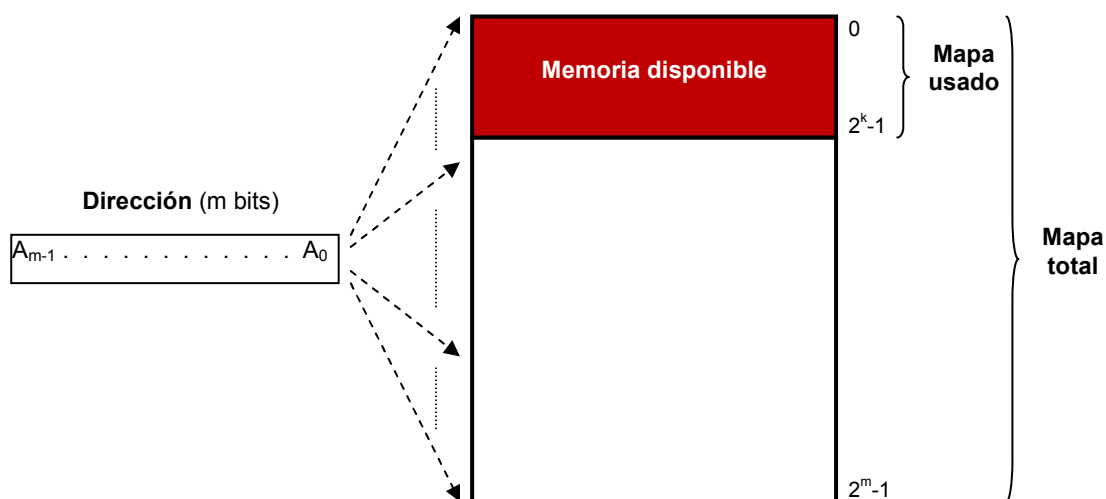


Figura 3.2. Mapa de memoria y memoria principal disponible.

MAPA DE MEMORIA COMÚN Y NO COMÚN

Al igual que a las posiciones de memoria, a los dispositivos de E/S se les debe asignar direcciones. El tratamiento de las direcciones asignadas a los distintos módulos de E/S, según las señales de control del sistema, puede realizarse de dos formas perfectamente diferenciadas:

- 1) Recibiendo el mismo tratamiento que cualquier posición correspondiente a cualquier módulo de memoria.
- 2) Recibiendo un tratamiento diferenciado.

En el primer caso tendremos lo que se denomina sistema con **MAPA DE MEMORIA COMÚN** y, en el segundo caso sistemas con **MAPA DE MEMORIA NO COMÚN**.

Los sistemas de **mapa de memoria común** no tienen instrucciones específicas de E/S, la diferencia entre una operación de E/S y una operación de memoria viene dada por la pertenencia de la dirección, referenciada por el bus de direcciones, a un módulo de E/S o a un módulo de memoria. Los sistemas de **mapa de memoria no común** tienen instrucciones específicas para operaciones de E/S. Por tanto, los módulos (tanto de memoria como de E/S) conectados al bus de direcciones del sistema necesitarán de alguna otra señal que les aclare cuando se les está referenciando. Dicha señal vamos a especificarla como IO/M#, o bien vamos a emplear dos señales IORQ# (Input-Output-Request) e MREQ# (Memory Request).

Los sistemas con **mapa de memoria común** deben asignar algunas posiciones del único mapa existente a los módulos de E/S. Los sistemas con **mapa de memoria no común**, tienen un mapa para los dispositivos de memoria y otro distinto para los

dispositivos de E/S. Estos dos mapas podrían tener, en principio, el mismo número de posiciones pero, generalmente, el mapa de direcciones de E/S tiene menos posiciones, no siendo necesario emplear para la selección de una posición de E/S todas las líneas del bus de direcciones; se suele emplear un conjunto de líneas contiguas del mismo bus de direcciones, empezando por la de más (A_{m-1}) o menos (A_0) peso del mismo (ver [Figura 3.3](#)).

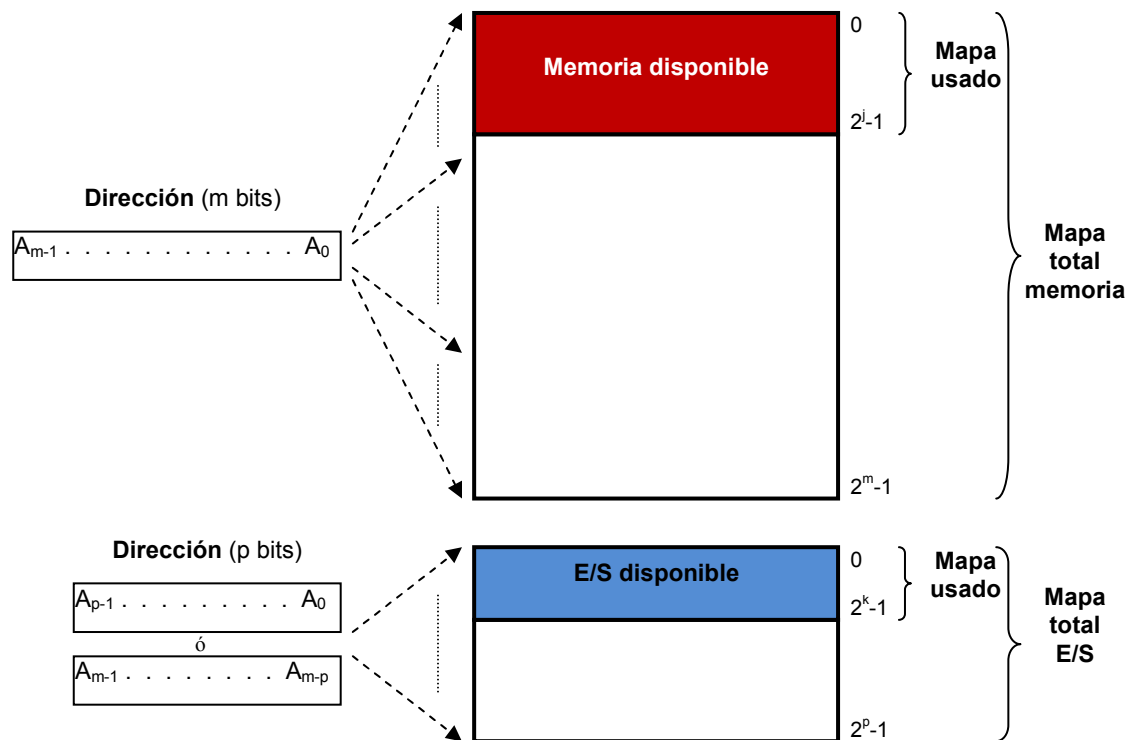


Figura 3.3. Mapas de direcciones en un sistema con *mapa de memoria no común*.

3.2.2. Configuración de la memoria principal

Los dos puntos de partida principales para realizar la configuración de la memoria principal de un computador son el ancho m de las direcciones que es capaz de generar y el ancho n de la palabra que se emplea. El ancho de las direcciones establece el mapa de memoria (número de posiciones de memoria) y deben ser empleadas para seleccionar la palabra deseada, lo que exige realizar una decodificación para activar los puntos de memoria deseados. El ancho de palabra establece el número de puntos de memoria que deben ser activados por cada dirección (número de bits en cada posición de memoria).

Ahora bien, dado que la memoria principal se construye con circuitos integrados, que típicamente tienen una configuración de $pK \times 1$, $pK \times 4$, $pK \times 8$ ó $pK \times 16$ (de manera general $pK \times q$), habrá que establecer la configuración de la memoria en virtud de las pastillas empleadas; teniéndose que completar la *decodificación interna* de la dirección (se realiza

dentro de la pastilla) con una *decodificación externa* para seleccionar las pastillas de memoria adecuadas (la tiene que resolver el usuario añadiendo circuitos decodificadores y según el diseño que se proponga).

DEFINICIÓN DE LA POSICIÓN DE MEMORIA (DEFINICIÓN DE LAS FILAS)

Si se desea una memoria con palabras de n bits y se parte de pastillas con un ancho de q bits, se necesitarán n/q pastillas en paralelo (se activan conjuntamente) para alcanzar el ancho de palabra deseado. Ver [Figura 3.4](#).

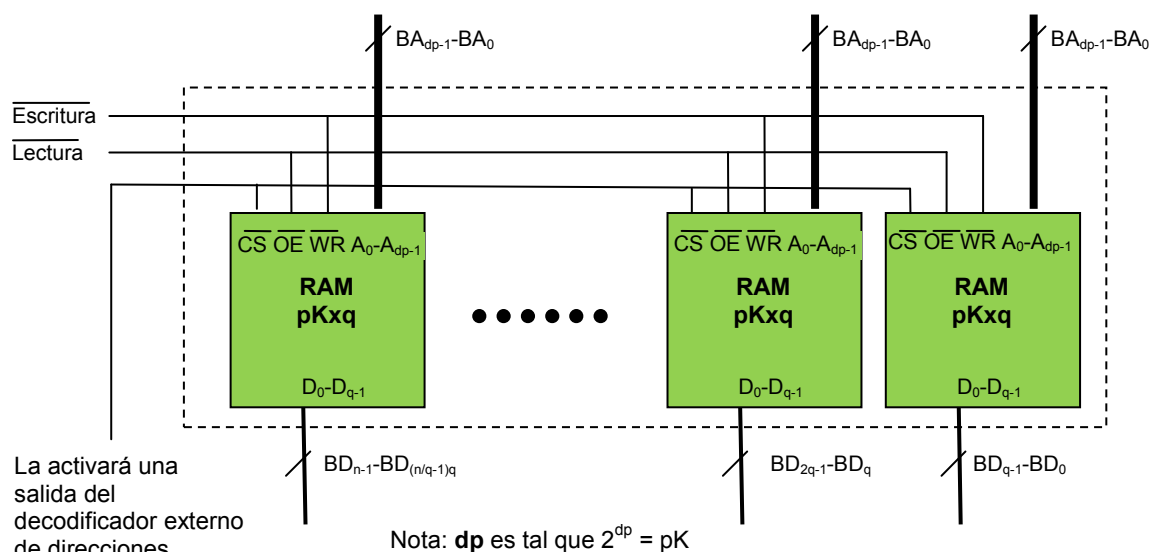


Figura 3.4. Configuración de pK posiciones de memoria de n bits con módulos de memoria de organización $pK \times q$.

DEFINICIÓN DEL NÚMERO DE POSICIONES (DEFINICIÓN DEL NÚMERO DE FILAS)

Una vez construida la fila de módulos de memoria, según el número de bits que deseemos en cada posición (en general n bits) y partiendo de un módulo de memoria con un número determinado de bits en cada posición (en general q bits), hay que definir el número de filas necesarias de ese tipo, partiendo de dos parámetros: número de posiciones que de ese tipo se desean ($m_x K$ posiciones) y número de posiciones que tiene el módulo de memoria a considerar (pK posiciones).

El número de filas necesarias es $m_x K / pK$; ver [Figura 3.5](#). Cada una de las líneas $CS\#$ de selección de fila vendrá definida por una/s u otra/s patilla/s de salida (según qué direcciones del mapa de memoria se hubieran asignado a dichos módulos) del módulo decodificador que realiza la función de decodificación externa de la dirección.

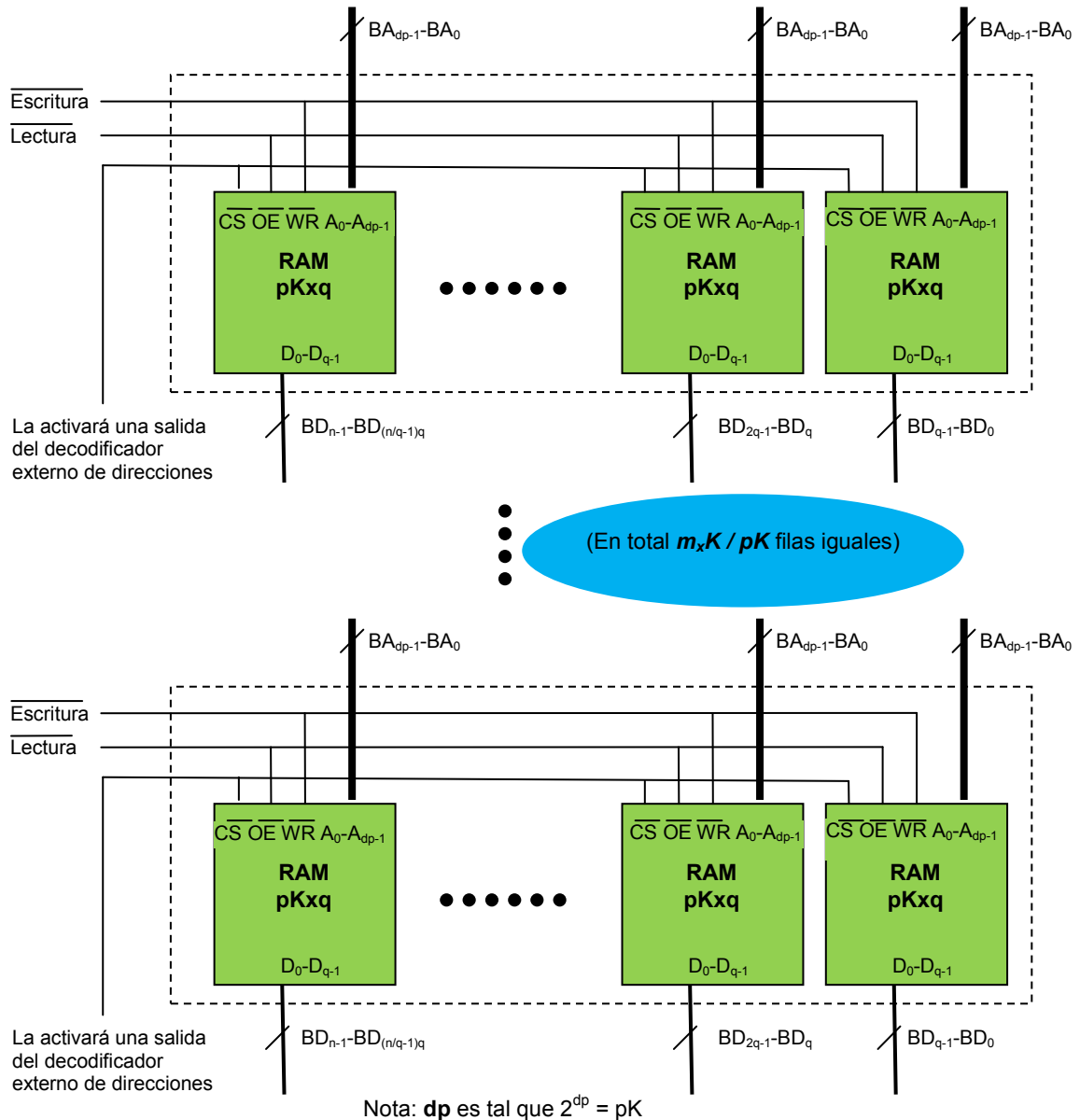


Figura 3.5. Configuración de $d_x K$ posiciones de memoria de n bits con módulos de memoria de organización $pK \times q$.

DISEÑO DE LA MEMORIA PRINCIPAL DE UN SISTEMA COMPUTADOR

La selección de cualquier posición de memoria en un sistema computador se realiza mediante la decodificación de la dirección definida en un momento determinado. Esta función se resuelve mediante una *decodificación interna*, que la resuelven los módulos de memoria internamente, y mediante una *decodificación externa*, que hay que resolver cuando se diseña la memoria del sistema computador correspondiente. Por lo tanto, el diseño de la memoria principal de un sistema computador consiste en resolver la función de *decodificación externa* de las direcciones; además habrá que establecer las conexiones directas de las líneas del bus de direcciones, de datos y control oportunas con las correspondientes de los distintos módulos de memoria.

Por lo tanto, el problema de diseño de la memoria principal se centra principalmente en la resolución de un problema de decodificación. Esta decodificación depende de la distribución de los distintos módulos de memoria en el mapa de memoria correspondiente, es decir, de la asignación de direcciones del mapa de memoria a los distintos dispositivos de memoria.

