



Universidad de Buenos Aires



Diseño, validación e implementación de una arquitectura RISC

por

Luciano César Natañe

Tesis presentada para optar al Título de

Ingeniero Electrónico

por la

Facultad de Ingeniería de la Universidad de Buenos Aires

Director:

Ing. Nicolás Alvarez

Co-Director:

Ing. Octavio Alpago

Miembros del Jurado:

Ing. XXXXXXXXXXXXXXXXXXXXXXXX

Ing. XXXXXXXXXXXXXXXXXXXXXXXX

Ing. XXXXXXXXXXXXXXXXXXXXXXXX

Calificación: _____

Fecha: _____

Universidad de Buenos Aires, Facultad de Ingeniería

Diseño, validación e implementación de una arquitectura RISC

por

Luciano César Natale

Resumen

El presente trabajo constituye la Tesis de Grado necesaria para obtener el título de Ingeniero Electrónico de la Facultad de Ingeniería de la Universidad de Buenos Aires.

El objetivo de este trabajo es el diseño, la validación y la implementación de una arquitectura RISC con el objetivo de generar un núcleo de procesamiento, sintetizable en FPGA, altamente configurable y suficientemente flexible y sencillo para ser utilizado en distintas aplicaciones dentro del ámbito de la investigación en los Laboratorios de Microelectrónica y de Sistemas Embebidos. Se presenta la teoría e historia necesaria para poder explicar las decisiones de diseño adoptadas en el desarrollo de la arquitectura. Se presentan vectores de prueba para las validaciones posteriores. Se verifica el diseño generado mediante emuladores. Se implementa el diseño en un lenguaje de descripción de hardware. Se sintetiza en FPGA el diseño y se contrastan los resultados con las emulaciones realizadas. Finalmente se analizan los resultados obtenidos, se contrastan los mismos contra trabajos similares, y se proponen trabajos futuros.

Agradecimientos

Agradecimientos Lucho.

Dedicatoria Lucho

Índice

1. Introducción al trabajo de tesis	1
1.1. Objetivo	2
1.2. Alcance	3
1.3. Organización del trabajo	4
2. Marco teórico: arquitectura de procesadores	5
2.1. Perspectiva histórica	6
2.1.1. La máquina de Turing	7
2.1.2. La revolución digital	7
2.1.3. Silicio, dispositivos semiconductores y transistores	8
2.1.4. Circuitos integrados y tecnología CMOS	9
2.1.5. Microprocesadores	11
2.2. Evolución hacia las arquitecturas modernas	15
2.2.1. Clasificación de las arquitecturas según el acceso a la memoria de instrucciones y datos	15
2.2.2. Clasificación de las arquitecturas según su conjunto de instrucciones	20
2.2.3. Taxonomía de Flynn	21

Índice de Tablas

Índice de Figuras

2.1.	Esquema simplificado de un “sistema de procesamiento”	6
2.2.	Primer transistor desarrollado por John Bardeen y Walter Brattain bajo la dirección de William Shockley en los Laboratorios Bell de AT&T en el año 1947	9
2.3.	Primer circuito integrado presentado por Texas Instruments, Inc. el 12 de Septiembre de 1958	11
2.4.	Grumman F14 Tomcat, avión de combate de la US Navy cuyos sistemas de vuelo eran controlados por la CADC.	12
2.5.	Imágen microscópica de la AL1 desarrollada por Four-Phase Systems Inc..	13
2.6.	Imágen microscópica del PICO1/GI250- Desarrollo colaborativo entre Pico y General Instruments de <i>chip</i> único para la Monroe/Litton Royal Digital III Calculatator.	14
2.7.	Encapsulado del Intel 4004; que es considerado el primer microprocesador comercial en ser introducido en el mercado.	15
2.8.	Carátula del escrito de von Neumann First Draft of a Report on the EDVAC	16
2.9.	Cinta de papel perforado con instrucciones de la Mark I	20

Capítulo 1

Introducción al trabajo de tesis

El presente trabajo se encuentra enfocado en el contexto del diseño de hardware digital. El mismo fue motivado por la necesidad de contar un con núcleo de procesamiento altamente configurable y suficientemente flexible y sencillo para distintas aplicaciones dentro del ámbito de la investigación en los Laboratorios de Microelectrónica y de Sistemas Embebidos; sintetizable en ¹. La arquitectura a desarrollar será del tipo RISC².

Desde la aparición de los microprocesadores a mediados de los años 70, la tendencia fue el aumento de la complejidad de las arquitecturas, generando un efecto de “bola de nieve”, al ir superponiendo capas sobre un núcleo central. Existió, entonces, una reacción adversa a esta tendencia. Por ejemplo, la arquitectura experimental de IBM 801; y también en Berkeley, Patterson y Ditzel fueron los primeros en acuñar el término RISC, para describir una nueva clase de arquitectura que deshacía el camino del resto de las arquitecturas hasta el momento, conocidas, en contraposición, como CISC³. A partir de este antecedente, los principales fabricantes de microprocesadores han lanzado al mercado sus propias implementaciones basadas en los principios establecidos en IBM y Berkeley.

