

Departamento de Arquitectura Instituto de Computación
Universidad de la República
Montevideo - Uruguay

Notas de Teórico

Arquitectura de von Neumann

Arquitectura de Computadoras
(Versión 4.2 - 2011)

13 ARQUITECTURA DE VON NEUMANN

13.1 Introducción

Esta arquitectura fue la propuesta por el matemático John von Neumann para la construcción de la computadora EDVAC en 1945 (la máquina se terminó de construir en 1949), sucesora de la que se considera la primer computadora electrónica, la ENIAC (1946).

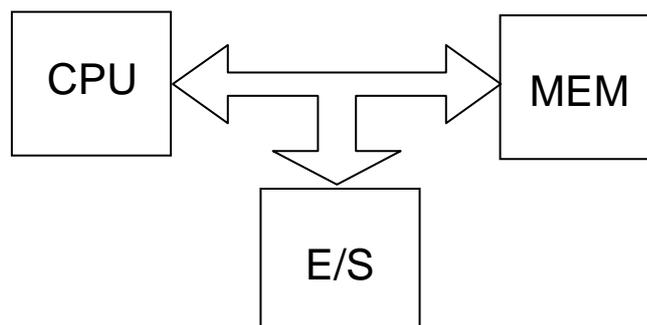
Los conceptos que propuso von Neumann han tenido una vigencia mucho más allá de la esperada en una industria como la de las tecnologías de la información. Por más que llame la atención todos los computadores modernos disponibles comercialmente poseen en el fondo la misma arquitectura que la EDVAC, definida en 1945. Mucho se ha adelantado en materia de los circuitos lógicos (la organización) que implementan la arquitectura, pero casi nada en los principios de diseño que planteó von Neumann hace más de 60 años (que en materia tecnológica representan más de una decena de generaciones completas).

Se puede decir que la idea fundamental de von Neumann se apoya en el hecho que una operación compleja normalmente se puede dividir como una secuencia ordenada de operaciones más simples. En otras palabras lo que propuso fue construir una máquina capaz de ejecutar algoritmos en forma explícita. Para ello introdujo el concepto de "programa almacenado" como una "secuencia lógicamente ordenada de instrucciones", siendo las "instrucciones" las operaciones básicas que implementa el hardware a través de sus circuitos lógicos.

Una Arquitectura de von Neumann tiene tres bloques constructivos básicos: la Unidad Central de Proceso (ó CPU por su sigla en inglés), la Memoria y la Entrada/Salida.

Las funciones de cada bloque son:

- **CPU:** se encarga de ejecutar los programas.
- **Memoria:** almacena el programa (conjunto de instrucciones ordenado lógicamente) y los datos (operadores y resultados de la ejecución de las instrucciones).
- **Entrada/Salida:** comunica el computador con el mundo exterior, permitiendo la interacción con los usuarios y con otras computadoras.



La **Memoria** se implementa como ya hemos visto en los capítulos anteriores. En un próximo capítulo veremos en detalle el sub-sistema de **Entrada/Salida**. En este capítulo nos concentraremos, entonces, en analizar la **CPU**.

Estos tres sub-sistemas se interconectan por medio de un **bus**, que contiene líneas de **datos**, de **dirección** y de **control**.

13.2 CPU

La **Unidad Central de Proceso** es la encargada de la ejecución de las **instrucciones**.

Para ello la arquitectura de von Neumann propone que la misma cuente con los siguientes elementos básicos:

- **Unidad Aritmética-Lógica (ALU):** circuito lógico que implementa operaciones de aritmética binaria (típicamente la operaciones básicas para representaciones "binario" y "complemento a 2" y lógicas (típicamente AND, OR, EXOR y NOT bit a bit). También es habitual que implemente operaciones de desplazamiento y/o rotación de bits.
- **Unidad de Control:** circuito secuencial que implementa el denominado "ciclo de instrucción", permitiendo acceder a la siguiente instrucción de un programa, leer sus operandos, efectuar la operación indicada en la ALU y guardar el resultado de la misma.
- **Banco (Set / Conjunto) de Registros:** una serie de posiciones especiales de memoria, ubicadas físicamente dentro de la propia CPU, que permiten un acceso a operandos y lugares de almacenamiento de resultados mucho más veloz que si estuvieran en el sistema de memoria normal. Algunos de estos registros son de uso interno de la propia CPU (más precisamente de su Unidad de Control) y otros son accesibles y utilizables por el programador.