Con el fin de simplificar el diseño de la memoria principal de un sistema computador, y en la medida de lo posible sistematizar, se va a proponer una metodología de diseño válida para sistemas con mapa de memoria común y no común y considerando dispositivos de memoria de semiconductor asíncronos (la temporización la establecen las señales de control). Evidentemente todo esto con independencia de que se aplique cualquier otra metodología específica para el diseño particular que se considere.

Para configurar la memoria principal de un computador, hay que partir de dos datos principales:

- Ancho m de las direcciones que es capaz de generar → Establece el mapa de memoria del sistema computador.
- Ancho n de las palabras que emplea → Define el número de módulos de memoria de un tipo determinado a utilizar por cada fila.

El tratamiento que se da a las posiciones asignadas a los distintos módulos de E/S en un sistema computador difiere de unos sistemas a otros dependiendo de las señales de control existentes. Se puede realizar de dos formas perfectamente diferenciadas:

1. Recibiendo el mismo tratamiento que cualquier posición correspondiente a cualquier módulo de memoria → **Mapa de memoria común**. En estos sistemas, que no tienen instrucciones específicas de E/S, la diferencia entre una operación de E/S y una operación de memoria viene dada por la pertenencia de la dirección referenciada por el bus de direcciones a un módulo de E/S o a un módulo de memoria.
2. Recibiendo un tratamiento diferente al de las posiciones de memoria → **Mapa de memoria no común**. En estos sistemas, con instrucciones específicas para E/S, está claramente diferenciada las operaciones de memoria y las de E/S. Por lo tanto, además de las líneas de dirección, es necesario alguna que otra señal de control del sistema que aclare cuando se quiere referenciar a un módulo de memoria o a un módulo de E/S. Dependiendo del sistema microcomputador al que se haga referencia, y de forma prácticamente general, o existe una señal: **IO/M#** (a nivel alto cuando se realiza una operación de E/S y a nivel bajo cuando se

realiza una operación en memoria) o existen dos señales que indican operaciones de memoria o de E/S de forma independiente: **MREQ#** (**M**emory **R**equ**e**st) e **IOREQ#** (**I**nput **O**utput **R**equ**e**st), ambas activas a nivel bajo. En estos sistemas existen dos mapas de direcciones, uno exclusivo para los dispositivos de memoria y otro independiente para los dispositivos de E/S. Aunque inicialmente la capacidad de los dos mapas podría ser la misma (los dos con 2^m posiciones), normalmente la capacidad del mapa para los dispositivos de E/S se reduce, por no ser necesario tener tantas posiciones disponibles para los dispositivos de E/S y de esta forma, por necesitarse menos líneas del bus de direcciones para seleccionar una posición en un espacio más reducido ($m_{E/S}$ líneas para seleccionar una posición entre $2^{m_{E/S}}$ posiciones que existan en el mapa de E/S, siendo $m_{E/S} < m$), simplificar en tamaño el campo de dirección de operando en el formato de las instrucciones de E/S. El subconjunto de líneas del bus de direcciones que se emplean para seleccionar una posición de E/S suelen ser contiguas, empezando en la de mayor o menor peso del bus de direcciones (subconjunto $A_{m-1} A_{m-2} \dots A_{m-m_{E/S}}$ ó subconjunto $A_{m_{E/S}-1} A_{m_{E/S}-2} \dots A_0$).

La resolución del problema de decodificación según el método que se va a plantear, lleva a tres soluciones diferentes según la división inicial de la memoria en “*trozos de igual tamaño*” que se haga. Dicho tamaño vendrá definido por el tamaño de los módulos de memoria y/o de E/S que se hayan seleccionado para formar parte de la implementación de la memoria principal de dicho sistema computador. Este tamaño elegido puede ser el correspondiente al número de posiciones del:

1. Módulo de memoria y/o de E/S con menor número de posiciones.
2. Módulo de memoria y/o de E/S con mayor número de posiciones.
3. Módulo de memoria y/o de E/S que predomine al considerar el número de filas que hay que emplear de los mismos, independientemente del número de módulos en cada fila, y sin tener en cuenta si son o no del mismo tipo.

Por diferenciarse los distintos tipos de métodos únicamente en el número de posiciones de los distintos módulos que van a formar parte de la memoria del sistema computador, cabría pensar en la existencia de tantos métodos de diseño (sistemáticos) como módulos de desigual número de posiciones se estén empleando en el sistema para la implementación de la memoria principal. Esto, aunque sería en principio posible, no tendría ningún sentido ya que, si bien los dos primeros métodos marcan un carácter perfectamente diferenciado, el tercero es el caso que podría denominarse mixto, por emplear las dos filosofías de diseño conjuntamente, teniéndose que englobar cualquier

otro método que utilizara cualquier otro tamaño de módulo existente en el grupo de diseño mixto, y resultando peor solución (necesitaría mayor circuitería lógica) que la correspondiente al tercer método comentado, por no acoger las ventajas de los dos primeros con el grado más adecuado; el mejor grado de utilización será, para estos métodos considerados como mixtos, el correspondiente al seleccionado como tercer método (los ejemplos que se consideren lo justificarán).

Para la aplicación de cualquiera de los métodos, será necesario que se suministren los siguientes datos:

1. Tipo de mapa de memoria: **común o no común**
2. Número de bits del bus de direcciones del sistema: **m**

En caso de haberse especificado en el apartado anterior mapa de memoria no común, se tendrá que especificar también:

- 2'. Número de líneas del bus de direcciones empleadas para la selección de una posición en el mapa de E/S, y qué líneas de ese conjunto, las de mayor o menor peso: **$m_{E/S}$** y subconjunto **$A_{m-1} A_{m-2} \dots A_{m-m_{E/S}}$** ó subconjunto **$A_{m_{E/S}-1} A_{m_{E/S}-2} \dots A_0$**
3. Número de bits del bus de datos del sistema: **n**
4. Número de posiciones de la memoria dedicadas a ROM, RAM, EPROM, NOVRAM, módulos de E/S (con sus tipos), etc.: **$2^{m_{RAM}}$, $2^{m_{ROM}}$, $2^{m_{EPROM}}$, $2^{m_{NOVRAM}}$, $2^{m_{E/S1}}$, $2^{m_{E/S2}}$** , etc.
5. Organización de cada uno de los tipos de módulos de memoria y de E/S a emplear en el sistema: **$p_{RAM} Kxq_{RAM}$, $p_{ROM} Kxq_{ROM}$, $p_{EPROM} Kxq_{EPROM}$, $p_{NOVRAM} Kxq_{NOVRAM}$, $p_{E/S1} Kxq_{E/S1}$, $p_{E/S2} Kxq_{E/S2}$** , etc.
6. Distribución de la memoria en el mapa: **conjuntos de direcciones asociadas a los distintos tipos de memoria**

Una vez se han especificado todos los datos anteriores, hay que comprobar si es posible la realización de la denominada decodificación incompleta. En cuyo caso, sólo si es posible aplicar la decodificación incompleta, el usuario tendrá que decidir qué decodificación aplicar:

7. Decodificación **completa** o **incompleta**. Siendo necesario además, en caso de elección de decodificación incompleta, definir el **número de bits a eliminar** (por supuesto siempre dentro del rango permitido)

Cuando tenemos un sistema con mapa de memoria no común, la posibilidad o no de decodificación incompleta hay que tratarla de forma independiente para los mapas de memoria y de E/S.

Una vez decidida decodificación completa o incompleta, y el número de bits a eliminar en el caso posible, se conoce el espacio real de memoria libre direccionable por el sistema computador. Hay que preguntar el fin de este espacio libre:

8. Hay que preguntar si estas posiciones de memoria libre deben contemplarse como una **futura ampliación** o **no**; y, en el caso de contemplarse como ampliación, hay que especificar qué **tipo de módulos** ocuparán en un futuro esas posiciones o parte de esas posiciones.

Para mostrar el método de diseño de forma general vamos a empezar mostrando las líneas del procesador que se van a considerar:

- a) **Mapa de memoria común:** Bus de direcciones ($A_{m-1}, A_{m-2}, \dots, A_1, A_0$); Bus de Datos ($D_{n-1}, D_{n-2}, \dots, D_1, D_0$); Indicación de operación de lectura (**RD#**); Indicación de operación de escritura (**WR#**).
- b) **Mapa de memoria no común:** Bus de direcciones ($A_{m-1}, A_{m-2}, \dots, A_1, A_0$); Bus de Datos ($D_{n-1}, D_{n-2}, \dots, D_1, D_0$); Indicación de operación de lectura (**RD#**); Indicación de operación de escritura (**WR#**); Indicación de operación en memoria (**MREQ#**); Indicación de operación de E/S (**IOREQ#**).

En la [Figura 3.6](#) se muestra el esquema de las líneas de los procesadores correspondiente a ambos tipos de mapas.

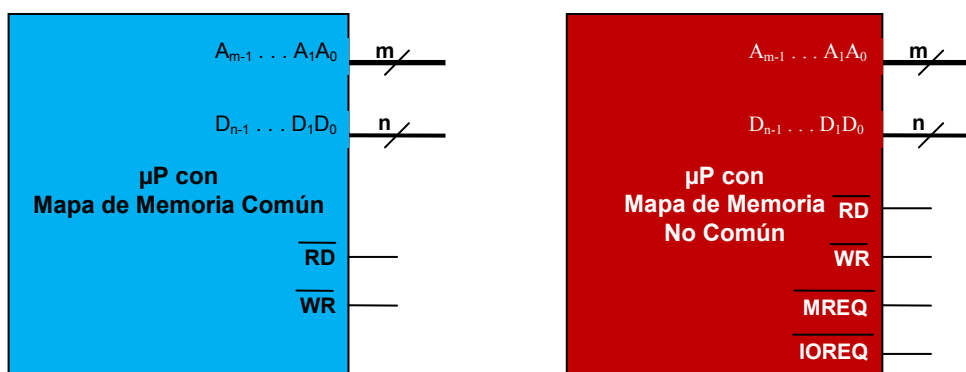


Figura 3.6. Esquema de líneas de microprocesadores con mapa de memoria común y no común.

Si se tiene **mapa de memoria común**, existe un único espacio de direcciones, que tendrán que compartir los dispositivos de memoria y de E/S. En este caso hay que considerar conjuntamente los tamaños de todos los dispositivos, de memoria y de E/S.

Cuando se habla de un sistema con **mapa de memoria no común**, se consideran dos espacios de direcciones, uno exclusivo para los dispositivos de memoria y otro exclusivo para los dispositivos de E/S. Para el mapa de memoria hay que considerar los tamaños de los dispositivos de memoria y, para el mapa de E/S hay que considerar exclusivamente los tamaños de los dispositivos de E/S. El diseño se hará de forma independiente, mapa de memoria y mapa de E/S. Una vez se han definido independientemente cada mapa de direcciones, así como se han diseñado sus sistemas correspondientes, se unen haciendo uso de las líneas *MREQ#* e *IOREQ#*; la *MREQ#* habilitará el decodificador primero correspondiente a la decodificación de los módulos de memoria y, de la misma forma, *IOREQ#* habilitará el primer decodificador que resuelve la función de decodificación de los módulos de E/S.

Haciendo referencia a un único espacio de direcciones, hay que empezar, una vez se conocen todos los datos mostrados anteriormente, seleccionando el parámetro que hemos denominado **tamaño básico (T_{bas})**. El tamaño básico T_{bas} está relacionado con el número de posiciones de memoria de cada módulo de forma que coincide o es un múltiplo o un divisor, potencia de dos.

Una vez seleccionado el tamaño básico T_{bas} , se divide el espacio de direcciones existente, de 2^m posiciones o de $2^{mE/S}$, en trozos de tamaño T_{bas} (el tamaño básico correspondiente al mapa de direcciones al que se esté haciendo referencia). A continuación, con el fin de resolver la función de decodificación, se define un **decodificador teórico** con tantas salidas como trozos de igual tamaño (T_{bas}) existan en el mapa de direcciones. El decodificador teórico tendrá un número de entradas d tal que 2^d es igual al número de salidas del decodificador.

Si a las entradas del decodificador teórico llevamos las líneas de más peso del bus de direcciones que correspondan (d líneas), cada una de las salidas de dicho decodificador se activará cuando las direcciones, referenciadas con el bus de direcciones, contengan la combinación de los d bits correspondiente. Y, como consecuencia, cada una de esas salidas del decodificador se van a utilizar para la habilitación de el/los módulo/s asociados a dicho trozo de memoria de tamaño T_{bas} . Se va a concretar y aclarar con un ejemplo.

Supongamos un bus de direcciones de 16 líneas ($A_{15} \dots A_1 A_0$) (2^{16} posiciones) y supongamos un tamaño básico T_{bas} de 8 Kposiciones (2^{13} posiciones). El número de trozos de tamaño T_{bas} serán $2^{16}/2^{13} = 2^3 = 8$. Por lo tanto, el decodificador teórico a emplear tendrá 8 salidas y, como consecuencia 3 entradas ($2^3 = 8$). Si a las líneas de entrada del decodificador teórico llevamos las 3 líneas de más peso del bus de direcciones de forma ordenada, cada salida del decodificador se habilitará según qué combinación esté presente en las líneas de entrada: si la combinación presente es 000 se habilitará la salida 0, si es 001 la salida 1, ..., y así sucesivamente. La característica común a todas las direcciones de memoria que pertenecen al primer trozo de 8 Kposiciones es que los 3 bits de más peso $A_{15}A_{14}A_{13}$ tienen el valor 000; por lo tanto, al haberse definido las 3 entradas del decodificador con los bits $A_{15}A_{14}A_{13}$ del bus de direcciones, la salida 0 del decodificador se empleará para activar el/los módulo/s de memoria asociados a ese trozo de tamaño 8 K. En la [Figura 3.7](#) se muestra gráficamente el ejemplo planteado.

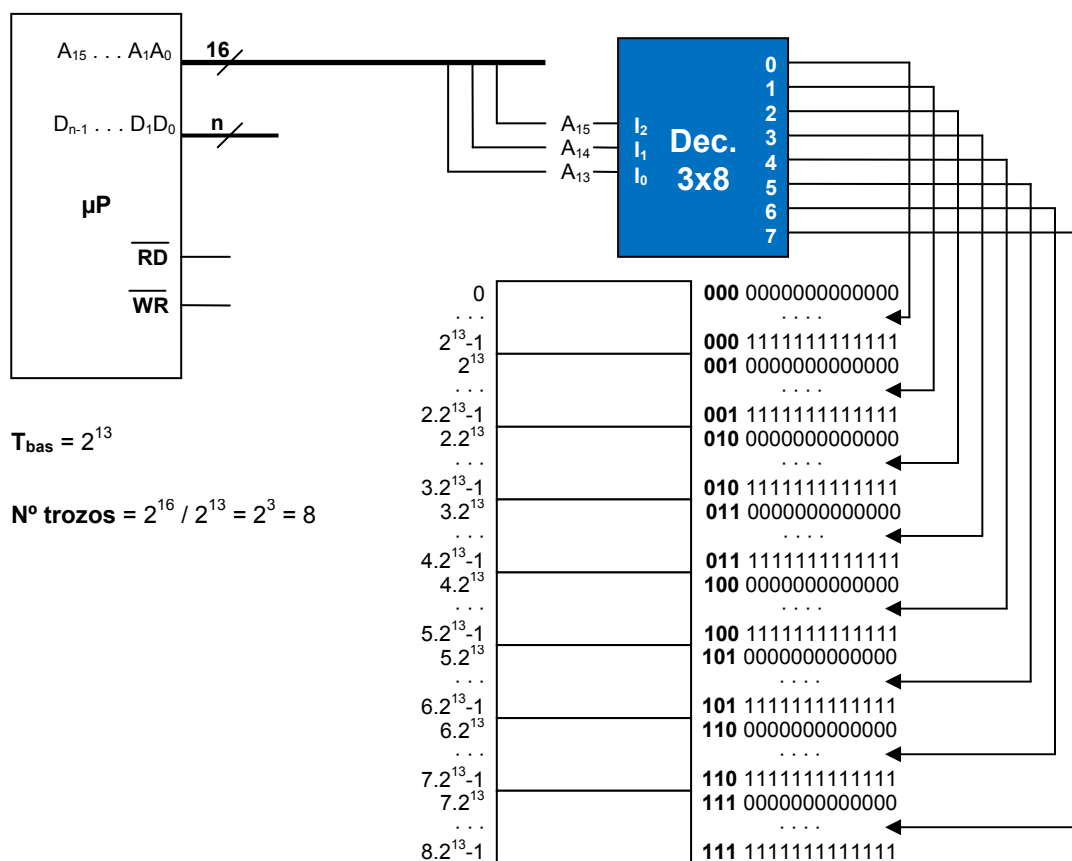


Figura 3.7. Esquema representativo de ejemplo de un decodificador de direcciones.