El concepto de las arquitecturas RISC se basa, principalmente, en el hecho de que al simplificar la lógica necesaria para la ejecución de una instrucción permite aumentar la frecuencia de operación de las compuertas que componen la lógica. Además, es posible dividir la ejecución de las instrucciones en etapas sencillas y consecutivas, permitiendo de esta manera implementar fácilmente optimizaciones como, por ejemplo, una arquitec-

¹ *Field Programmable Gate Array*: dispositivo electrónico de compuertas lógicas reprogramables en campo

² *Reduced Instruction Set Computer*. Técnica de diseño de unidades de procesamiento basada en el hecho de que un conjunto de instrucciones simple provee una mayor performance al ser combinado con una arquitectura capaz de ejecutar dichas instrucciones en algunos pocos ciclos de máquina.

³ *Complex Instruction Set Computer*. Técnica de diseño de unidades de procesamiento basadas en el hecho de que el conjunto de instrucciones debe ser lo más poderoso posible.

tura de *pipeline*⁴. Es por esto que el conjunto de instrucciones es sencillo, permitiendo solamente operaciones básicas entre registros internos del microprocesador. El trabajo realizado por cada instrucción, en general, es menor que el generado por una instrucción CISC, pero se hace de manera sencilla y rápida. Es importante notar que no solamente la ganancia radica en poder aumentar la frecuencia de operación de la lógica, sino que estas condiciones facilitan el desarrollo de diseños de bajo consumo, característica muy valorada en el nicho de los sistemas embebidos.

El mercado de los sistemas embebidos es excesivamente amplio y está inserto en todas las industrias. En un automóvil, por ejemplo, podemos encontrar microprocesadores en el sistema de frenos, en la central de inyección electrónica, en el sistema de entretenimiento y navegación, etc. La otra arista de vital importancia para el mercado de los sistemas embebidos, es el de los dispositivos móviles, donde se vuelve vital el requerimiento de bajo consumo. Estamos viviendo la revolución de IoT⁵, que se trata básicamente de sistemas embebidos autónomos que están conectados a “la nube” y pueden ser monitoreados y controlados remotamente a través de *Internet*.

Dentro del universo de las arquitecturas RISC, actualmente se destacan dos: MIPS y ARM. La primera, fue desarrollada por un grupo de investigadores de la Universidad de Stanford (entre ellos John L. Hennessy, pionero del concepto RISC junto a David Patterson, coautores de la bibliografía más relevante del área). Esta arquitectura, por su sencillez, es la predilecta al momento del desarrollo de cursos enfocados en la enseñanza de arquitectura de computadoras. Si bien MIPS posee gran relevancia académica, es muy popular en el mercado de los microprocesadores en sistemas embebidos como equipos de telecomunicaciones, decodificadores de TV digital, y consolas de entretenimiento, con ejemplos muy conocidos como *Nintendo* y *PlayStation*. ARM, por otro lado, ha ganado una importante porción del mercado de los sistemas embebidos (con un gran aporte de los dispositivos móviles), basando su modelo de negocios en la venta de la propiedad intelectual (IP, *intellectual property*) del diseño de los microprocesadores a las empresas que finalmente producen el microprocesador.

1.1. Objetivo

La Tesis tiene como objetivo principal el diseño, la validación e implementación de una arquitectura RISC y su conjunto de instrucciones.

El enfoque de la tesis se basará en un desarrollo teórico del conjunto de instrucciones y de las características de la arquitectura; y en el desarrollo práctico del emulador y la implementación en lenguaje descriptor de hardware.

El concepto central detrás del desarrollo será el de **ortogonalidad**. Esto implica, por una

⁴ Técnica de diseño de arquitecturas de computadoras en la que se segmenta la ejecución de las instrucciones en múltiples etapas, permitiendo que múltiples instrucciones estén ejecutándose en paralelo.

⁵ Es un concepto que se refiere a la interconexión digital de objetos cotidianos con internet.

parte, que los bloques constructivos de la arquitectura que se repiten sean independientes e indiferenciables entre sí. Por otra parte, los formatos de las instrucciones, en la medida de lo posible, se diseñaran de manera tal que se pueda mantener el mismo ancho de campo para los datos inmediatos y los desplazamientos (excepto en los casos donde es explícitamente conveniente agrandarlos sin penalizar la complejidad del diseño).