Notemos que el Banco (Set / Conjunto) de Registros no es un elemento esencial de la arquitectura. Es un elemento de carácter tecnológico que tiene que ver con dos aspectos:

- por razones vinculadas a la electrónica, la velocidad de transferencia de información entre dos circuitos está limitada por la distancia geográfica entre los mismos
- la memoria rápida (usada en los pocos registros de la CPU) es más cara que la memoria más lenta (usada en las mucho más abundantes posiciones del sistema de memoria principal).

De allí que se utilice poca memoria rápida (por el costo) y que se ubique en las cercanías de donde se va a utilizar (dentro de la CPU) para minimizar su tiempo de acceso.

13.3 Caracterización de una Arquitectura von Neumann

Como ya expresamos todas las computadoras actuales tienen en esencia la misma arquitectura. Sin embargo dentro de la idea general de von Neumann que todas respetan, se diferencian entre sí por decisiones de diseño que afectan la cantidad y calidad de sus elementos componentes.

Así una arquitectura particular (ej: Intel x86, PowerPC, SPARC, MIPS, etc) establece en forma diferenciada los siguientes elementos característicos, los que deben ser conocidos por los programadores "de bajo nivel" para poder escribir programas para una de esas arquitecturas:

- **Set de Instrucciones:** la cantidad de instrucciones disponibles y la calidad y complejidad de las operaciones implementadas en el hardware de la CPU.
- **Formato de Instrucción:** la forma en que se codifican las instrucciones.
- **Set de Registros:** la cantidad de registros disponibles al programador, así como la función que pueden cumplir (ej: origen o destino de operaciones, almacenamiento de direcciones de operandos).

- **Modos de Direccionamiento:** formas de generar las direcciones para hallar los operandos o almacenar los resultados de las operaciones.
- **Manejo de la Entrada/Salida:** forma de comunicación con los "periféricos" (dispositivos que implementan la interacción con el mundo exterior).
- **Manejo de Interrupciones:** manejo de una forma particular de invocar a ciertas sub-rutinas de los programas que estudiaremos mas adelante.

13.3.1 Set de Instrucciones

Históricamente desde la EDVAC hasta principios de la década del 80 la tendencia fue ir, progresivamente, aumentando la cantidad de instrucciones disponibles en el hardware y su complejidad. Es así que las arquitecturas más difundidas a fines de los 70 (Intel 8080, Motorola 6800, Digital VAX, IBM 370) disponían de cientos de instrucciones y soportaban por hardware operaciones complejas tales como la comparación de strings y la búsqueda de un elemento en un array.

En ese momento se pensaba que cuanto más "potente" era el hardware (en cuanto a la variedad de operaciones disponibles), más eficientemente se podrían ejecutar los programas, logrando una mayor "performance" del sistema.

Sin embargo en 1980 se publicaron, casi simultáneamente, trabajos de investigación de las universidades de Berkeley y Stanford (ambas en California, EE.UU.) que proponían un enfoque radicalmente distinto: lo mejor era disponer de un conjunto mínimo de instrucciones que estuvieran implementadas en forma óptima. Estos trabajos acuñaron el término **RISC (Reduced Instruction Set Computer)** para referirse a los diseños basados en este concepto. A partir de ese momento el resto de las arquitecturas pasaron a denominarse, por contraposición, **CISC (Complex Instruction Set Computer)**.

Durante muchos meses corrieron ríos de tinta en publicaciones de todo tipo discutiendo cuáles arquitecturas eran mejor: si las basadas en la filosofía RISC o las CISC. La historia de esta batalla "tecnológica" tiene aspectos muy interesantes y un resultado doblemente paradójico: si bien las primeras implementaciones prácticas de RISC (incluyendo un chip diseñado por un grupo de estudiantes de Berkeley) demostraron que las RISC eran mucho mas eficientes (considerando su desempeño en relación con la frecuencia del reloj utilizado), el mercado siguió siendo de las CISC, en particular de la que se convirtió en un estándar de facto: la arquitectura x86 de Intel (la de los actuales procesadores "Pentium" y "Core"). La segunda paradoja es que para lograr los niveles actuales de desempeño, Intel debió incorporar conceptos de RISC en su diseño. Hoy se puede decir que un procesador Pentium es, en el fondo, una máquina RISC que ejecuta un conjunto de "micro-rutinas" que "emulan" las instrucciones de una máquina de arquitectura x86.