En el esquema mostrado, cada línea de salida del decodificador nos permite habilitar los módulos que estén asociados al trozo del mapa de memoria correspondiente; este trozo es de 8 Kposiciones. En principio nos podemos encontrar tres situaciones:

- I. El tipo de módulo de memoria que va a ocupar el trozo del mapa de memoria considerado tiene un número de posiciones que coincide con el tamaño básico considerado. En este caso, la línea de salida del decodificador que corresponde se llevaría directamente a la patilla de selección de chip (CS#) del módulo de memoria asociado. (Se consideran decodificadores con salidas activas a nivel bajo).
- II. El tipo de módulo de memoria a considerar tiene un número de posiciones mayor que el tamaño básico. En este caso, habrá que asociar al módulo de memoria más de un trozo del mapa de memoria y, como consecuencia hay que hacer uso de tantas líneas de salida del decodificador como trozos de tamaño básico del mapa ocupe. Para conectar las líneas de salida del decodificador correspondiente con la de selección del chip de memoria hay que resolver una función combinacional simple; se llevarán las salidas del decodificador correspondientes a las entradas de una puerta AND, y su salida se conectará a la entrada CS# de habilitación del módulo de memoria considerado.
- III. El tipo de módulo de memoria referido tiene un número de posiciones menor que el tamaño básico considerado. En este caso, si llevamos a la línea de habilitación CS# del chip de memoria directamente la línea de salida del decodificador asociada al trozo del mapa de memoria correspondiente, estaríamos asociando a cada posición real del módulo de memoria tantas posiciones lógicas del mapa como combinaciones distintas existan de las líneas del bus de direcciones que no se hayan tenido en cuenta para la selección de una posición de dicho módulo de memoria (se estaría hablando de decodificación incompleta). Si le queremos asociar una única posición lógica a cada posición real (decodificación completa), para este caso habrá que hacer una segunda decodificación; a este segundo decodificador se llevarán las líneas del bus de direcciones que aún no se hayan empleado.

En la [Figura 3.8](#) se muestra la resolución de las tres situaciones planteadas, para el ejemplo que se ha considerado en la [Figura 3.7](#) y seleccionando para cada caso los siguientes tamaños de los módulos de memoria:

- I. Módulo de memoria con 8 Kposiciones.
- II. Módulo de memoria con 16 Kposiciones.
- III. Módulo de memoria con 4 Kposiciones.

Por lo tanto, según la división del mapa de memoria en “*trozos de igual tamaño*” que hagamos nos vamos a encontrar con:

1. Módulo de memoria y/o de E/S con menor número de posiciones → Todas las posibilidades se resolverán según las situaciones I y II, directamente o con puertas AND.
2. Módulo de memoria y/o de E/S con mayor número de posiciones → Las distintas posibilidades se resolverán según las situaciones I y III, directamente o con una segunda decodificación.
3. Módulo de memoria y/o de E/S que predomine al considerar el número de filas que hay que emplear de los mismos, independientemente del número de módulos en cada fila, y sin tener en cuenta si son o no del mismo tipo → Las posibilidades que se plantean se resuelven según las situaciones I, II y III, directamente, con puertas AND o con una segunda decodificación.

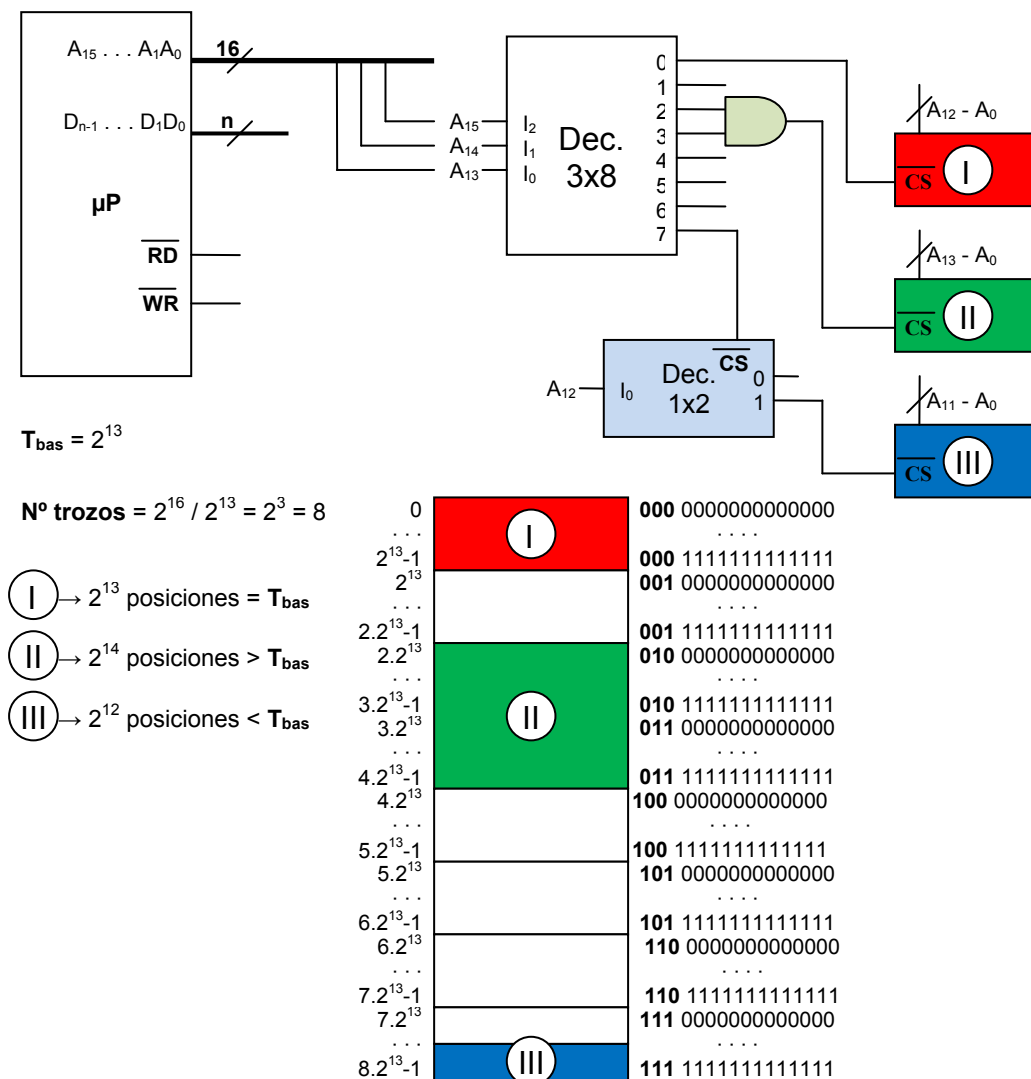


Figura 3.8. Ejemplo mostrando la decodificación de las direcciones para módulos de memoria de tamaños = T_{bas} , > T_{bas} y < T_{bas} .

Todo lo que se ha aplicado para un mapa de direcciones, si tenemos dos espacios de direcciones diferenciados (mapa de memoria no común) es aplicable de igual forma a ambos espacios de direcciones, resolviéndose cada uno de ellos de manera independiente según la metodología mostrada. La compatibilidad de ambos esquemas, cada uno de ellos asociado a un mapa, se resuelve haciendo uso de las señales MREQ# e IOREQ#; la primera se conectará a la línea de selección de chip (CS#) perteneciente al decodificador principal asociado a los módulos de memoria que ocupan posiciones en el espacio de direcciones de memoria, y la segunda se conectará a la línea de selección de chip (CS#) del decodificador principal asociado a los módulos de E/S que ocupan posiciones en el espacio de direcciones de E/S. Ver [Figura 3.9](#).

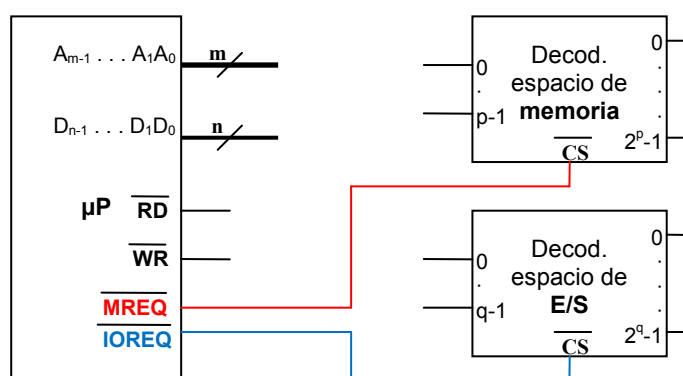


Figura 3.9. Conexión de las señales MREQ# e IOREQ#, en un sistema con mapa de memoria no común, a los decodificadores principales de ambos mapas.

3.3.RECURSOS PARA MEJORAR LAS PRESTACIONES DE LA MEMORIA PRINCIPAL

3.3.1.Introducción: aumento de la capacidad y aumento de la velocidad de acceso

La memoria principal es uno de los componentes básicos y esenciales de todo computador, es el área de trabajo del computador.

La memoria principal ha de cumplir con los siguientes requisitos:

1. Tener suficiente capacidad. Para poder albergar cualquier aplicación o aplicaciones.
2. Tener un tiempo de acceso a sus posiciones de memoria en consonancia con el tiempo de procesamiento por parte de la Unidad de Control, lo más cercano posible.

El nivel de jerarquía de MEMORIA VIRTUAL tiene como objetivo *aumentar la capacidad real de la memoria principal*.

El nivel de jerarquía de MEMORIA CACHE tiene como objetivo *incrementar la velocidad de los accesos*, aproximar lo máximo posible el tiempo de acceso a los tiempos de procesamiento en la Unidad Central de Proceso, cuando menos que sean de órdenes de magnitud similares.

3.3.2. Memoria Virtual

La *memoria virtual* permite a la CPU la ejecución de programas que no están totalmente en Memoria Principal, de forma transparente. Esta necesidad puede venir impuesta por la capacidad excesivamente grande de una aplicación o por el número elevado de usuarios que estén ejecutando aplicaciones en un momento determinado en el mismo computador.

Se trata de introducir mecanismos de gestión que asignen espacio de memoria a los usuarios o aplicaciones a medida que lo van necesitando, buscando una ejecución eficiente de los programas.

Antes de aparecer el mecanismo de memoria virtual, los programas que necesitaban mucho espacio empleaban la técnica de los OVERLAYS (trozos en los que se divide un programa). Era necesario que el usuario diseñara un programa denominado CONTROLADOR DE OVERLAY (siempre debe estar presente en la memoria).

En 1961, un grupo de la Universidad de Manchester propone un método para realizar la técnica de los overlays de forma automática, liberando al programador de esta tarea. A este mecanismo se le conoce con el nombre de MEMORIA VIRTUAL.

Un sistema informático con memoria virtual pone a disposición del programador una enorme cantidad de memoria que, en realidad, reside en un dispositivo de almacenamiento masivo (perteneciente al nivel de memoria secundaria).

La implantación de la MEMORIA VIRTUAL conlleva:

- a) La incorporación de HARDWARE, que configura en el microprocesador la MMU (Memory Management Unit – Unidad de Manejo de Memoria).
- b) La incorporación de funciones especiales al SOFTWARE del Sistema Operativo. Este Sistema Operativo debe incorporar las siguientes funciones:
 - Poseer una lista completa con la descripción de todos los objetos residentes tanto en la memoria virtual como en la memoria principal.
 - Cargar objetos en la memoria principal en tiempo de ejecución.
 - Manejar de los espacios libres de la memoria principal y, cuando no existan, decidir los objetos que no se prevé usar posteriormente para devolverlos a la memoria virtual y así poder utilizar el espacio que ocupaban.

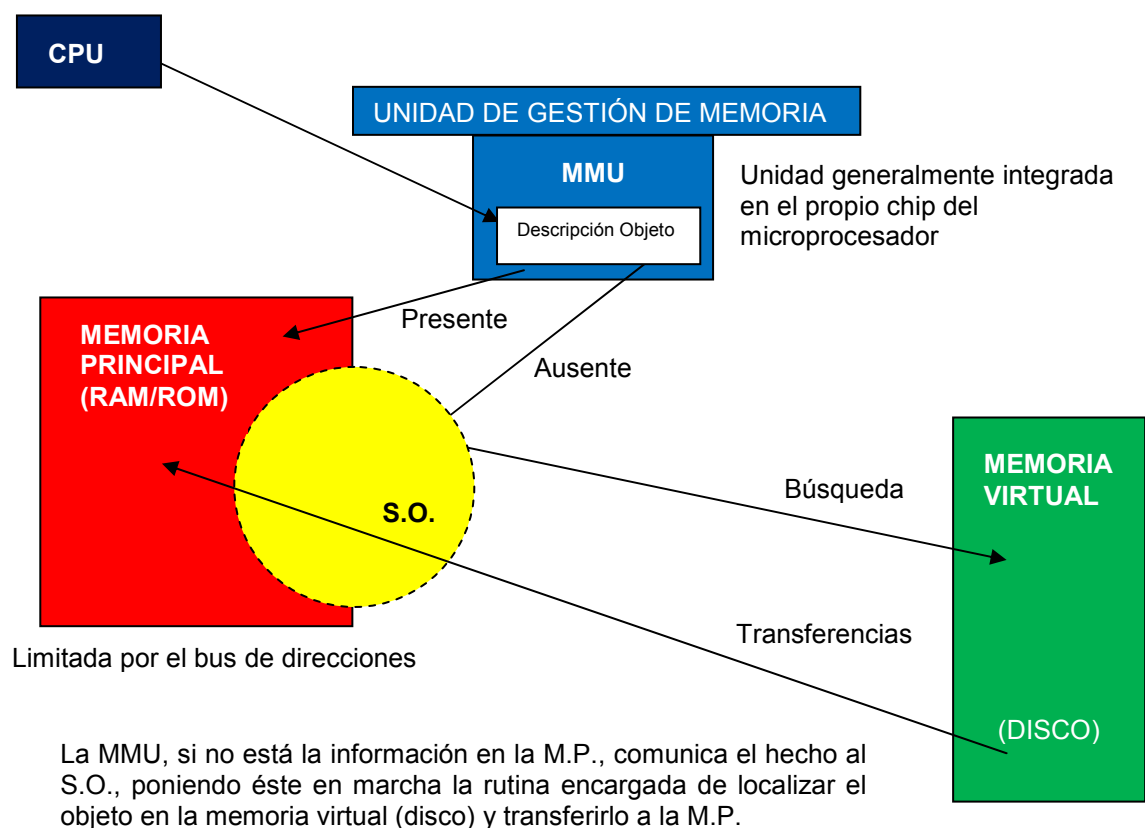


Figura 3.10. Sistema con Memoria Virtual.

3.3.3. Memoria caché o inmediata

Se trata de una MEMORIA AUXILIAR, de pequeña capacidad y alta velocidad, que se añade a la memoria principal para acelerar su funcionamiento.

En la memoria caché se guarda la información correspondiente a las posiciones de memoria principal que más frecuentemente se prevé se van a usar. *El éxito de la caché está en anticiparse a las peticiones del Procesador.*

La memoria caché y principal se dividen en *bloques o líneas de unos pocos bytes*, que hacen la función de páginas. Estos bloques son las unidades mínimas de transferencia entre la Memoria Caché y la Memoria Principal.

El empleo apropiado de la caché puede dar una tasa de aciertos superior al 90%. De esta forma, se acerca mucho el tiempo de acceso a la información al correspondiente tiempo de acceso a memoria caché.