El objetivo perseguido va a ser el de mantener la sencillez y la ortogonalidad, favoreciendo así la simplificación de la implementación. Se trabajará en el desarrollo de la definición de la arquitectura y su conjunto de instrucciones en favor de este objetivo. Se definirá la interfaz física para la conectividad con periféricos, los tipos de datos que maneja la arquitectura, la cantidad y tipos de registros internos, el acceso a memoria de programa y de datos con su organización y modo de direccionamiento, la interfaz con la ALU⁶ y la FPU⁷, mecanismos de manejos de excepciones e interrupciones, modos de operación y manejo de periféricos. Luego se definirá el conjunto de instrucciones que ejecutará la arquitectura.

Una vez definida la arquitectura y su conjunto de instrucciones, se procederá a diseñar los vectores de prueba para poder validar las implementaciones. Se desarrollará un emulador de la arquitectura que deberá validar los vectores de prueba diseñados. Una vez concluida esta etapa, se implementará a nivel RTL el diseño en *Verilog*. Este diseño será validado mediante simulaciones y utilizando dispositivos programables. Se validará también contra los vectores de prueba. Se analizarán los recursos utilizados en dispositivos FPGA. Se realizará un análisis comparativo entre la arquitectura desarrollada y otras arquitecturas RISC.

1.2. Alcance

Como resultados a obtener de la tesis se tienen los siguientes:

- Especificación completa de la arquitectura
- Vectores de prueba
- Emulador de la arquitectura
- *IP Core* codificado en el lenguaje *Verilog* de la arquitectura completa
- Resultado de los vectores de prueba tanto en el emulador como en el *IP Core*
- Análisis comparativo entre la arquitectura desarrollada y otras arquitecturas RISC
- Proposición de trabajos futuros y/o mejoras.

⁶ *Arithmetic Logic Unit*. Bloque constructivo encargado de realizar las operaciones aritmético lógicas sobre los datos.

⁷ *Floating Point Unit*. Bloque constructivo encargado de realizar las operaciones en punto flotante sobre los datos.

1.3. Organización del trabajo

En esta sección se describe la organización de la presente tesis. Con el objetivo de que la misma sea autocontenida, los primeros capítulos se ocupan de presentar las bases o conocimientos necesarios para comprender la totalidad del trabajo.

El desarrollo de la tesis se organiza de la siguiente forma:

- En el capítulo 2 se presentará la teoría general de las arquitecturas de procesadores y una revisión histórica sobre el tema. Se estudiará la diferenciación entre los universos de procesadores CISC y RISC y se justificará la elección de diseñar una arquitectura RISC para la tesis. Se presentarán las técnicas de diseño de arquitecturas estudiadas. Además se presentarán reseñas de otras arquitecturas actuales y sus decisiones de diseño, para luego contrastarlas con los objetivos perseguidos por el presente trabajo.
- En el capítulo 3 se presentará la especificación completa de la arquitectura diseñada, explicitando los criterios y las decisiones de diseño tomadas. Además se presentará el diseño de los vectores de prueba que se utilizarán para validar las implementaciones de la arquitectura.
- En el capítulo 4 se desarrollarán las implementaciones del emulador de la arquitectura y del *IP Core* en RTL. Dicho RTL cumplirá con ciertas condiciones de portabilidad y legibilidad del código, para que el mismo sea efectivamente un IP core. Se evaluará y validará el *IP Core* utilizando simuladores. Se explicitarán las decisiones de diseño necesarias para pasar de la abstracción del diseño a la implementación real.
- En el capítulo 5 se validarán las implementaciones del capítulo 4 mediante los vectores de prueba diseñados para el capítulo 3. El *IP Core* será sintetizado para distintos dispositivos FPGA. Se analizará en cada caso el consumo de recursos utilizados, máxima frecuencia de operación y la potencia consumida.
- En el capítulo 6 se extraerán las conclusiones pertinentes sobre los resultados obtenidos y se propondrán futuras mejoras de la arquitecturas a partir del análisis realizado.

Capítulo 2

Marco teórico: arquitectura de procesadores

En líneas generales un procesador es un sistema que permite, por un lado ingresar instrucciones y datos obteniendo en consecuencia los resultados de operar lo indicado en las instrucciones sobre los datos. No es necesario para tal fin, definir ninguna tecnología que soporte este comportamiento; se trata más bien de un desarrollo teórico. Dicho desarrollo implica distinguir las distintas partes que lo componen. En una primera aproximación, en un sistema de procesamiento podemos distinguir cuatro componentes:

- Datos de entrada
- Instrucciones
- Unidad de procesamiento
- Datos de salida

Estos componentes se relacionan de la forma mostrada en la figura 2.1. Debe notarse que los datos de entrada y las instrucciones ingresan a la unidad de procesamiento, la cual genera datos de salida como resultado de operar las instrucciones sobre los datos de entrada.

Los datos de entrada representan el dominio sobre el cual puede operar la unidad de procesamiento. Para definir un sistema de procesamiento debemos especificar, entonces, cuál es el dominio que puede manejar. Dicha definición deberá especificar el formato y cantidad de datos que acepta la unidad de procesamiento, así como los mecanismos para ingresarlos al sistema.

