

Subsistema de memoria

Organización de computadoras 2015

Universidad Nacional de Quilmes

1 Subsistema de memoria

El sistema de computos necesita tener varios tipos de memorias, para ser utilizadas en diferentes situaciones, cumpliendo distintos objetivos. Antes de mostrar porque es esto cierto, veamos cómo caracterizar las memorias según los diferentes aspectos.

En primer lugar, una unidad puede ser interna o externa al sistema según si es fija o desmontable (esto es, portable a otro sistema).

Además puede ser volátil, si su contenido se pierde al perder la alimentación eléctrica, o persistente, si no necesita electricidad para mantener la información.

Puede ser de lectura y escritura, como se requiere en una memoria principal, o bien de solo lectura para almacenar información estática que no necesita modificarse.

Por último, se las puede clasificar según su método de acceso, en secuencial, directo, aleatorio o asociativo. El **método de acceso secuencial** requiere que la dirección (o identificación) de cada dato esté almacenada junto con él, y por lo tanto el dispositivo debe recorrerse secuencialmente hasta encontrar la identificación buscada. El **método de acceso directo** es el utilizado por los discos magnéticos, donde la dirección del dato se basa en su ubicación física, en particular la búsqueda se da en dos etapas: se accede primero a la zona que incluye al dato y luego se busca secuencialmente dentro de esa zona. En el **método aleatorio** cada dato tiene un mecanismo de acceso único y cableado físicamente (cada acceso es independiente de los accesos anteriores). Finalmente, el **método asociativo** organiza cada unidad de información con una etiqueta que la describe en función de su contenido. Entonces, para recuperar un dato se debe analizar un determinado conjunto de bits dentro de la celda, buscando la coincidencia con la clave o patrón buscado. Esta comprobación del contenido de las celdas se lleva a cabo de manera simultánea en todas las celdas.

1.1 Memorias ROM

Para algunas aplicaciones, el programa no necesita ser modificado y entonces puede ser almacenado de manera permanente en una memoria de solo lectura ó ROM (*Read Only Memory*). Ejemplos de memorias ROM pueden encontrarse en videojuegos, calculadoras, hornos de microondas, computadoras en automóviles, etc. Además casi todas las computadoras personales necesitan una memoria ROM donde almacenar el primer programa que da arranque al

sistema operativo a partir del acceso (en la mayoría de los casos) a un disco rígido primario.

1.2 Jerarquía de memorias

Como se dijo, la **memoria principal** es volátil y de acceso aleatorio. Si es volátil entonces se necesita otra posibilidad de almacenamiento donde los programas puedan almacenarse de manera persistente y desde donde se recuperen bajo demanda del usuario (cuando dispara la ejecución de un programa).

Esta **memoria secundaria** puede ser resuelta con un disco rígido o un disco de estado sólido (SSD) donde se mantengan persistentes (instalados) los programas. Cuando se necesita de la ejecución de un programa, hecho que puede darse a partir del requerimiento de un usuario o por invocación de otro programa, este debe ser cargado en memoria y permanece allí hasta que la memoria se sobrescribe o se apaga el sistema.

La pregunta que puede surgir es ¿porque no montar la memoria principal en una tecnología persistente, como puede ser un disco de estado sólido? Para responderla es importante destacar que existe una relación de compromiso entre el tiempo de acceso, el costo por bit y la capacidad de almacenamiento (ver figura 1). Un disco tiene mayor capacidad (y por lo tanto menor costo por bit) pero un tiempo de acceso mucho mayor. Esto impacta directamente en el desempeño de la CPU pues el tiempo de ejecución de una instrucción está condicionado por el tiempo de acceso a memoria principal. Además, en este punto es importante notar que algunos programas pueden no ser ejecutados nunca y algunos datos pueden no ser accedidos, aunque se los tenga disponibles en el sistema.



Figure 1: Relacion de compromiso

Es por esto que se incorpora una memoria mas pequeña y rápida (de acceso aleatorio) donde se cargan los programas a la hora de ser ejecutados y los datos cuando se los

necesita. Si lo que se necesita es asegurar una velocidad aceptable y no se necesita gran capacidad, ¿Porqué no implementar la memoria principal con registros de la CPU? Para responder esto se debe tener en mente que el costo sea razonable. En general las arquitecturas cuentan con pocas decenas de registros debido al costo que eso implica. Por otro lado, si se implementa la memoria principal con una RAM se otorga flexibilidad para extender su tamaño pues se trata de un componente externo a la CPU y en cambio los registros suelen estar definidos en el diseño electrónico de la misma. Este esquema se describe en la figura 2.

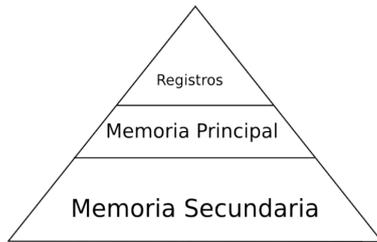


Figure 2: Jerarquía de memoria

2 Memoria Caché

Desde hace un par de décadas, el progreso tecnológico de las computadoras personales esta marcando una tendencia a duplicar, cada año y medio, la velocidad de los procesadores y el tamaño de la memoria principal pero la velocidad de las memorias apenas crece un 10%. Esto ocasiona un crecimiento de la brecha entre la velocidad del procesador y la velocidad de la memoria, causando un mayor impacto en el tiempo de ejecución de los programas, pues la velocidad con la que la CPU puede ejecutar instrucciones está limitada por el tiempo de acceso a memoria, dado que al menos una vez es necesario accederla dentro del ciclo de ejecución de instrucción.

Un enfoque para buscar una solución a esta brecha es incorporar una memoria mas pequeña que la memoria principal pero mas rápida, teniendo en mente que no es viable construir la memoria principal a partir de registros internos a la CPU. Esta idea se basa en que los accesos que se solicitan no son al azar, sino que respetan una lógica que tiene relación con la naturaleza de la ejecución de los programas, pudiendose distinguir dos tipos de accesos: lectura de instrucciones y lectura de datos. La lectura de instrucciones es, en su mayoría, un recorrido de celdas consecutivas de memoria, y así mismo la lectura de los datos de un arreglo también lo es. Por otro lado, las instrucciones de una rutina o de una repetición se utilizan mas de una vez. Este patrón de acceso se describe con los **principios de localidad**: el principio de localidad espacial enuncia que las posiciones de memoria cercanas a alguna accedida recientemente son mas probables de ser requeridas que las mas distantes, y el principio de localidad temporal dice que cuando un programa hace referencia a una posición de memoria, se espera que vuelva a hacerlo en poco tiempo.

A partir de esta idea, se busca tener las celdas mas usadas (y no todas) en una memoria mas inmediata, intermedia a la cpu y la memoria principal (ver figura 3), denominada **Memoria Caché**.



Figure 3: Jerarquía con Memoria Caché

La memoria caché tiene como objetivo proveer una velocidad de acceso cercana a la de los registros pero al mismo tiempo proveer un mayor tamaño de memoria. Para esto, la caché contiene una copia de porciones de la memoria principal, y en el momento de leer una celda, primero se verifica si se encuentra en la caché. **Si esto no ocurre** debe traerse de la memoria principal un bloque de celdas a la caché y luego la celda se entrega a la CPU. El principio de localidad dice que cuando un bloque de datos es traído a la caché porque se necesitó una de sus celdas, entonces es probable que próximamente se necesiten otras celdas del mismo bloque.

Fuente

1. *Organización y Arquitectura de Computadoras* , William Stallings, Capítulo 4
2. http://es.wikipedia.org/wiki/Unidad_de_estado_sólido