TIPOS DE CONEXIÓN DE LA CACHÉ

- CONEXIÓN EN SERIE

El microprocesador ve a la memoria a través de la caché. Todas las peticiones se envían en primer lugar a la caché y, si da lugar, desde ésta a la M.P.

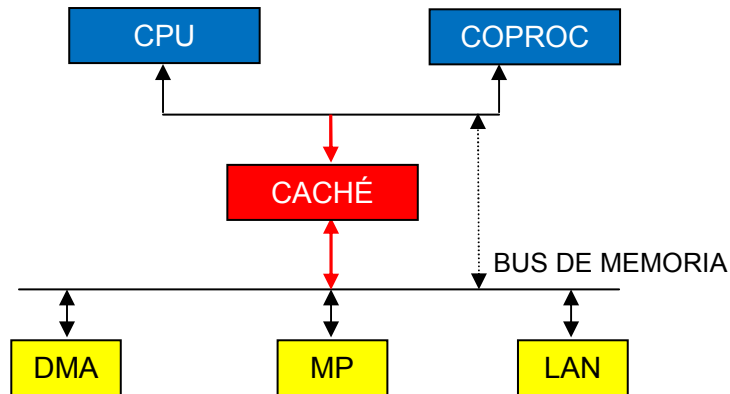


Figura 3.11. Conexión en serie.

Ventajas:

- Disminuye el número de peticiones de acceso a la M.P. → Además de reducir el tiempo de acceso, *reduce el tiempo de utilización del bus*.
- La CPU puede estar trabajando con la caché y otros maestros del bus del sistema pueden estar accediendo a la M.P. (Importante en sistemas multiprocesadores).

Inconvenientes:

- No se puede quitar del sistema.
- Hay penalizaciones en caso de ausencia.
- Hay dos buses independientes → *Complicación del hardware*.

- CONEXIÓN EN PARALELO

Todas las peticiones de memoria llegan simultáneamente a la M.P. y a la caché. Si ésta tiene el dato pedido, se lo suministra al microprocesador y envía una señal a la M.P. para que aborte el ciclo iniciado.

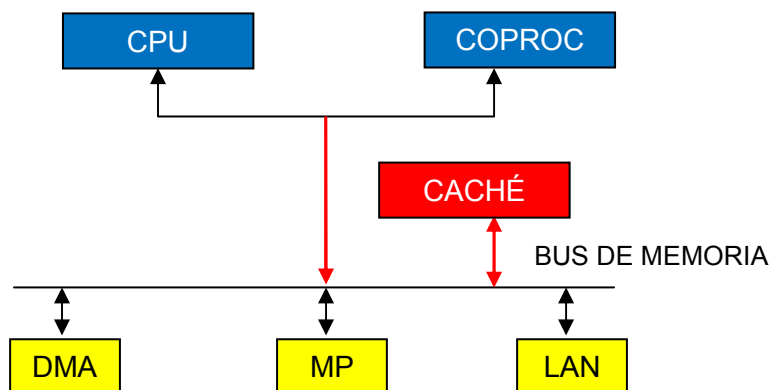


Figura 3.12. Conexión en paralelo.

Ventajas:

- Se puede quitar del sistema sin tener que hacer modificación alguna.
- No hay penalizaciones en caso de ausencia.
- Más sencilla de integrar (no hay dos buses independientes).

Inconvenientes:

- Todas las peticiones provocan accesos a memoria principal, aún estando la información en la caché.
- No se reduce la utilización del bus de memoria (los sistemas multiprocesadores no se benefician de este tipo de estructura).

OPCIONES PARA MEJORAR EL ACCESO A LA INFORMACIÓN DE LA CACHÉ

Se trata de reducir la penalización por fallo (más importante aún si cabe en la conexión serie que en la paralela). Hay dos formas de reducirla:

- **Incrementar el ancho físico o lógico del sistema de memoria.** Las transferencias desde la Memoria Principal a la Memoria Caché se realizan implicando a varias palabras de memoria a la vez, por permitirlo el ancho de los caminos de comunicación entre ambas.

- **Memoria entrelazada.** La Memoria Principal se construye con distintos bancos de memoria capaces de operar independientemente unos de otros; de esta forma, aunque las transferencias de información entre Memoria Principal y Memoria Caché tienen que ser en serie, las búsquedas de información en bancos de memoria independientes pueden solaparse en el tiempo.

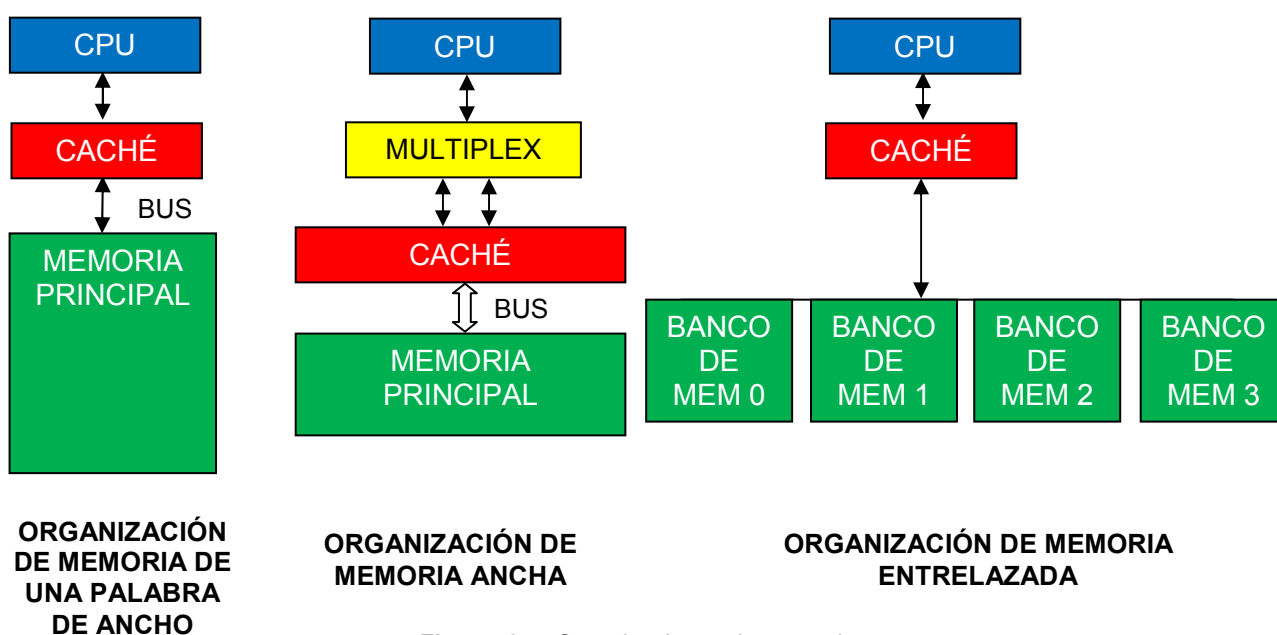


Figura. 3.4. Organizaciones de memoria.

3.3.4. Protección de la memoria principal

Una de las funciones primordiales del Sistema Operativo (S.O.) es la protección. Éste ha de garantizar la confidencialidad de la información de los usuarios y ha de asegurar que unos procesos no interfieren con otros, destruyendo, copiando o modificando su información. Ahora bien, el S.O. consiste en un conjunto de programas, por tanto cuando está ejecutando un proceso de usuario, no está ejecutando el S. O.. Esto significa que la vigilancia sobre los procesos de usuario, para comprobar que no realizan operaciones prohibidas, no puede hacerla el S.O., sino que esta vigilancia debe hacerse mediante circuitos especiales incluidos en el procesador y la memoria. La secuencia por tanto, es la siguiente:

- El hardware vigila las acciones del proceso de usuario. Si detecta el intento de una operación prohibida, no la ejecuta y genera una excepción.
- El S.O. entra en ejecución como consecuencia de la excepción. Analiza la excepción y decide la acción a tomar, que muchas veces será terminar el proceso.

En esta sección se verán los mecanismos de protección del procesador y de protección de la memoria.

PROTECCIÓN DEL PROCESADOR

Un procesador de propósito general cuenta con dos niveles de ejecución: el *nivel de usuario* y el *nivel de núcleo*. El nivel de núcleo es el más permisivo y está previsto para el S.O., mientras que el nivel de usuario tiene más restricciones y está previsto para los programas de usuario. El nivel de ejecución está indicado mediante uno o varios bits del registro de estado, ubicados en la zona de núcleo.

En base al esquema mostrado en la [Figura 3.13](#), los elementos de almacenamiento (registros, mapa de memoria y mapa de E/S), se observa que el nivel de usuario está limitado en comparación con el nivel de núcleo. Desde el nivel de usuario, no se puede acceder a todo el registro de estado, a los registros especiales, al mapa de E/S ni a todo el mapa de memoria. Además, tampoco es posible manejar todo el repertorio de instrucciones.

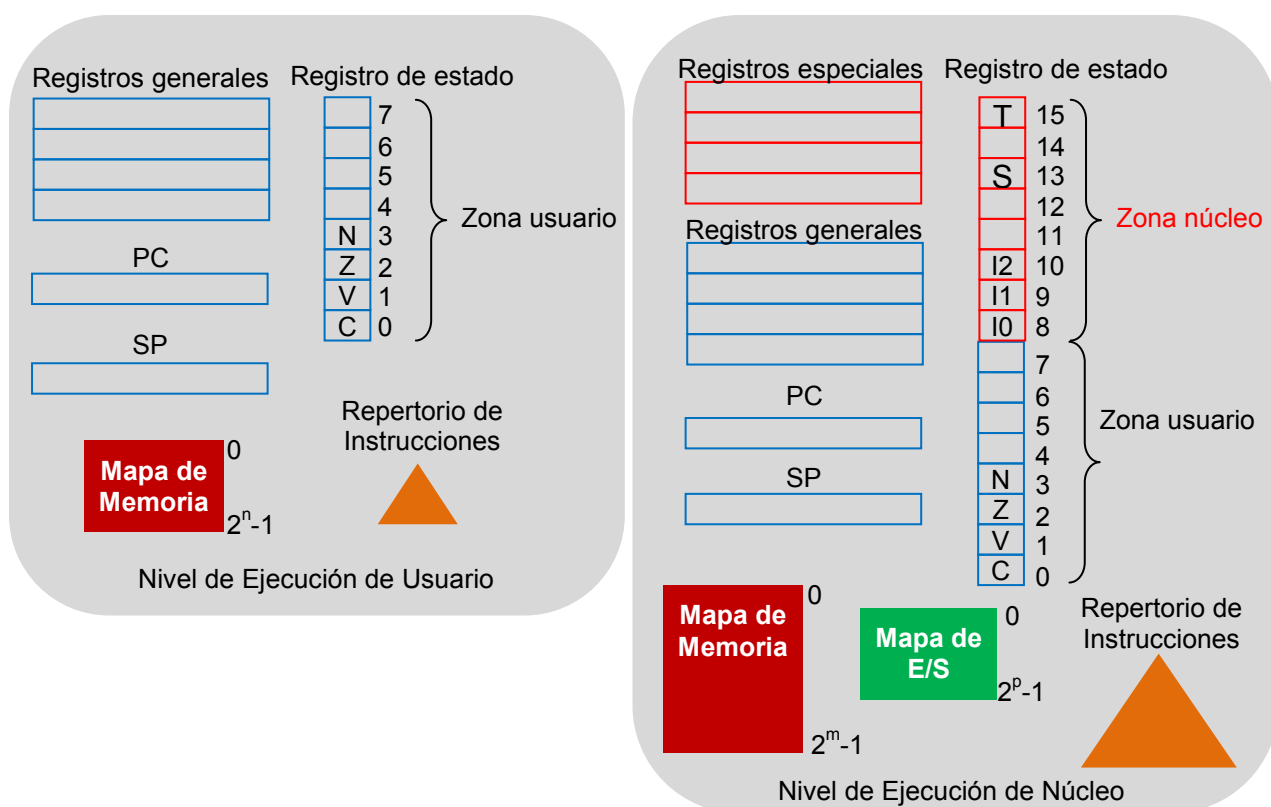


Figura 3.13. Niveles de ejecución de usuario y de núcleo.

Cuando la Unidad de Control decodifica cada instrucción máquina, comprueba que ésta contiene un código de operación permitido por el nivel de ejecución indicado por el registro de estado. Además comprueba que la instrucción no utiliza ningún elemento de almacenamiento prohibido. Si todo es correcto, ejecuta la instrucción. Si detecta un posible error de permisos genera una excepción. Esta excepción provoca que se ejecute el S.O., quien decide la acción a tomar por el proceso causante del error.

Desde el nivel de usuario, no se tiene acceso a los bits del registro de estado que establecen el nivel de ejecución, por lo que ningún programa de usuario se podrá ejecutar en nivel de núcleo; quedando así garantizada la protección del procesador.

Otra cuestión a resolver desde el punto de vista de la protección es, ¿cómo puede el S.O. ejecutar en modo núcleo?. Para ello se dota de un mecanismo de interrupción basado en la siguiente secuencia ([Figura 3.14](#)):

- Se produce una interrupción.
- La Unidad de Control la acepta y:
 - o Lee el vector de interrupción.
 - o Pone el computador en nivel de ejecución de núcleo.

- Guarda el contenido del registro contador de programa PC y del registro de estado en la zona de memoria denominada pila.
- Modifica el registro PC mediante una indirección en la tabla de interrupciones con el vector de interrupción.

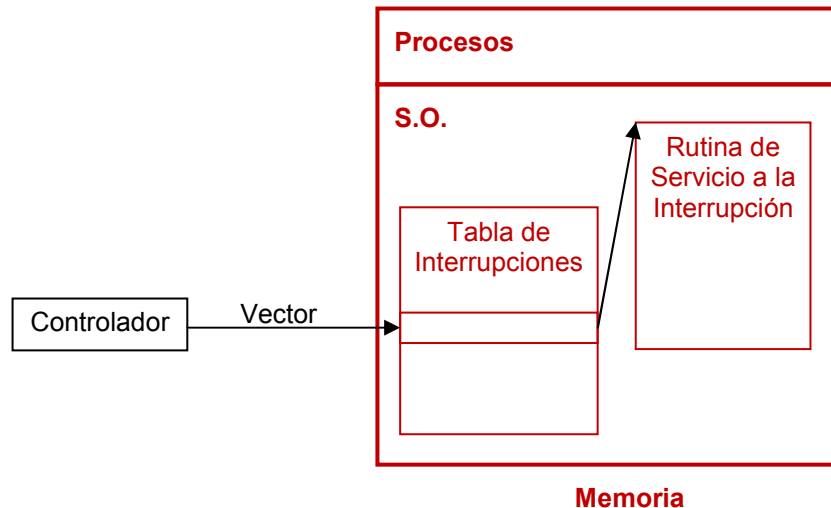


Figura 3.14. Acceso a la Rutina de Servicio a la Interrupción.

Observar que de nada serviría esta secuencia si el programa de usuario fuese capaz de modificar la tabla de interrupciones, y poner un valor de su espacio de direcciones, o si la tabla de interrupciones tuviese alguna entrada que apuntase al espacio de memoria reservado a los procesos. En estos casos, al producirse una interrupción, se pasaría a ejecutar un programa de usuario con nivel de ejecución de núcleo, y que tendría todos los permisos para hacer lo que quisiera. Por tanto, el S.O. debe ser el único encargado de tratar las interrupciones, y la única forma de pasar del nivel de usuario a nivel de núcleo es mediante una interrupción. Es decir, si el S.O. debe ejecutar en modo núcleo, la única forma de activarlo debe ser mediante una interrupción.

PROTECCIÓN DE LA MEMORIA

La protección de la memoria es necesaria para evitar que un proceso, de forma inadvertida o por malicia, trate de acceder a una zona de memoria que no le ha sido asignada. La protección de la memoria es especialmente importante en sistemas multiusuarios, ya que es imprescindible evitar interferencias entre los usuarios y entre estos con el S.O.