Las instrucciones definirán las operaciones que la unidad de procesamiento puede realizar sobre los datos de entrada. Por lo tanto, para definir un sistema de procesamiento, debemos definir qué operaciones será capaz de realizar sobre el dominio de entrada del

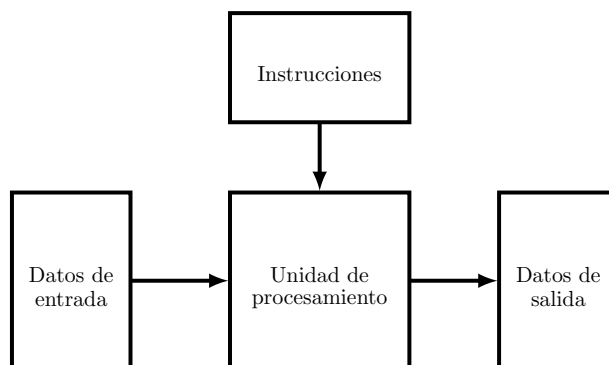


Figura 2.1: Esquema simplificado de un “sistema de procesamiento”

mismo. A su vez, se debe definir la salida esperada de las instrucciones; es por eso que al definir las instrucciones estamos definiendo intrínsecamente el dominio de salida del sistema. Es importante notar que la información sobre las instrucciones también representa una entrada para el sistema, es por eso que, como en el caso de los datos de entrada, se deberá especificar formato y cantidad que acepta la unidad de procesamiento, así como los mecanismos para ingresarlas.

Los datos de salida representan el dominio sobre el cuál las instrucciones vuelcan el resultado de las operaciones realizadas sobre los datos de entrada. El dominio de salida, como fue notado antes, queda definido al definir las instrucciones, pero debe definirse el formato y cantidad de los mismos, así como los mecanismos necesarios para extraerlos del sistema.

En este contexto la unidad de procesamiento es la encargada de recibir las instrucciones y datos, ejecutar las operaciones para finalmente presentar los resultados.

Un microprocesador es -en efecto- una implementación de un sistema de procesamiento. El sustento físico de dicha implementación es la microelectrónica. En las siguientes secciones de este capítulo, se repasará el marco histórico en el que surgen los microprocesadores y las arquitecturas asociadas.

2.1. Perspectiva histórica

Desde épocas remotas, el hombre se ha destacado en el mundo animal por su capacidad de modificar su entorno para resolver problemas recurrentes. Es esta capacidad la fortaleza de la especie en la naturaleza. La película “2001: Odisea del espacio”¹ narra, desde un enfoque particular, parte de la evolución del ser humano, o al menos la interpretación de los autores sobre la misma. En la misma, ubicándose temporalmente varios millones de años atrás, un clan de cavernícolas prehumanos intentan sobrevivir en condiciones

¹2001: Odisea en el espacio es un film del año 1968 dirigida por Stanley Kubrick basada en la novela de Arthur C. Clarke.

extremas. Comen los pocos hierbajos que pueden encontrar en el desolado paisaje, hierbajos que para colmo han de compartir con una manada de tapires que habita la misma zona. La única fuente de agua del clan —un simple charco— les es arrebatada por un clan rival. Por si fuera poco, este desdichado clan vive permanentemente amenazado por un leopardo que domina la región y que de vez en cuando caza a alguno de sus miembros. En resumen: este grupo de homínidos padece hambre, frío y miedo, y parecen condenados a una segura extinción. En ese contexto, y por motivos que no vienen al caso, aparece uno de los cavernícolas contemplando el esqueleto de un animal. Parece reflexionar sobre lo que tiene delante, como si estuviese viéndolo desde una nueva perspectiva. Hay algo nuevo en aquellos huesos. Algo que hasta entonces ni él ni ninguno de sus congéneres habían visto. Los huesos que hay tirados por el suelo pueden ser usados. El cavernícola toma el más robusto de los huesos y empieza a golpear el esqueleto; primero con precaución, más tarde con fuerza, hasta que termina consumido por un frenesí violento. Este cavernícola acaba de descubrir el primer arma —la primera herramienta— de la historia. O dicho de otro modo, acaba de aparecer el primer ser humano sobre la faz de la tierra. Gracias al uso del hueso —o de herramientas similares como palos o piedras— el clan que estaba a punto de extinguirse descubre que puede cazar a los tapires con los que convive y comérselos. Así que sus problemas de hambre han terminado. También gracias a sus armas pueden atacar al clan rival y recuperar el charco de agua, lo que soluciona también sus problemas de sed. Y deducimos que serán capaces incluso de defenderse del peligroso leopardo. Los miembros del clan ya no son prehumanos indefensos; ahora son humanos armados.