13.3.2 Formato de Instrucciones

El formato de las instrucciones refiere a la codificación de las distintas instrucciones para su almacenamiento en la memoria del sistema.

Al igual que para los demás tipos de datos manipulados que vimos (caracteres, números) los computadores trabajan con una representación de las instrucciones mediante un código binario.

El código binario reserva una serie de bits para identificar la operación realizada por la instrucción, otros indican los operandos a utilizar y sus direcciones, así como la indicación de dónde se almacena el resultado.

Los atributos que definen el formato de las instrucciones incluyen aspectos tales como si los códigos binarios asociados son de largo fijo o variable y si las instrucciones tienen operandos y destino independientes (en cuyo caso se habla de "arquitectura de 3 direcciones") ó solapados (correspondiente a "arquitectura de 2 direcciones").

Las arquitecturas RISC promocionaron la utilización de formatos de instrucción de largo fijo, como forma de simplificar la circuitería de decodificación y permitir la utilización eficaz de técnicas de optimización (pipeline, superescalaridad) que se verán más adelante. Las arquitecturas CISC normalmente tienen instrucciones de largo variable.

Las arquitecturas RISC también utilizan formatos de "3 direcciones", mientras que en las CISC es más común encontrar implementaciones de "2 direcciones" (como es el caso de Intel x86).

13.3.3 Set de Registros

Los registros son, como mencionamos, posiciones "especializadas" de memoria, ubicadas dentro de la propia CPU y que poseen una manera de direccionarlas distinta de la memoria "normal" del sistema.

Los primeros diseños siguieron el ejemplo de la EDVAC de von Neumann y utilizaron registros "con personalidad". En la propuesta original de von Neumann la computadora poseía, entre otros, un registro **Acumulador** para las operaciones aritméticas, un registro **Contador** para operaciones que implicaran contar y un registro **Indice** para ser utilizado para contener direcciones de operandos en memoria. Esas ideas son las que llevaron a Intel a nombrar los registros de su primer microprocesador 8080 con letras: A (Acumulador), B (Base, para direcciones), C (Contador, en instrucciones de string), D (Data, para almacenar datos, es decir operandos en general) y así sucesivamente. Luego la arquitectura x86 de 16 bits tomó esos nombres y le agregó una X (de eXtended, en referencia a que pasaron de ser de 8 bits a 16 bits), quedando AX, BX, CX y DX, a los que se agregaron SI (Source Index) y DI (Destination Index) entre otros.

La característica de los registros "con personalidad" es que su función dentro de las instrucciones (como operandos fuente, operandos destino, contadores o parte del direccionamiento de los operandos) está condicionada y no todos sirven para cualquier función. Su aplicabilidad es restringida y específica.

Las arquitecturas RISC promovieron desde el comienzo el uso de registros "despersonalizados", de uso general, aplicables todos y cada uno a cualquiera de las funciones antedichas. Intel incorporó la idea a partir del procesador 80386, primer representante de 32 bits de su familia arquitectónica. Los registros de 32 bits de la arquitectura x86 se denominan EAX, EBX, ECX, EDX, ESI, EDI, etc (la "E" significa "Enhanced") y pueden intercambiarse casi completamente en su uso en las instrucciones.

13.3.4 Modos de Direccionamiento

Los modos de direccionamiento establecen las formas en que se puede, a nivel de las instrucciones, especificar la dirección de un operando o del lugar donde colocar el resultado de la operación correspondiente a la instrucción.