Como se ha visto antes, la Unidad de Control puede hacer una primera protección de memoria, prohibiendo que los programas de usuario accedan a una parte del mapa de memoria, que queda así reservada para el S.O. Esta solución es útil para evitar

interferencias entre los procesos de usuario y el S.O. A continuación se verán mecanismos para proteger a los procesos de usuario entre sí.

- PROTECCIÓN DE MEMORIA REAL

Los procedimientos que se verán en este caso se emplearán cuando el computador no tenga sistema de memoria virtual.

- **Registros de borde:** El mecanismo de registros de borde puede emplearse cuando se hace asignación de memoria principal continua a los procesos. Cada proceso tiene asignada una franja de memoria, y la protección viene garantizada por dos registros de borde que contienen los valores de las direcciones de memoria que determinan el borde de la franja del proceso activo. El comparador verifica cada acceso a memoria realizado por el proceso, y en caso de fallo, el comparador invalida el acceso a memoria y genera una excepción, pasando a ejecutarse el S.O.. Es el S.O. el encargado de cancelar el proceso y de generar los correspondientes mensajes de error.

La función de cargar los registros de borde con los valores correspondientes según el proceso, corresponde al S.O.. Esta carga se realiza mediante instrucciones privilegiadas, que los procesos comunes no pueden ejecutar, para evitar posibles casos de autoampliación de la zona de memoria asignada.

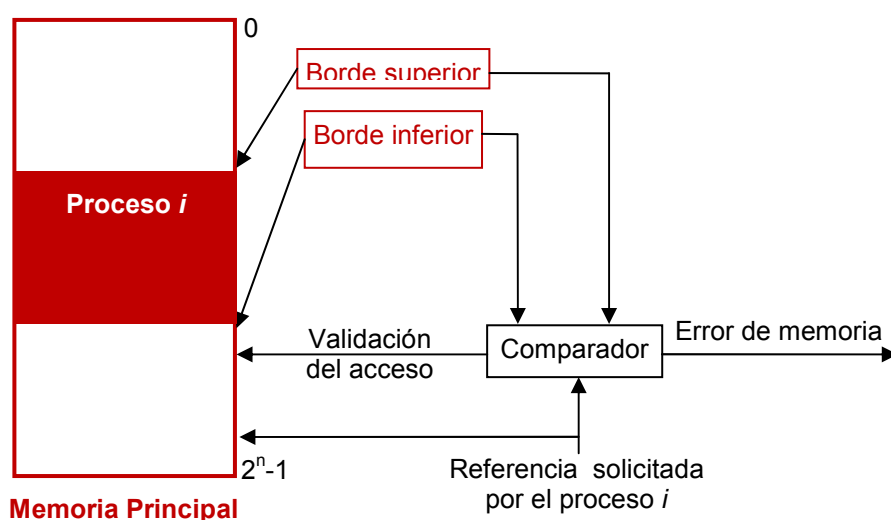


Figura 3.15. Protección de memoria con Registros de Borde.

- **Cerrojo y clave:** El mecanismo de *cerrojo* y *clave* permite establecer niveles de protección. Estos niveles de protección pueden ser de gran interés en los siguientes casos:
 - Acceso de lectura y escritura (por ejemplo a zonas de datos propias).

- Acceso de lectura (por ejemplo a programas propios o a datos de otros usuarios).
- Acceso de escritura (por ejemplo de un usuario en un buffer de E/S).
- Acceso de ejecución (por ejemplo de un programa propiedad de otro).

En este caso, la memoria principal se divide en bloques, de tamaño fijo o variable, a los que se asignan unos códigos de acceso para los distintos niveles de protección previstos. Para que un proceso pueda emplear uno de estos bloques, debe conocer la clave correspondiente.

Dado que la comprobación de la protección hay que hacerla en cada acceso a memoria, se debe establecer un procedimiento *hardware* para realizarla, que suele consistir en lo siguiente:

- Se divide la memoria en bloques de tamaño fijo, a los que se asigna una serie de bits de cerrojo.
- Cada proceso tiene asignado un código que el sistema operativo se encarga de introducir en un registro específico para ello, antes de darle el control.
- En cada acceso, se compara el código del proceso con el cerrojo del bloque direccionado. En caso de coincidencia del cerrojo con el código del proceso y el tipo de acceso solicitado, se atiende la petición.
- En caso contrario, se rechaza el acceso y se genera una excepción de intento de acceso a memoria no permitida. Como con toda interrupción, pasa a ejecutarse el S.O. que decide la acción a tomar con el proceso.

- PROTECCIÓN DE MEMORIA VIRTUAL

Cuando se dispone de memoria virtual, pueden utilizarse los mecanismos propuestos para memoria física, aprovechándose para ello las posibilidades brindadas por el mecanismo de traducción de dirección propio del sistema de memoria virtual. La partición natural de la memoria en páginas o segmentos (cuando existe memoria virtual), permite que la protección se haga a este nivel.

Dado que el aislamiento entre los procesos no puede ser total, puesto que a veces es necesario que estos compartan información, se tratarán simultáneamente los métodos de protección y los mecanismos de comunicación interproceso mediante memoria compartida.

- **Mapa de memoria virtual independiente por proceso:** Una alternativa extrema de protección, consiste en dotar a cada proceso de un mapa virtual independiente, y hacer que el acceso a mapas diferentes que no comparten direcciones, sea

completamente imposible ([Figura 3.16](#)). Esta independencia de mapas virtuales se consigue haciendo que cada proceso tenga su propia tabla de páginas y que el S.O. cambie el contenido de la tabla destinada a albergar los elementos de la tabla de páginas que han sido empleados recientemente (esta tabla se denomina *TLB Translation Lookaside Buffer*), cuando da el control a un proceso, o que la tabla TLB guarde referencia del proceso.

La necesaria comunicación entre procesos puede hacerse mediante una parte especial del mapa de memoria del S.O., supervisando éste todos los accesos a esa zona común. El inconveniente de esta solución es que es muy costosa en tiempo de procesador, ya que cada acceso a la memoria requiere activar el S.O.

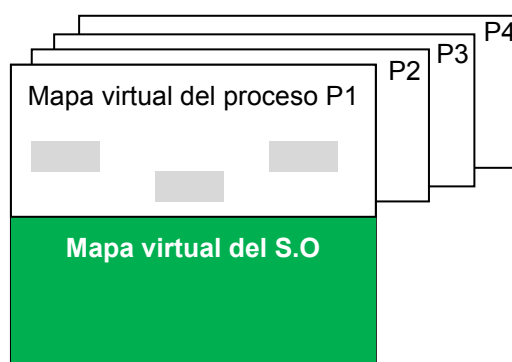


Figura 3.16. Espacios virtuales independientes.

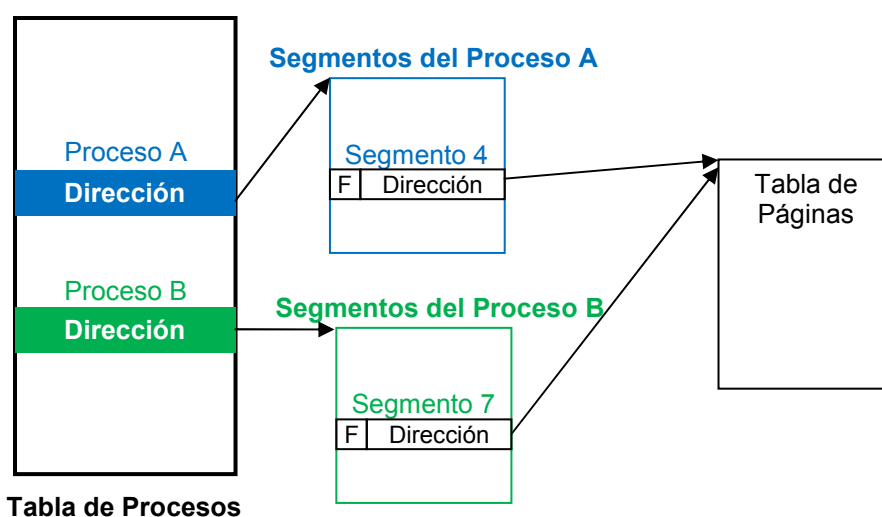


Figura 3.17. El proceso A comparte el segmento 4 con el segmento 7 del proceso B.

También se puede compartir memoria virtual simplemente proyectando segmentos de los mapas de varios procesos sobre la misma memoria física. Para el caso de memoria virtual segmentada paginada, esto se reduce a hacer que los descriptores de todos estos segmentos apunten a la misma tabla de páginas. Todos los procesos

comparten entonces esa tabla de páginas y, por tanto, los marcos de página asignados ([Figura 3.17](#)).

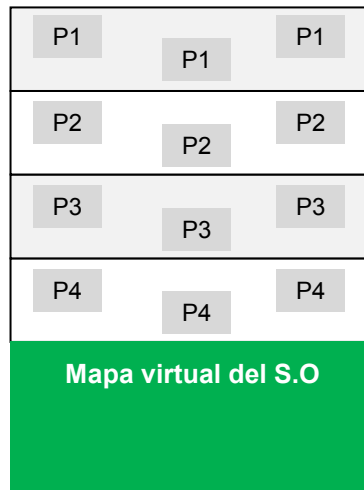


Figura 3.18. *Espacio virtual único.*

- **Mapa de memoria virtual común** ([Figura 3.18](#)): En este caso, la protección de la memoria puede realizarse por segmento o por página, según el caso. La técnica generalmente empleada es la de cerrojo y clave, que se aplica a los segmentos o páginas, añadiendo la información pertinente en cada elemento de la correspondiente tabla de traducción. También es conveniente incluir esta información en la tabla TLB, para que las traducciones estén convenientemente validadas. Este mecanismo permite proteger y compartir el nivel de operación que se desee, no siendo necesario emplear otros métodos para compartir información.

3.4.DISPOSITIVOS DE ALMACENAMIENTO

INTRODUCCIÓN

En esta sección se va a realizar una revisión de los dispositivos de memoria empleados para implementar dos de los niveles de memoria vistos en la sección 3.1, como son el nivel de Memoria Principal y el nivel de Memoria Secundaria.

DISPOSITIVOS DE ALMACENAMIENTO DE L/E PARA MEMORIA PRINCIPAL. DISPOSITIVOS DE MEMORIA COMERCIALES

En esta sección se van a comentar los siguientes dispositivos de memoria:

- **RAM** (Random Access Memory): Si se especifica de esta forma se refiere a la SRAM.
- **DRAM** (Dynamic RAM).

- **SRAM** (Static **RAM**): Son 5 veces más rápidas que las **DRAM**, 2 veces más grandes que las **DRAM** y 2 veces más caras que las **DRAM**. La **Memoria Caché** se construye con **SRAM**.
- **SIMM** (Single In-Line Memory Module). Un **SIMM** típico consiste en varios chips de **DRAM** instalados en una pequeña **placa de circuito impreso o PCB**, la cual calza en un receptáculo **SIMM** en la placa del sistema. Los **SIMMs** vienen con varios formatos, incluso los de 30 y de 72 contactos. Los de 72 contactos nos dan acceso de 32 en 32 bits.
- **DIMM** (Dual In-Line Memory Module). Se parecen bastante a la memoria de tipo **SIMM**. Al igual que los **SIMMs**, la mayoría de los **DIMMs** se instalan verticalmente en los conectores de expansión. La diferencia principal entre los dos consiste en que, en un chip **SIMM**, los contactos de cada fila se unen con los contactos correspondientes de la otra fila para formar un solo contacto eléctrico; en un chip **DIMM**, los contactos opuestos permanecen eléctricamente aislados para formar dos contactos separados. Los módulos **DIMMs** se utilizan frecuentemente en las configuraciones que brindan soporte para un bus de memoria de 64 bits o más amplio.
- **SO DIMM** (Small Outline **DIMM**). Otro tipo de memoria que se usa comúnmente en los ordenadores “laptop” y portátiles se llama **Small Outline DIMM** (de contorno pequeño) o **SO DIMM**. Un **DIMM** de contorno pequeño es como un **SIMM** de 72 contactos en un paquete de dimensiones reducidas, sin embargo, existen algunas diferencias técnicas importantes. El **DIMM** de contorno pequeño y el **SIMM** de 72 contactos se diferencian en la disposición de los contactos. Tiene aproximadamente la mitad de longitud de un **SIMM** típico de 72 contactos.
- **FPM DRAM** (Fast Page Mode **DRAM**). Paginación rápida. Permite que la CPU obtenga datos nuevos en la mitad de tiempo del acceso normal siempre y cuando se encuentre en la misma página que la solicitud anterior. Accede más rápido a los datos en la misma fila o filas, donde ya haya obtenido datos anteriormente.
- **EDO RAM** (Extended Data Output **RAM**). Modelo avanzado de la **FPM DRAM**, a veces se le denomina modo **HIPER PAGE**. Realiza el mismo proceso anterior también con columnas. Requieren diseño especial. Permite incrementar la velocidad de acceso entre un 10 y un 30% con respecto a la **FPM DRAM**, y viene a costar sólo un 5% más.

- **BEDO RAM** (Burst Extended Data Output **RAM**). Es una memoria **EDO RAM** que contiene una entrada pipeline y un contador de ráfagas (“burst”) de 2 bits. Obtiene un aumento de rendimiento de un 100% respecto a la **FPM DRAM** y de un 33 a un 50% respecto a la **EDO RAM**.
- **SDRAM** (Synchronous **DRAM**). Se sincroniza con el reloj del sistema (en el flanco positivo) que controla la CPU. Al estar sincronizado con la CPU, elimina los tiempos de espera y hace el proceso de recuperación de datos más eficaz. Se presentan en módulos tipo DIMM SDRAM a 3.3 V y con 168 contactos (PC66, PC100, PC133).
- **SLDRAM** (SyncLinck **DRAM**). Si la **SDRAM** emplea cuatro bancos de memoria, ésta emplea 16. Es la competencia para las memorias **RAMBUS**. Desarrollada por un consorcio de fabricantes de computadoras llamado *Consortio Synclinc*.
- **SGRAM** (Synchronous Graphic **RAM**). (Memoria gráfica síncrona de acceso aleatorio). Se utiliza en adaptadores de vídeo y aceleradoras gráficas. Al igual que la **SDRAM** puede sincronizarse con el reloj del bus de la CPU hasta velocidades de 100 MHz. Además, la **SGRAM** emplea algunas otras técnicas, como escritura con máscara y bloques de escritura, para aumentar el ancho de banda para funciones gráficas intensivas. Al contrario que la **VRAM** y la **WRAM**, la **SGRAM** es de puerto único; sin embargo, puede abrir dos páginas de memoria a la vez, lo que simula la naturaleza de puerto doble de otras tecnologías de **RAM** vídeo.
- **RDRAM** (Rambus **DRAM**). Es sumamente rápida y requiere cambios significativos en los sistemas de memoria. Alcanza velocidades aproximadamente cuatro veces más rápida que la **DRAM** normal. Es muy costosa. Tiene un bus de datos de 16 bits, pero funciona a velocidades mucho mayores (266, 356 y 400 MHz). Además es capaz de aprovechar cada señal doblemente, de forma que en cada ciclo de reloj envía 4 bytes en lugar de 2. Algunos tipos de **RDRAM**:
 - PC600→ a 266 MHz físicos: $2 \times 2 \text{ bytes/ciclo} \times 266 \text{ MHz} = 1.06 \text{ GB/s}$.
 - PC700→ a 356 MHz físicos: $2 \times 2 \text{ bytes/ciclo} \times 356 \text{ MHz} = 1.42 \text{ GB/s}$.
 - PC800→ a 400 MHz físicos: $2 \times 2 \text{ bytes/ciclo} \times 400 \text{ MHz} = 1.6 \text{ GB/s}$.
- **Concurrent RDRAM**.
- **Direct RDRAM**. Pretende abrir un camino de datos más ancho para acelerar las transferencias de datos. Se implementa en placas muy similares a los actuales módulos **DIMM**.