Esta cita cinematográfica pretende graficar la diferenciación del ser humano en la cadena alimenticia. Gracias al poder de la observación y el razonamiento hemos sido capaces de modificar nuestro entorno para asegurar la supervivencia de la especie en un mundo en el que la misma se encontraba en clara desventaja. En este contexto, el ser humano ha sido capaz de generar un desarrollo tecnológico, que hoy en día es vertiginoso.

2.1.1. La máquina de Turing

Alan Mathison Turing² es considerado el padre de la computación y la inteligencia artificial.

2.1.2. La revolución digital

La revolución digital es considerada la tercera revolución industrial. Sus comienzos se remontan a fines de los años 50. La adopción y proliferación de las computadoras digita-

²Alan Mathison Turing (23 June 1912 – 7 June 1954): Típicamente considerado el padre de las ciencias de la computación y de la inteligencia artificial, fue un biólogo, matemático, lógico, científico de la computación y de la criptografía que formalizó los conceptos de algoritmo y computación mediante la “Máquina de Turing”, que es un módulo de una computadora de propósito general

les y el mantenimiento de registros digitales de información son las características que la definen. De manera implícita, esta revolución está relacionada con los cambios radicales provocados por la computación y las tecnologías de telecomunicaciones. Esta revolución dió origen a lo que hoy conocemos como la “era de la información”. En el corazón de este proceso encontramos dos componentes tecnológicas. Por un lado está la teoría, que puede remontarse a los primeros análisis matemáticos de la lógica, como el álgebra de Boole introducida por primera vez en un pequeño folleto publicado en 1847 bajo el nombre *The Mathematical Analysis of Logic* y la posterior publicación del libro *An Investigation of the Laws of Thought on Which are Founded the Mathematical Theories of Logic and Probabilities*, publicado en 1854. En estas publicaciones George Boole³ pretendió utilizar técnicas algebraicas para tratar expresiones de la lógica proposicional. En la actualidad, el álgebra de Boole se aplica de forma generalizada en el ámbito del diseño electrónico. Claude Shannon⁴ fué el primero en aplicarla en el diseño de circuitos de conmutación eléctrica biestables, en 1948. Por el otro lado, se encuentra la revolución del silicio y la producción en masa y uso generalizado de circuitos lógicos digitales que permitieron el desarrollo en gran escala de circuitos electrónicos que implementan funciones lógicas basados en los desarrollos de Boole. Aproximadamente cien años pasaron desde que Boole publicó su teoría, hasta que la tecnología encontró el camino para implementar esos conocimientos de forma práctica y útil. La invención del transistor data del año 1947. En la figura 2.2 se muestra una imagen de una réplica de este primer transistor de la historia. Este invento fué el que permitió la creación de equipos digitales avanzados. En el contexto de la lógica digital, los transistores son utilizados como llaves de conmutación que permiten o no el paso de una señal eléctrica al ser excitados por otra señal. Previo a los transistores, la lógica era implementada con componentes electromecánicos y válvulas termoiónicas de vacío; tecnologías que por su naturaleza eran de dimensiones y consumos energéticos elevados y poco fiables. Para poner esta problemática en perspectiva, comparamos la primera computadora de propósito general electrónica, llamada ENIAC⁵ implementada con 18000 válvulas, con un consumo de 160 kW, capaz de realizar 5000 sumas/s, 385 multiplicaciones/s, con 5 millones de soldaduras y un peso de 30 Tn, contra un procesador Intel Core I7 que permite 177000 MIPS con un consumo de aproximadamente 100 W, es decir, 35 millones de veces más rápido. Para hacer la misma cantidad de sumas por segundo se requerirían 5 GW con la ENIAC; considerar que la capacidad actual instalada en argentina es de 30 GW.

2.1.3. Silicio, dispositivos semiconductores y transistores

Los semiconductores son materiales, que como bien indica su nombre, no son del todo conductores. En ellos, la capacidad de conducir una corriente eléctrica puede ser manipulada de diversas maneras. En 1931 Wolfgang Pauli -quien en 1945 fue premiado con

³Nota sobre George Boole

⁴Nota sobre Claude Shannon

⁵Descripción de la ENIAC

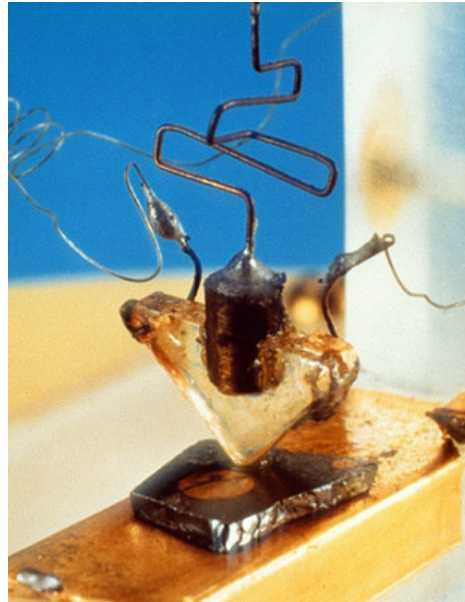


Figura 2.2: Primer transistor desarrollado por John Bardeen y Walter Brattain bajo la dirección de William Shockley en los Laboratorios Bell de AT&T en el año 1947