Existen tres modos de direccionamiento básicos: Inmediato, Directo e Indirecto. A continuación presentaremos las características de cada uno:

- **Inmediato:** en este modo en la instrucción se encuentra el propio operando (su valor). Se utiliza típicamente para constantes (ya que las instrucciones no se pueden modificar).
- **Directo:** en este modo en la instrucción se encuentra la dirección del operando. Se pueden distinguir dos tipos:
 - Directo a Registro: el operando está almacenado en un registro y la instrucción contiene el identificador del registro
 - Directo a Memoria: el operando está almacenado en memoria y la instrucción contiene la dirección donde se encuentra
- **Indirecto:** en este modo en la instrucción se encuentra la dirección del lugar

donde se encuentra la dirección del operando. También se pueden distinguir dos tipos:

- Indirecto por Registro: el operando está almacenado en una posición de memoria cuya dirección se encuentra en un registro y la instrucción contiene el identificador del registro. De los dos tipos de indirecto éste es el habitualmente implementado, utilizando mas de un registro (sumando sus contenidos) para formar la dirección. En algunos casos la arquitectura prevé multiplicar el contenido del registro por el tamaño en bytes del operando. En ese caso se habla de direccionamiento **indizado**.
- Indirecto por Memoria: el operando está almacenado en una posición de memoria cuya dirección está en otra posición de memoria y la instrucción contiene la dirección de esta última. Esta variante no es implementada en las arquitecturas prácticas disponibles.

Los modos Directo e Indirecto (especialmente en su modalidad Indizado) se suelen combinar para establecer la dirección de un operando o de un destino. La dirección se forma sumando el contenido de uno o mas registros identificados en la instrucción con la "dirección base" contenida en la misma. Esto permite implementar fácilmente el acceso a estructuras tipo array. Se coloca la dirección de comienzo de la estructura en la instrucción (modo directo) y se cambia el contenido de un registro "índice" (que también debe estar referenciado en la instrucción) para ir recorriendo el array haciendo un "loop" y reutilizando la misma instrucción.

13.3.5 Manejo de Entrada/Salida

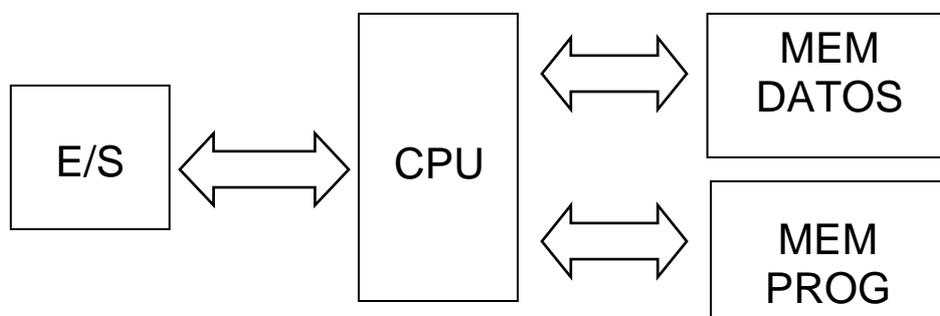
Más adelante veremos un capítulo completo dedicado a este tema. Por ahora interesa marcar que las distintas arquitecturas se distinguen en este punto básicamente por el hecho de tener o no instrucciones específicas, con su propio espacio de direccionamiento, para la comunicación con los dispositivos que permiten la interacción con el "mundo exterior".

13.3.6 Manejo de Interrupciones

Nos referiremos a este tema en un capítulo específico.

13.4 Arquitectura Harvard

La arquitectura Harvard se caracteriza porque existen, a diferencia de la von Neumann, dos unidades de memoria separadas: una para los datos y otra para las instrucciones:



Muchas veces se presenta la arquitectura Harvard como algo diferente de la arquitectura von Neumann. Sin embargo en el fondo no cambia la esencia de la von Neumann: siguen existiendo “programas” formados por una secuencia lógica de “instrucciones” las que son ejecutadas por una “CPU”.

Actualmente muchos procesadores implementan una **arquitectura Harvard modificada**, que consiste en mantener ciertos caminos separados y paralelos para acceder a las instrucciones en forma simultánea con los datos, pero permitiendo que los datos y las instrucciones están almacenados en una memoria común. Los detalles de estas implementaciones se verán más adelante en el curso.