- **MDRAM** (Multibank **DRAM**) (**DRAM** Multibanco). Abreviatura de Multibank **DRAM** (**DRAM** multibanco), una tecnología de memoria relativamente nueva desarrollada por MoSys Inc. La MDRAM utiliza pequeños bancos de **DRAM** (de 32 KB cada uno) en una matriz, donde cada banco tiene su propio puerto I/O que se comunica con un bus interno común. Gracias a este diseño, los datos pueden ser leídos o escritos simultáneamente a varios bancos, lo que la hace mucho más rápida que la DRAM convencional. Otra ventaja de la **MDRAM** es que la memoria puede ser configurada en incrementos menores, lo que reduce el coste de algunos componentes. Por ejemplo, es posible producir chips de **MDRAM** de 2.5 MB, que es lo que necesitan los adaptadores de vídeo con 24-bits de color a una resolución de 1024x768. Con las arquitecturas de memorias convencionales, es necesario utilizar directamente 4 MB. Actualmente la **MDRAM** se utiliza en algunos adaptadores de vídeo y aceleradoras gráficas.
- **VRAM** (Vídeo **RAM**) (**RAM** de Vídeo). Memoria de las tarjetas de vídeo; que debe trabajar muy rápido para mantener la pantalla nítida (60-70 veces por segundo). Al contrario que la **RAM** convencional, la **VRAM** resulta accesible para dos dispositivos a la vez. Esto permite a un monitor acceder a la **VRAM** para refrescos de pantalla al mismo tiempo que un procesador gráfico proporciona nuevos datos. La **VRAM** tiene un mejor rendimiento gráfico pero es más cara que la **RAM** convencional. Algunas aceleradoras utilizan **DRAM** convencional, pero otras emplean un tipo especial de **RAM** de vídeo (**VRAM**), que permite acceder a la memoria simultáneamente tanto a la circuitería de vídeo como al procesador.
- **DDR-SDRAM** (Doble Data Rate SDRAM). Los módulos de memoria DDR-SDRAM han venido sustituyendo a los módulos SDRAM tradicionales gracias a que permiten el doble de tasa de transferencia de datos.

¿Cómo es físicamente un módulo de memoria DDR-SDRAM? Los módulos de memoria DDR-SDRAM (o DDR, como los llamaremos en adelante) son del mismo tamaño que los DIMM de SDRAM, pero **con más conectores: 184 pines** (los módulos DDR2 y DDR3 tienen 240 pines) en lugar de los 168 pines de los módulos DIMM-SDRAM tradicionales, este aumento de pines es necesario para implementar el sistema DDR. Además, para que no exista confusión posible a la hora de instalarlos (lo cual tendría consecuencias sumamente desagradables), los módulos DDR tienen **1 única muesca** en lugar de las 2 muescas de los módulos DIMM ([Figuras 3.19](#) y [3.20](#)). En relación al consumo, los módulos DDR emplean niveles de tensión de **2,5 V**

(bajando hasta 1,5 V en el caso de los módulos DDR3), lo que supone una reducción del 25% respecto a los 3,3 V de los módulos DIMM. Esto permite el aumento de autonomía de ordenadores portátiles.

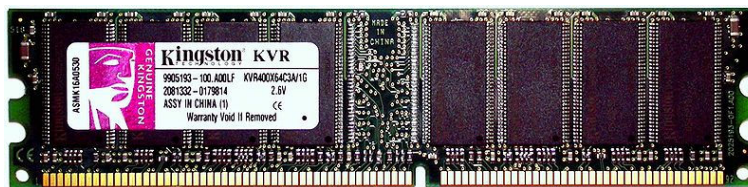


Figura 3.19. Módulo de memoria DDR.



Figura 3.20. Módulos de memoria DIMM.

La [Figura 3.21](#) muestra una comparación gráfica entre memorias de tipo DDR, DDR2 y DDR3.

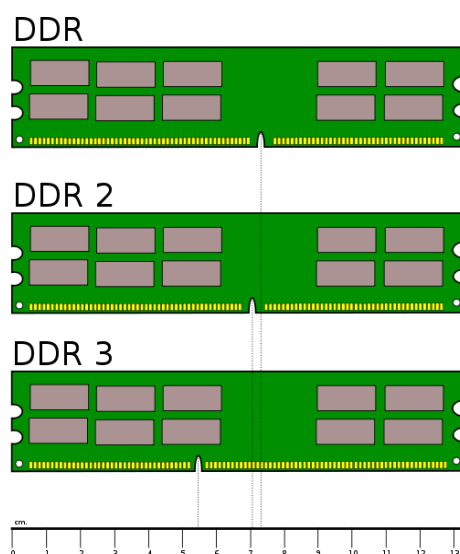


Figura 3.21. Comparación gráfica de módulos DDR.

¿Cómo funciona un módulo de memoria DDR-SDRAM? Los módulos de memoria DDR-SDRAM envían datos 2 veces por cada señal de reloj, una vez en cada flanco de señal (flanco ascendente y descendente), en lugar de enviar datos sólo en uno de los flancos de la señal. De esta forma, un dispositivo con tecnología DDR que funcione con una señal de reloj “real”, “física”, de por ejemplo 100 MHz, enviará tantos datos como otro sin tecnología DDR que

funcione a 200 MHz. Por ello, las velocidades de reloj de los dispositivos DDR se suelen dar en lo que podríamos llamar **“MHz efectivos o equivalentes”** (en nuestro ejemplo, 200 MHz, “100 MHz x 2”). ¿Y por qué se hace esto? ¿No es más fácil subir el número de MHz? intelectualmente es más sencillo, pero **sucede que cuanto más rápido vaya un dispositivo (en MHz “físicos”), más difícil es de fabricar**. Precisamente éste es uno de los problemas de la memoria Rambus: funciona a 266 MHz “físicos” o más, y resulta muy difícil y cara de fabricar.

La tecnología DDR comenzó a emplearse por primera vez en los microprocesadores Intel Pentium 4 y AMD Athlon.

Tipos de DDR-SDRAM y nomenclatura. Existen módulos de memoria DDR de diferentes clases, categorías y precios. Las velocidades de reloj pueden variar entre 100 MHz (DDR-200) y 266 MHz (DDR-533). En el primer caso, el módulo es capaz de transmitir 1600 MB/s, conocido como PC1600, y el segundo 4264 MB/s, conocido como PC4300. En la [Tabla 3.2](#) se recoge de forma cronológica la evolución de los módulos de memoria más comunes empleados para la memoria principal

Año	Tecnología	Velocidad
1987	FPM	20 MHz
1995	EDO	20 MHz
1997	PC66 SDRAM	66 MHz
1998	PC100 SDRAM	100 MHz
1999	RDAM	800 MHz
1999/2000	PC133 SRAM	133 MHz
2000	DDR SDRAM	266 MHz
2001	DDR SDRAM	333 MHz
2002	DDR SDRAM	434 MHz
2003	DDR SDRAM	500 MHz
2004	DDR2 SDRAM	533 MHz
2005	DDR2 SDRAM	800 MHz
2006	DDR2 SDRAM	667-800 MHz
2007	DDR3 SDRAM	1066-1333 MHz

Tabla 3.2. Comparación de módulos de memoria más comunes (Fuente: <http://www.kingston.com>).

En la [Figura 3.22](#) se muestra cómo la capacidad de transferencia de información ha ido aumentando con el desarrollo de la tecnología.

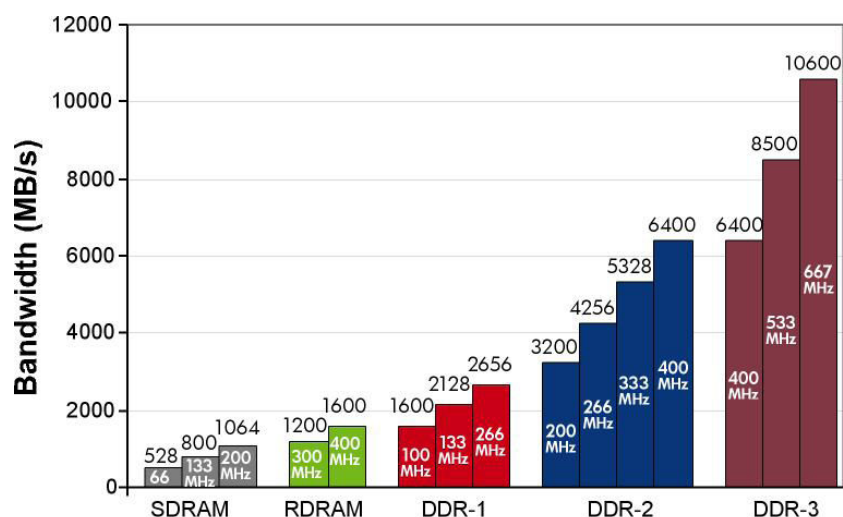


Figura 3.22. Comparación de ancho de banda de módulos de memoria desde SDRAM hasta DDR
(Fuente: www.hp.com).

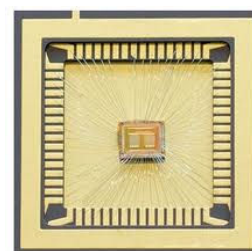
Precios de módulos de memoria. En la [Tabla 3.3](#) se realiza una comparación de precios de distintos módulos de memoria (desde el módulo FPM RAM hasta el módulo DDR3 SDRAM). Como es previsible, con el avance de la tecnología, se han conseguido módulos de mayor capacidad a igual o menor coste que los primeros módulos de memoria. Estos precios son orientativos y pueden variar en función de factores como la demanda, la aparición de nuevos tipos de módulos, etc.

Tipo de Memoria	Capacidad	Fabricante	Precio (€)
FPM RAM	128 MB	HP	42,78
EDO RAM	128 MB	Kingston	23,92
SDRAM	128 MB	HP	28,27
RDAM	512 MB	Kingston	18,12
DDR SDRAM	2 GB	Kingston	28,27
DDR2 SDRAM	2 GB	Kingston	23,20
DDR3 SDRAM	4 GB	Kingston	22,28

Tabla 3.3. Comparación de precios de módulos de memoria (Fuente: <http://www.shopper.cnet.com>, Noviembre 2011).

ÚLTIMOS AVANCES EN MATERIA DE ALMACENAMIENTO DE INFORMACIÓN

Memoria PCM (Fuente: www.ibm.com, Junio 2011): IBM ha desarrollado un nuevo tipo de chip de memoria flash que puede almacenar varios bits por célula. Tradicionalmente, los datos se almacenan como una carga eléctrica. Sin embargo, esta nueva tecnología de chips de memoria, llamada memoria de cambio de



fase (Phase-Change Memory) almacena los datos al cambiar el estado de un material de amorfo a cristalino. Los chips PCM almacenan los datos cuando hay un cambio de resistencia en el material. Requiere menos consumo y permite almacenar datos hasta 100 veces más rápido que una memoria actual. El chip presenta una anchura de 90 nm, y es capaz de almacenar datos aun cuando no está alimentado eléctricamente (memoria no volátil).

PCRAM, Phase-Change RAM (Fuente: Y. Xie. Modeling, Architecture, and Applications for Emerging Memory Technologies. *IEEE Design & Test of Computers*, 2011). Esta tecnología se basa en una aleación de $\text{Ge}_2\text{Sb}_2\text{Te}_5$ (GST). La capacidad de almacenamiento se realiza a partir de la diferencia de resistencia entre la fase amorfa (alta resistencia) y la fase cristalina (baja resistencia) de la aleación del material. Para las operaciones de lectura/escritura, el material se cristaliza mediante la aplicación de un pulso eléctrico que calienta una porción significativa de la célula por encima de su temperatura de cristalización. Para el borrado de la celda, se aplica de forma abrupta una corriente eléctrica tan elevada que es capaz de fundir el material y, seguidamente se corta, dejando el material en estado amorfo. Este tipo de memorias es muy similar a la memoria PCM comentada anteriormente.

STT-RAM, Spin Torque Transfer RAM (Fuente: Y. Xie. Modeling, Architecture, and Applications for Emerging Memory Technologies. *IEEE Design & Test of Computers*, 2011). Consiste en un nuevo tipo de memorias RAM magnéticas (MRAM) no volátiles, alta velocidad de lectura/escritura (<10 ns), capaz de soportar más de 10¹⁵ ciclos de borrado/escritura (programación), y consumo nulo en estado de standby. La capacidad de almacenamiento o programación de este tipo de memorias MRAM deriva de la unión magnética a través de un túnel (Magnetic Tunneling Junction), en el que una delgada capa en forma de túnel de un dieléctrico, por ejemplo, óxido de magnesio (MgO) es intercalada entre dos capas ferromagnéticas. Una unión MTJ tiene baja resistencia si las dos capas ferromagnéticas están magnetizadas con misma polaridad; de lo contrario la unión MJT presentará alta resistencia.

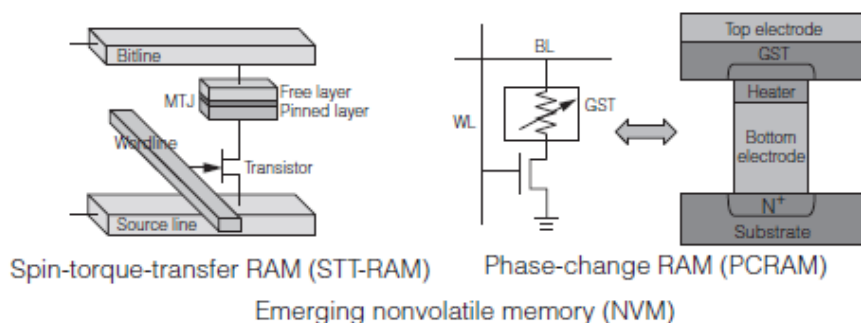


Figura 3.23. Tecnologías emergentes STT-RAM y PCRAM (Fuente: www.ieee.org).

RRAM, Resistive RAM (Fuente: Global Semiconductor Alliance, <http://www.gsaglobal.org>). En una celda RRAM, los datos se almacenan en forma de dos (single-level cell, SLC) o más estados (o niveles) de resistencia (multi-level cell, o MLC) del dispositivo de conmutación resistiva (Resistance Switch Device, RSD). El cambio de resistencia en los óxidos de metales fue descubierto hace décadas en una capa fina de NiO. A partir de entonces, se ha desarrollado una gran variedad de materiales de óxido de metal con características de conmutación resistiva. En base a los mecanismos de almacenamiento, los materiales para construir las memorias de tipo RRAM pueden ser catalogados como basado en filamento, basado en interfaz, o células de metalización programable (Programmable Metallization Cell, PCM). En base a la propiedad eléctrica de conmutación resistiva, los dispositivos RSD se pueden dividir en dos categorías: unipolar o bipolar. Las memorias RRAM basadas en filamento son un ejemplo típico de conmutación unipolar. El material aislante entre dos electrodos se puede hacer conductor aplicando un nivel de tensión suficientemente alto. El almacenamiento de datos se puede lograr mediante la ruptura (RESET) o vuelta a conectar (SET) del camino conductor. Por otro lado, las memorias basadas en PMC representan la tecnología de conmutación bipolar. Su mecanismo de conmutación puede ser explicado como la formación o ruptura del “nanocable” moviendo los iones de metal entre dos electrodos de metal.

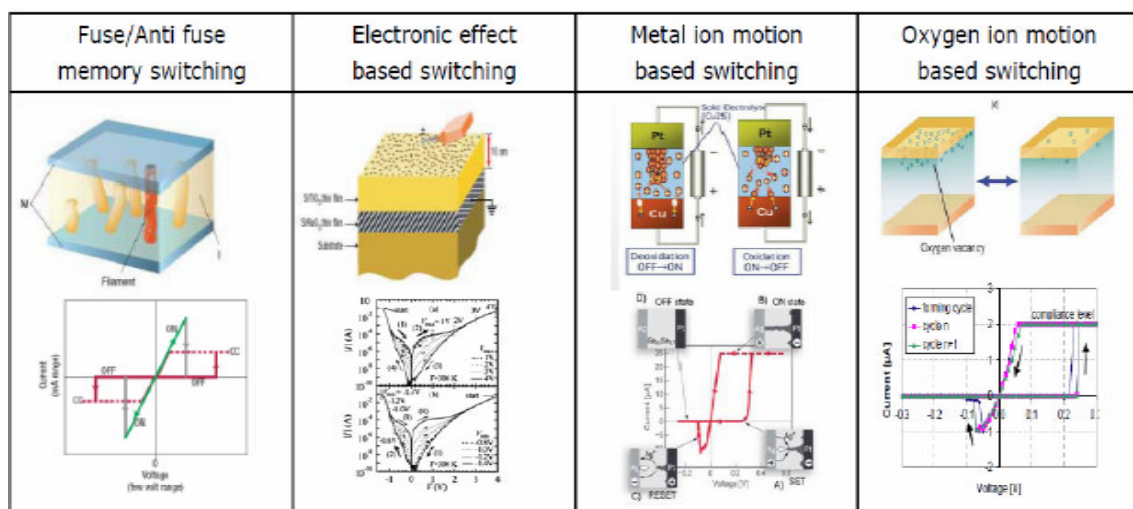


Figura 3.24. Técnicas de conmutación resistiva para memorias RRAM (Fuente: <http://www.gsaglobal.org/events/2010/0316/docs/7.GMC-PierreFazan.pdf>).