el Nobel de Física- enunció: “Uno no debería trabajar en semiconductores, eso es un lío deleznable, quién sabe si realmente existen”. Nada más alejado de la realidad que hoy nos rodea. Los materiales semiconductores permitieron la creación de los transistores y posteriormente los circuitos integrados, dispositivos que iniciaron la revolución digital. Los primeros dispositivos fueron fabricados sobre germanio y arseniuro de galio. Incluso, el primer circuito integrado, fue realizado en germanio. Pero fue el silicio el material que realmente revolucionó la industria, por su alta disponibilidad en la naturaleza y relativa fácil manipulación para la fabricación de dispositivos semiconductores en circuitos integrados. Puede afirmarse que el silicio es uno de los materiales mejor conocidos por el ser humano. Hace más de 50 años que se lo estudia en detalle para mejorar la tecnología, logrando importantes avances. También puede afirmarse que los transistores hoy en día es el bien más abundante en el mundo. Una comparativa del año 2012 muestra que durante año se produjeron en el orden de 10^{17} granos de arroz en la tierra, mientras que en el mismo período se produjeron en el orden de 10^{19} transistores, según declaraciones de la *Semiconductor Industry Association* de los Estados Unidos de América. La clave de semejante número en la producción son los circuitos integrados.

2.1.4. Circuitos integrados y tecnología CMOS

En la figura 2.3 puede verse el primer circuito integrado de la historia. Medía aproximadamente media pulgada de ancho e implementaba dos transistores montados en una

barra de germanio. Nace así el concepto del *chip* o *microchip*: en inglés, *chip* significa “corte o fracción pequeño de un material duro” directamente relacionado con la técnica de fabricación de circuitos integrados donde, a grandes rasgos, una barra cilíndrica de cristal de silicio puro, es cortada en muy delgadas láminas en forma de discos, sobre las cuales se “imprimen” los circuitos integrados, repitiendo el mismo patrón múltiples veces en un mismo disco para finalmente cortar ese disco en diminutos fragmentos que contienen el diseño. La evolución de las técnicas de fabricación permitieron integrar en un mismo *chip* de silicio más de un dispositivo, permitiendo implementar circuitos relativamente complejos en una pequeña área de silicio, disminuyendo así los riesgos de fallas por interconexión entre dispositivos. A su vez, las técnicas de fabricación evolucionaron permitiendo escalar el tamaño de los dispositivos fabricados incrementando la cuenta de dispositivos integrados por unidad de área. Como consecuencia buscada de disminuir el tamaño de los dispositivos, se logró aumentar la velocidad máxima de conmutación que los mismos pueden lograr, redundando en generar lógica cada vez más compleja y rápida en un mismo *chip*. En simultáneo con estos avances se comenzaron a fabricar transistores MOS⁶, transistores de efecto de campo eléctrico, que como gran ventaja sobre los “antiguos” transistores de juntura, evitaban la disipación de potencia al mantenerse en un estado definido (encendido o apagado), es decir que sólo disipaban potencia al cambiar de estado. Es así que en el año 1978 aparece la tecnología CMOS⁷ la cual permitió elevar el nivel de integración en forma masiva, manteniendo bajos niveles de disipación de potencia. Debido a la relativa sencillez geométrica del diseño de los transistores MOS, se pueden reutilizar diseños escalándolos para las nuevas generaciones de tecnología. El 19 de abril de 1965 Gordon E. Moore⁸, cofundador de Intel⁹, estableció de forma empírica que la cantidad de dispositivos integrados en un circuito integrado se duplicaría cada año. Más tarde, en 1975, modificó su propia ley al corroborar que el ritmo bajaría, y que la capacidad de integración no se duplicaría cada 12 meses sino cada 24 meses aproximadamente. Esta progresión de crecimiento exponencial, duplicar la capacidad de los circuitos integrados cada dos años, es lo que se denomina ley de Moore. Sin embargo, en 2007 el propio Moore determinó una fecha de caducidad: “Mi ley dejará de cumplirse dentro de 10 o 15 años”, no obstante también aseveró que una nueva tecnología vendrá a suplir a la actual. El cumplimiento se ha podido constatar hasta hoy.