MEMRISTOR (Fuente: Artículo “Design Methodologies and Circuit Techniques for Emerging Non-Volatile Memories”). Esta tecnología consiste en hacer que el comportamiento resistivo del material sea función del tiempo. Es decir, la Memresistencia es una función de la carga, que depende del comportamiento histórico de la corriente (o tensión).

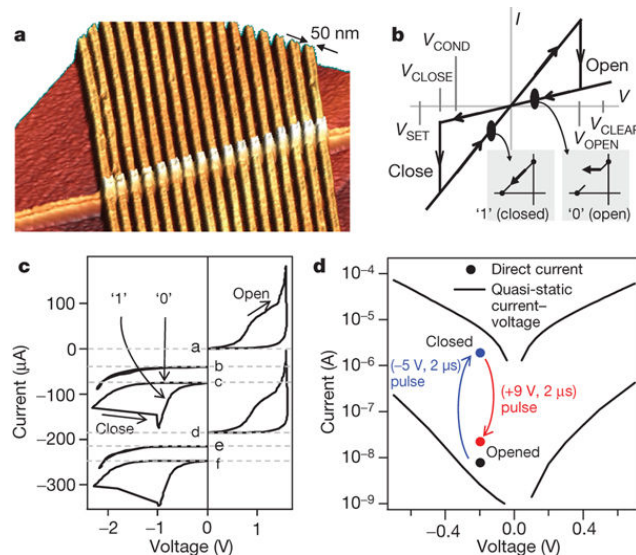


Figura 3.25. Caracterización de una celda Memristor (Fuente: J Borghetti et al. *Nature* 464, 873-876 (2010) doi:10.1038/nature08940).

Básicamente se trata de un componente pasivo de dos terminales eléctricos entre los que existe una relación entre la carga eléctrica y el flujo magnético. Cuando la corriente fluye en una dirección a través del material, aumenta la resistencia eléctrica. Y cuando la corriente fluye en la dirección opuesta, la resistencia disminuye. Cuando la corriente se detiene, el material conserva la última resistencia que había, y cuando el flujo de carga se inicia de nuevo, la resistencia del circuito será la que era cuando estuvo activo por última vez.

DISPOSITIVOS DE ALMACENAMIENTO PARA MEMORIA SECUNDARIA

En esta sección se verán dispositivos empleados para la implementación de la memoria secundaria en el sistema computador. Actualmente, la clasificación entre el nivel de memoria secundaria y el de memoria auxiliar no es tan clara como puede serlo entre el nivel de memoria principal y el de memoria secundaria, por lo que puede considerarse un que el nivel de memoria secundaria y auxiliar forman un mismo nivel. Los dispositivos que se describirán a continuación son: Cinta magnética, Disco magnético y Memorias ópticas.

CINTA MAGNÉTICA. La primera unidad de cinta magnética fue introducida en 1951, modelo UNISERVO para sistema computador UNIVAC 1, con una capacidad de 1,44 Mb y velocidad de 2,45 m/s. Las unidades de cinta magnética convencionales emplean como soporte una cinta Mylar sobre la que se deposita una capa de óxido de hierro, óxido de cromo o partículas de metal. Su longitud suele ser de unos 800 m. La cinta se divide en 9 pistas, cada una asignada a su correspondiente cabeza de lectura-escritura. Se leen 9 bits en paralelo: 8 de datos y 1 de paridad impar. No toda la longitud de la cinta puede

emplearse para almacenar datos. Por un lado se debe almacenar información relativa a la dirección para poder alcanzar la zona de datos deseada, y por otro se deben dejar claros o zonas muertas entre los registros que se quieren leer de forma independiente, es decir, parando la cinta entre ellos. Estos claros reciben el nombre de IRG, siglas que corresponden a *Inter Record Gap*.

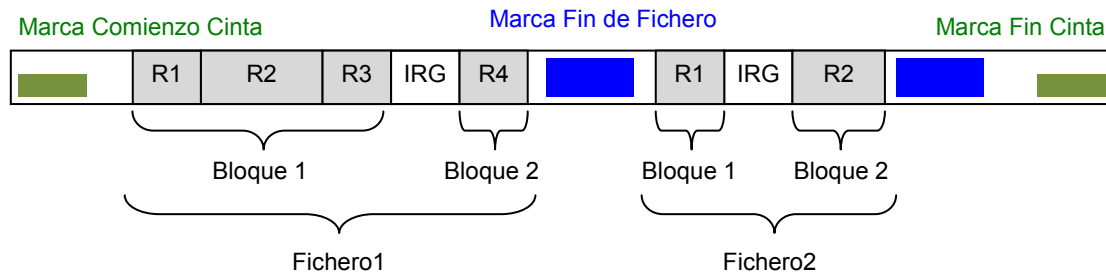


Figura 3.26. Formato de grabación de cinta magnética.

La información está organizada en ficheros, y cada fichero en varios registros que pueden formar un bloque. Cada fichero tiene un registro de cabecera con su nombre y características que permite identificarlo. A su vez cada bloque tiene 3 zonas diferenciadas: Cabecera, con información relativa a la sincronización e identificación, Datos y Cola, con código de detección de error y fin de registro. La unidad sólo permite leer y grabar información en un sentido, pero puede retroceder un número de registros determinado.

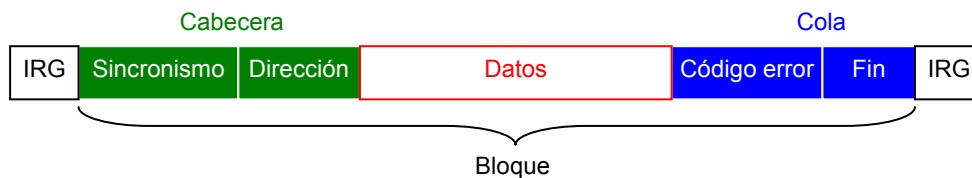


Figura 3.27. Organización de bloques en una cinta magnética.

Las características principales de estos dispositivos son:

- Memoria dinámica, la cinta ha de moverse delante de la cabeza para lectura/escritura.
- Existe contacto físico entre la cinta y las cabezas de lectura/escritura.
- Pueden leer o grabar datos en paralelo.
- Soporte económico.
- Grabación horizontal.
- Acceso secuencial, hay que leer toda la cinta hasta llegar al lugar deseado.
- No puede intercalar información adicional, por lo contrario hay que regrabar todo el resto de la cinta.

- Se emplean fundamentalmente para almacenamiento de archivo muerto o copia de seguridad.
- Memoria no volátil.
- Memoria de lectura no destructiva.

Como último de los avances en cintas magnéticas, en 2010 las compañías Fujifilm e IBM anunciaron un nuevo tipo de cinta magnética empleando nanotecnología basada en partículas de BaFe con capacidad de almacenamiento de 4.57 billones de bits/cm², lo que permitiría conseguir una capacidad de almacenamiento de 35 TB (Fuente: http://www.fujifilmusa.com/press/news/display_news?newsID=879807)

DISCOS MAGNÉTICOS. En 1956 aparece el IBM 350, primer disco con brazo móvil y cabeza flotante. Este disco tenía una capacidad de 5 MB y un tiempo de acceso de aproximadamente 0.5 s. A partir de aquí, el disco duro ha venido tomando cada vez más relevancia, hasta convertirse hoy en día, en uno de los periféricos claves en todo sistema computador. Estos dispositivos son de acceso directo, frente al acceso secuencial de las cintas magnéticas comentadas anteriormente.

El soporte ([Figura 3.28](#)) consiste en un disco recubierto por una película magnética, que gira a gran velocidad (desde 3600 rpm hasta 10000 rpm). El sistema puede tener una o varias cabezas de lectura/escritura por superficie. En el primer caso se hablará de discos de cabeza fija y en el segundo de discos de cabeza móvil. El sistema puede tener una o varias superficies o platos que se dividen en: Pistas, Sectores y Cilindros.

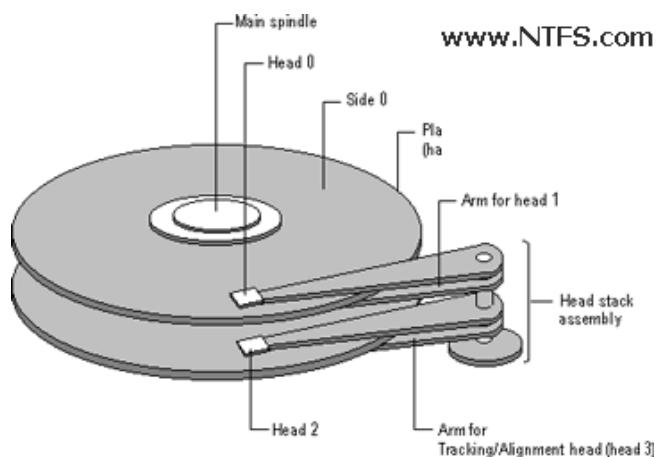


Figura 3.28. Detalle de partes de un disco duro (Fuente: <http://www.ntfs.com/hard-disk-basics.htm>).

Una **pista** es la tira del soporte de almacenamiento que gira delante de la cabeza. En los dispositivos de cabeza fija, cada pista tiene asociada una cabeza, mientras que en los discos de cabeza móvil cada cabeza en una posición determinada define una pista. Cada pista se divide en **sectores**, siendo el sector la unidad de información que se transfiere en un acceso. En los sistemas de varias superficies, aquellas pistas alineadas verticalmente a las que se accede en cada posición del brazo constituyen un **cilindro**.

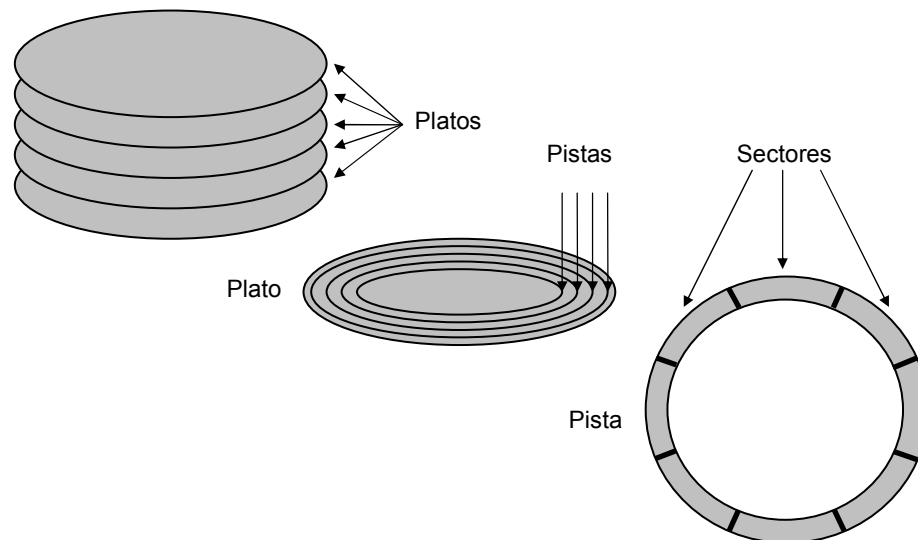


Figura 3.29. Platos, pistas y sectores en un disco duro.

Inicialmente, todas las pistas tenían el mismo número de sectores, normalmente 32; posteriormente, para optimizar la capacidad de almacenamiento se construyen discos con mayor número de sectores en las pistas externas que en las internas. Los discos de los grandes computadores de IBM permiten que el usuario seleccione el tamaño de los sectores. Aunque, casi todos los demás sistemas fijan el tamaño que corresponde a un sector (normalmente 512 bytes por sector).

El direccionamiento de las unidades de disco se descompone en las siguientes fases:

- Selección de la cabeza correspondiente (direccionamiento cableado).
- Posicionamiento del brazo en caso de ser móvil (direccionamiento mecánico).
- Interpretación de la información leída de la pista, para seleccionar el sector deseado (direccionamiento almacenado).

Debido a la inercia mecánica del disco, éste tarda un tiempo apreciable en alcanzar el régimen estable de su velocidad de rotación. Por ello, el disco está siempre girando, por lo que se ha de evitar que las cabezas de lectura/escritura toquen el soporte para evitar su desgaste. Las cabezas se montan sobre un “deslizador” que hace que éstas “vuelen” a una “altura” del orden de una fracción de mm de la superficie del disco. Esto se consigue gracias al aire que desplaza el disco al girar. La corriente de aire hace que la cabeza planee sobre la superficie adaptándose a las irregularidades que inevitablemente tiene el soporte.

Otro problema a considerar en las unidades de disco son los cortes de alimentación. La disminución de la velocidad del disco, disminuye la fuerza ascensional de la cabeza, por lo que ésta choca con la superficie o plato. Adicionalmente, decir que las partículas

de polvo suponen otro problema importante; tener en cuenta que una partícula de polvo es 3000 veces superior a la distancia a la que se deslizan las cabezas sobre la superficie.

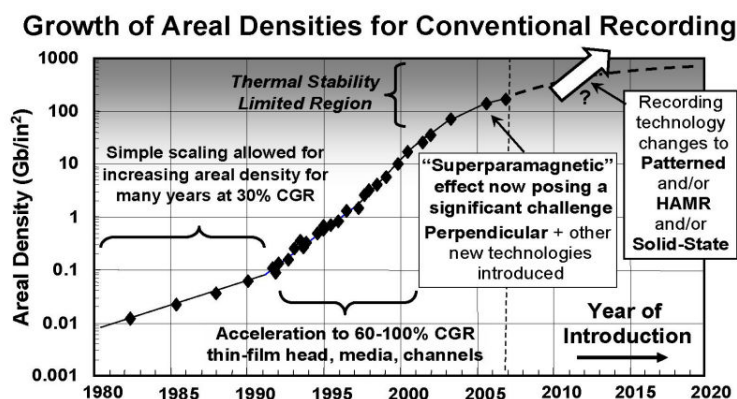


Figura 3.29. Previsión de la evolución de la densidad de grabación en disco duro (Fuente: *A Perspective on the Future of Hard Disk Drive (HDD) Technology*, P. Frank, R. Wood, IEEE Asia-Pacific Magnetic Recording Conference, 2006).

La [Figura 3.29](#) muestra una previsión de la evolución de la densidad de área de grabación y en la [Tabla 3.4](#) se recogen algunas características de discos duros recientes.

Factor de forma	Anchura (mm)	Altura (mm)	Capacidad (año)	Platos
3.5"	102	19 ó 25.4	4 TB (2011)	5
2.5"	69.9	7, 9.5, 11.5 ó 15	1.5 TB (2010)	4
1.8"	54	5 u 8	320 GB (2009)	3

Tabla 3.4. Características de discos duros recientes.

A continuación se comentarán algunos aspectos concretos relacionados con la tecnología de los discos duros como son: Tipos de cabeza, Formato de grabación, Grupos de discos, Tiempo de acceso, Velocidad de transferencia, Caché de disco, Operaciones de varios sectores.