La ley de Moore, no es una ley en el sentido científico, sino más bien una observación del ritmo de avance de la industria de aquellos momentos. Al momento de publicar esas declaraciones, Moore trabajaba en los Laboratorios de Fairchild Semiconductor, donde trabajaba junto a Robert Noyce. Ellos fueron los fundadores de Intel en 1968. El ritmo de crecimiento de la industria de los Semiconductores dió lugar así, entre otras cosas, a la creación de los microprocesadores.

⁶Transistor Metal-Oxido-Semiconductor: Nota

⁷CMOS: Nota. Complementariedad. Sólo utiliza MOSFET-N y MOSFET-P para implementar cualquier función

⁸Nota sobre Gordon E. Moore

⁹Nota sobre Intel

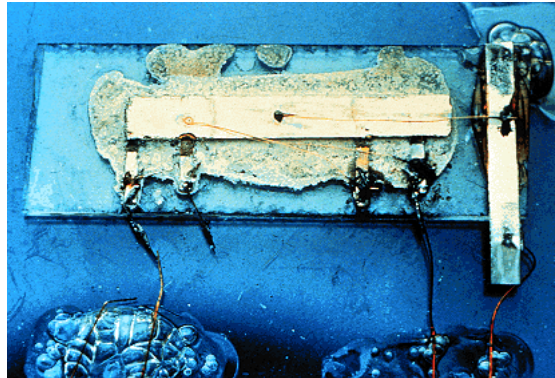


Figura 2.3: Primer circuito integrado presentado por Texas Instruments, Inc. el 12 de Septiembre de 1958

2.1.5. Microprocesadores

Los microprocesadores surgen como consecuencia del alto nivel de integración y complejidad que la tecnología de circuitos integrados permitieron alcanzar. Una computadora digital, podía ser construída en uno o algunos pocos *chips* de circuitos integrados. Los primeros diseños de microprocesadores propiamente dichos datan de fines de los años 60. Dentro de este contexto aparece el término *Central Processing Unit* (Unidad Central de Procesamiento) o CPU, que es la parte de la máquina encargada de tomar los datos de entrada y las instrucciones y generar los resultados, dentro del esquema del sistema de procesamiento planteado al principio de éste capítulo. El microprocesador integra además otros componentes y funcionalidades y se vale de ciertos periféricos que pueden estar o no integrados en el mismo circuito, como por ejemplo, memorias de sólo lectura (*Read Only Memory* o ROM¹⁰) y memorias de acceso aleatorio (*Random Access Memories* o RAM¹¹) y dispositivos de entrada / salida (*I/O, input output devices*). El diseño de arquitecturas de microprocesadores, como todo proceso de innovación tecnológica, no estuvo exento de discusiones y distintas vertientes y prácticas a seguir. Los primeros diseños de los que se tenga registro en orden cronológico, fueron:

CADC

En 1968, la armada de los Estados Unidos le encomienda a la empresa Garret Air Research¹² la producción de una computadora digital que pudiera competir con los sistemas electromecánicos que estaban bajo desarrollo para los sistemas de control de vuelo del nuevo avión de combate F-14 Tomcat. El diseño fue completado en 1970 y utilizaba un conjunto de chips (*chipset*) MOS como CPU, siendo aproximadamente 20 veces más pe-

¹⁰Nota sobre ROM

¹¹Nota sobre RAM

¹²Nota sobre Garret AiResearch



Figura 2.4: Grumman F14 Tomcat, avión de combate de la US Navy cuyos sistemas de vuelo eran controlados por la CADC.

queño que su contraparte electromecánica y a su vez, mucho más confiable. Este diseño fue utilizado en todos los primeros modelos de Tomcat. La armada se rehusó a publicar el diseño hasta 1997, es por eso que la CADC (Central Air Data Computer) y su *chipset* asociado MP944 no son muy conocidos y no son considerados en la mayoría de la bibliografía reelevante y es a partir de ese momento en que son reconocidos como los primeros diseños válidos de microprocesadores, y no sólo eso, si no que también son la primera especie de *Digital Signal Processor* (Procesador de Señales Digitales) o DSP¹³ puesto que además de las funciones de microprocesador, implementaba la funcionalidad de medir variables como altitud, velocidad vertical y número de *mach*¹⁴ a partir de datos de sensores pitot, presión estática y temperatura. Los responsables del diseño fueron Ray Holt¹⁵ y Steve Geller¹⁶

Four-Phase Systems AL1

Un diseño de Lee Boysel¹⁷ que data del año 1969 dentro de Texas Instruments¹⁸, era un componente de una implementación en partes de un microprocesador de 24 bits, compuesto por tres chips iguales: AL1. Durante un juicio por patentes, se demostró que una sólo AL1 junto con memorias ROM y RAM e I/O era capaz de funcionar como CPU.