- **Tipos de Cabeza.** La tecnología de fabricación de las cabezas de lectura/escritura ha ido evolucionando para conseguir mayores densidades de grabación, pudiendo diferenciar los siguientes tipos:
 - Cabeza de ferrita. Formada por un toro de ferrita, un entrehierro de cristal y un arrollamiento por el que circula la corriente de grabación o de lectura.
 - Cabeza MIG (Metal-in-Gap). Consiste en depositar una aleación metálica como el AlFeSi en los bordes interiores del entrehierro. Con esto se consigue aumentar el campo magnético generado por la cabeza.
 - Cabeza de película delgada. Se fabrican depositando, en pasos sucesivos y mediante técnicas de vacío o fotolitográficas, los materiales que constituyen el transductor. Su ventaja es que se consiguen transductores de muy reducido tamaño.

- Cabeza magnetorresistiva. Se trata de una cabeza de película delgada con dos transductores, uno para la lectura y otro para la escritura.
- **Formato de grabación.** Para poder reconocer la información del disco hay que añadir una información de direccionamiento y, a veces, de sincronismo. El formato de grabación especifica esta información, además de incluir unos claros o *gaps* como sucede en las cintas magnéticas. El formato puede ser de tipo *hardware* (en desuso) o *software*. El primer tipo consistía en hacer unas muescas en el borde de uno de los discos para indicar el comienzo de cada sector. En el segundo caso, los sectores quedan delimitados por el código de cabecera.
- **Grupos de discos.** Cuando se desea obtener una alta fiabilidad los discos se organizan en grupos y se añade algún tipo de representación redundante. Mediante un código corrector de error, se puede asegurar que la información no se perderá aunque falle algún disco. Estos grupos de discos se conocen como RAID (*Redundant array of Inexpensive Disks*). A continuación se muestran ejemplos de RAID 1, RAID 3 y RAID 5.
- **RAID 1.** Crea una copia exacta (o espejo) de un conjunto de datos en dos o más discos. Esto resulta útil cuando el rendimiento en lectura es más importante que la capacidad. Un conjunto RAID 1 puede ser tan grande como el más pequeño de sus discos. En este caso se incrementa exponencialmente la fiabilidad respecto a un solo disco; es decir, la probabilidad de fallo del conjunto es igual al producto de las probabilidades de fallo de cada uno de los discos (pues para que el conjunto falle es necesario que lo hagan todos sus discos).

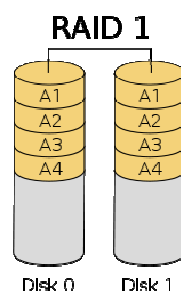


Figura 3.30. Diagrama de una configuración RAID 1.

- **RAID 3.** En este esquema se utilizan n discos más uno de redundancia. Un bloque de información se reparte en $n-1$ trozos iguales, que se graban en los discos de datos. El código corrector se graba en el último disco. Dado que lo mínimo que se puede grabar en una operación de acceso a

un disco es un sector, este esquema obliga a grabar bloques muy grandes, de $n-1$ sectores. Es un esquema interesante cuando los ficheros a almacenar son grandes. Nótese que la detección de la unidad que ha fallado se hace mediante el código detector de error incluido en el sector.

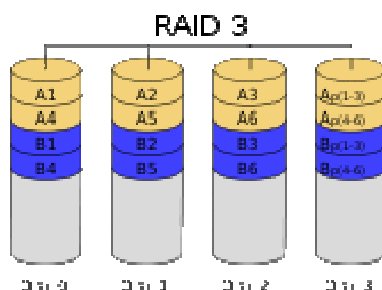


Figura 3.31. Diagrama de una configuración RAID 3. Cada número representa un byte de datos y cada columna un disco.

- **RAID 5.** En este esquema todos los discos tienen datos y redundancia de los otros discos. Funciona bien con bloques de tamaño reducido, por lo que es de uso general. Sin embargo penaliza las escrituras pequeñas, puesto que se ha de leerla redundancia, calcular su nuevo valor y volver a escribirla.

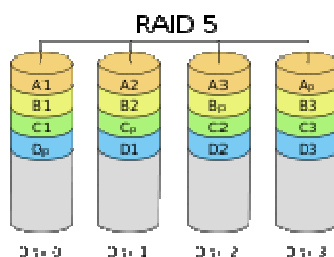


Figura 3.32. Diagrama de una configuración RAID 5.

- **Tiempo de acceso.** El tiempo de acceso de estos dispositivos viene dado por el tiempo que tardan en posicionar el brazo en la pista deseada que se llama **tiempo de búsqueda**, más que tarda la información de la pista en pasar delante de la cabeza que se llama **latencia**. En caso de tener el disco apagado hay que añadir el tiempo de arranque del mismo. El tiempo de búsqueda se divide en las siguientes fases:
 - **Aceleración.** El brazo es acelerado hasta que o bien alcanza la mitad de la distancia de la búsqueda o bien alcanza la máxima velocidad de movimiento. Esta fase puede durar entre 5-10 ms.
 - **Deslizamiento.** En búsquedas largas el brazo se mueve a la velocidad máxima.

- Frenado. El brazo es frenado hasta acercarse al cilindro deseado.
- Ajuste. El controlador del disco ajusta la posición al cilindro deseado. Esta fase dura del orden de 1-3 ms.
- **Velocidad de transferencia.** Consiste en la velocidad a la que se transmiten los bits de un sector, que depende de la velocidad de giro y de la densidad de grabación angular. Así por ejemplo, suponiendo una velocidad de grabación lineal de $d = 10000$ bits/cm, un radio de disco de $R = 5$ cm y una velocidad de giro de 3600 rpm, el número total de bits (datos más direccionamiento) de la pista será de $2\pi R d \approx 315000$ bits. Por tanto la velocidad de transferencia será de $315000 \text{ bits} \cdot 3600 \text{ bits/min} = 18.9 \text{ Mbits/s}$. Empleando bytes de 8 bits más bit de paridad, se obtiene una velocidad de transferencia de 2.1 MBytes/s. Observar que la velocidad de transferencia debe calcularse partiendo de la capacidad de almacenamiento bruta y no neta.
- **Caché de disco.** Muchos controladores de disco incluyen una memoria de semiconductor para acelerar su respuesta. Esta memoria funciona a modo de caché de disco y contiene la información más reciente accedida.
- **Operaciones de varios sectores.** En el acceso a un disco es muy frecuente realizar operaciones que afecten a una serie de sectores consecutivos, puesto que con un solo acceso se transfiere más información. En estas operaciones adquieren importancia los siguientes factores: Tiempo de conmutación entre cabezas y Tiempo de acceso a la primera pista del cilindro siguiente. Para optimizar los accesos consecutivos, los sectores de las pistas contiguas no se organizan alineados al mismo radio, sino que se intercalan (*interleave*) adecuadamente de forma que se pueda leer sin perder tiempo entre el último sector de una pista y el primero de la siguiente, de lo contrario se tendría que esperar un giro completo del disco.

Las características principales de estos dispositivos son:

- Memoria dinámica de acceso directo.
- Memoria no volátil.
- Memoria de lectura no destructiva.
- El proceso de grabado se hace en serie, las pistas empeladas son de 1 bit de ancho.
- Tiempo de acceso que puede llegar a ser de unos pocos milisegundos.
- Velocidad de transferencia de varios Mbytes/s.
- Capacidad de almacenamiento de cientos de Gigabytes por unidad.
- Existen modelos con soporte fijo e intercambiable.

MEMORIAS ÓPTICAS. Los discos ópticos almacenan la información binaria en una capa interna protegida, pudiéndose emplear distintas tecnologías para ello. En todos ellos la operación de lectura se realiza haciendo incidir un rayo láser que, al ser reflejado por la superficie, permite detectar variaciones microscópicas de propiedades reflexivas ópticas ocasionadas como consecuencia del proceso de grabación (escritura) en el mismo. Las tecnologías de grabación (escritura) que se emplean son:

- Por moldeado durante la fabricación, mediante un molde de níquel (CD-ROM y DVD-ROM).
- Por la acción de un haz láser (CD-R y CD-RW).
- Por la acción de un haz láser en conjunción con un campo magnético (discos magneto-ópticos, MO).

Comparando las características de los discos ópticos y de los discos magnéticos, se tiene:

- Los discos ópticos, además de ser medios removibles con capacidad para almacenamiento de grandes cantidades de datos, son portables y seguros en la conservación de los datos y son no volátiles.
- El coste por byte almacenado es pequeño.
- Los discos ópticos tienen mayor capacidad que los discos magnéticos puesto que el haz láser incide en la superficie del disco puntualmente y con la precisión característica del enfoque óptico del láser, permitiendo que tanto los bits en una pista como las pistas estén más próximos.
- Los discos ópticos son más seguros que los discos magnéticos, puesto que la capa que almacena los datos es inmune a los campos magnéticos externos que puedan rodearlo; estando incluso protegida de la corrosión ambiental y de los posibles desperfectos que se les puede ocasionar con la manipulación de los mismos. Esto es consecuencia de encontrarse esta capa entre dos capas transparentes de policarbonato.
- La cabeza móvil que porta la fuente láser y la óptica asociada nunca puede tocar la superficie del disco por encontrarse a 1 mm de distancia; no produciéndose desgaste por rozamiento, ni existe riesgo de aterrizaje como ocurre en los discos rígidos de cabeza flotante. Tampoco el láser que incide sobre la superficie que contiene la información puede afectarla de modo alguno dada su escasa intensidad.

Las aplicaciones de los discos ópticos son:

- Las bases de datos en CD-ROM para bibliotecas de datos.
- Para distribución de software.
- Para manuales de software.

- Para demostración de productos y aplicaciones.
- Para servidores de archivos en una red local.
- Para copias de resguardo seguras.
- Para bibliotecas de imágenes.

La permanencia de la información en un CD-ROM común se estima entre 10 y 15 años, puesto que la superficie de aluminio que contiene la información se va oxidando muy lentamente en ese tiempo; salvo que la superficie se someta a una protección antióxido especial (de oro). Por esa razón, la vida media de la información en un CD-R es bastante mayor (la fina capa metálica interior presenta oro).

En informática, los tipos de discos ópticos que se emplean son:

- Los grabados masivamente por el fabricante, para ser leídos por el usuario. En estos, a partir de un disco máster grabado con luz láser, se realizan múltiples copias obtenidas por inyección de material (sin usar láser). La fina capa de aluminio, entre dos capas transparentes protectoras, contiene en una cara los unos y ceros como surcos discontinuos, formando una sola pista en espiral; la espiral es leída con luz láser por la unidad de CD del usuario. Tipos:
 - El CD-ROM (Compact Disc ROM - Disco Compacto de sólo lectura).
 - El DVD-ROM (Digital Versatil Disc ROM - Disco Versátil Digital de sólo lectura).
 - El BD (Blu Ray Disc – Disco de Rayo azul. A Blu le falta la e de Blue por problemas de patente de nombre de colores en algunos países). Empleado para vídeo de alta definición y para almacenamiento de datos de alta densidad. El uso del láser azul para escritura y lectura permite almacenar más cantidad de información por área que los discos DVD, debido a que el láser azul tiene una menor longitud de onda que los láseres usados para almacenar en discos DVD. Su capacidad de almacenamiento llega a 50 Gigabytes a doble capa, y a 25 Gigabytes a una capa. El Blu-ray de 400 GB a 16 capas ya fue patentado y salió al mercado en 2010, así como un Blu-Ray de 1 Terabyte en 2011. El Blu Ray es un dispositivo de almacenamiento más eficiente que los DVD corrientes, en cuanto a capacidad de almacenamiento, en cuanto a fiabilidad de lectura y en cuanto a sensibilidad a posibles ralladuras en el disco.
- Los grabables una sola vez por el usuario. En la escritura el haz láser sigue una pista en espiral preconstruida en una capa de pigmento; donde incide el haz, su calor decolora para siempre el punto de incidencia. En la lectura, esta capa deja pasar el haz láser hacia la capa reflectora dorada que está debajo,

reflejándose de forma distinta, según el haz haya atravesado un punto decolorado o no, detectándose así los unos y los ceros. Ambas capas a su vez están protegidas por dos capas transparentes. Tipos:

- El CD-R (Compact Disc Recordable – Disco Compacto Grabable). Se puede grabar en varias sesiones, pero la información grabada no puede ser borrada ni reescrita; las sesiones posteriores utilizan el espacio libre dejado por las sesiones anteriores.
- El DVD-R y el DVD+R. (DVD Recordable – DVD Grabable). Disco óptico en el que se puede grabar o escribir datos con una mayor capacidad de almacenamiento que en un CD-R, normalmente 4.7 GB en lugar de los 700 MB de almacenamiento estándar de los CD-R. Esta mayor capacidad se debe a la mayor densidad de pistas y un haz láser con una longitud de onda más corta. Tanto el DVD-R como el DVD+R hacen referencia al mismo tipo de disco pero según distintos fabricantes. Están compuestos de dos discos de policarbonato de 0.6 mm de espesor, pegados con un adhesivo el uno al otro. En uno de los discos está el surco que guía el láser y está cubierto con el tinte grabador y un reflector y, la otra superficie sirve únicamente para asegurar mecánicamente el conjunto (para darle rigidez), además de conseguir compatibilidad de espesor con los discos CDs.
- Los borrables y regrabables. Existen los siguientes tipos:
 - El CD-RW (CD ReWritable – CD Reescribible). Con tecnología de grabación magneto-óptica. El haz de rayo láser calienta puntos (que equivale a “1s” lógicos de información binaria) de una capa previamente magnetizada uniformemente para que pierdan su magnetismo original (que corresponde a “0s”); al mismo tiempo que el haz de láser calienta la zona, se aplica un campo magnético para producir en esos puntos un campo magnético contrario al original (para grabar unos). Estas diferencias puntuales de magnetización son detectadas en la lectura por la luz láser (de menos potencia que la correspondiente al proceso de escritura), puesto que provocan distinta polarización de la luz que reflejan según estén polarizados.
 - El CD-PD (CD Phase Dual – CD de Dos Fases). La escritura se realiza por cambio de fase (de cristalina a amorfa o viceversa) de los puntos de la capa del disco que contiene los datos. Se trata de una tecnología puramente óptica, sin magnetismo, que requiere una sola pasada para escribir una porción o la pista en espiral completa. En la tecnología PD

(Phase Change/Dual), que también es por cambio de fase, la unidad escribe pistas concéntricas. La palabra Dual indica que responde tanto a unidades de lectura espiral como a unidades de lectura concéntrica.

- El DVD-RW y el DVD+RW (DVD ReWrite – DVD Reescribible) es de tecnología similar a la empleada en el CD-RW pero admite mayores densidades de grabación y, como consecuencia mayores capacidades de almacenamiento. La capacidad estándar es de 4.7 GB. El DVD+RW y el DVD-RW son parecidos pero existen algunas diferencias entre ellos. El surco que corresponde al DVD+RW ondula a mayor frecuencia que el que corresponde al DVD-RW y permite mantener constante la velocidad de rotación del disco. La mayor ventaja del DVD+RW es que permite mayor rapidez de escritura que el DVD-RW, por no tener que realizar un formateo previo (de 2 a 4 minutos) ni realizar la operación del cierre del disco (puede ser de hasta 30 minutos).
- El DVD RAM. Se trata de un disco DVD regrabable aprobado por el DVD-Forum (Es una organización compuesta por las compañías de hardware, software, medios de comunicación y contenidos que usan y desarrollan el formato DVD. Inicialmente se bautizó como el DVD Consortium cuando fue fundado en 1995. Se creó para facilitar el intercambio de información e ideas sobre el formato DVD y permitirle crecer gracias a las mejoras e innovaciones técnicas). Se diferencia del DVD-RW y del DVD+RW en que no hace falta borrar todo el disco para recuperar el espacio que ocupa la información que se quiere borrar y en que se puede escribir directamente en ellos (igual que en un disco duro) sin necesidad de programas de grabación de DVD, ni de programas controladores intermedios. Inicialmente los discos tenían una capacidad de 2.9 GB y estaban encerrados en una carcasa protectora llamada CADDY, poco práctica pero necesaria (los discos DVD RAM son bastante vulnerables a la suciedad y a las ralladuras). Actualmente los DVD RAM que se venden son de 4.7 GB y sin la carcasa protectora, existiendo discos que emplean las dos caras para obtener el doble de capacidad. Los DVD RAM son competidores de las cintas magnéticas para backups, siempre que el coste por byte almacenado lo justifique.