¹³Nota sobre DSP

¹⁴Nota sobre Mach number

¹⁵Nota sobre Ray Holt

¹⁶Nota sobre Steve Geller

¹⁷Nota sobre Lee Boysel

¹⁸Nota sobre TI

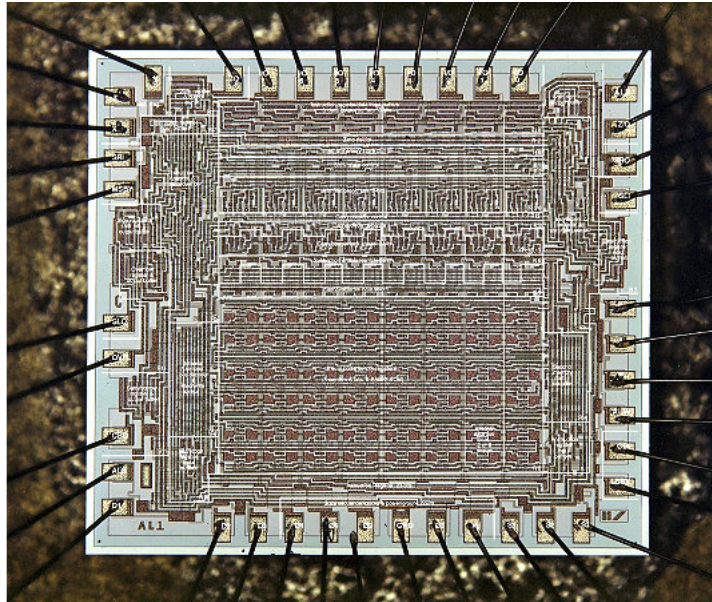


Figura 2.5: Imágen microscópica de la AL1 desarrollada por Four-Phase Systems Inc..

Pico/General Instruments

En 1971, Pico y General Instruments (GI) colaboraron en el diseño de circuitos integrados con el objetivo de fabricar una implementación de un diseño de *chip* único para la calculadora Monroe/Litton Royal Digital III Calcultator. Integraba en el mismo *chip*, memoria ROM y RAM llamado PICO1/GI250. Pico era un emprendimiento de cinco ingenieros de diseño de GI, que tenían la visión de crear esta arquitectura en un único *chip*. Contaban con experiencia previa en diseños de *chipsets* tanto de GI como de Marconi-Elliot¹⁹. Algunos de ellos habían trabajado para Elliot Automation para crear una computadora de 8 bits en tecnología MOS, y habían ayudado a establecer un laboratorio de investigación en tecnologías MOS en Escocia en el año 1967. El mercado de las calculadoras estaba en pleno auge, con lo cuál estos diseños supusieron un éxito comercial para Pico y GI. GI, por su parte continuó con la innovación en microprocesadores microcontroladores²⁰ y la división GI Microelectronics se convirtió en 1987 en Microchip PIC microcontroller.

Intel 4004

El Intel 4004 es generalmente el reconocido como el primer microprocesador comercial disponible en el mercado. El primer anuncio respecto de este dispositivo data del 15

¹⁹Nota sobre Marconi-Elliot

²⁰Nota sobre diferencia entre microprocesador y microcontrolador

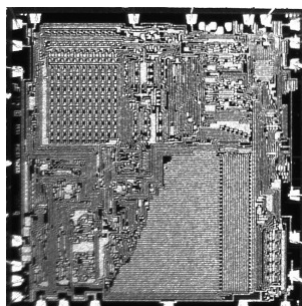


Figura 2.6: Imágen microscópica del PICO1/GI250- Desarrollo colaborativo entre Pico y General Instruments de *chip* único para la Monroe/Litton Royal Digital III Calculator.

de Noviembre de 1971, en la publicación *Electronic News*²¹. El diseño nace a partir del requerimiento de una compañía japonesa fabricante de calculadoras llamado Busicom²², que le encomienda a Intel la tarea de desarrollar un *chipset* para calculadoras de escritorio de alto rendimiento. El diseño original requerido por Busicom especificaba un *chipset* compuesto por 7 *chips* diferentes para la realización de una CPU de propósito específico cuyo programa estuviera almacenado en una ROM y sus datos guardados memoria de lectura y escritura implementada con registros de desplazamiento. Intel le asignó el proyecto a Ted Hoff²³, quien propuso simplificar el diseño utilizando almacenamiento en memoria RAM dinámica y una arquitectura de CPU de propósito general. El diseño propuesto por Hoff implementaba la solución en 4 *chips*: uno de ROM para almacenar el programa, uno de RAM dinámica para almacenar los datos, un dispositivo sencillo de I/O y una CPU de 4 bits. A pesar de no ser específicamente un diseñador de circuitos integrados, el supuso que se podía integrar todo el CPU en un único *chip*. Estas especificaciones surgieron de la interacción de Hoff con un empleado a su cargo, el ingeniero de software llamado Stanley Mazor²⁴ y un ingeniero de Busicom llamado Masatoshi Shima²⁵. Este proyecto se llamó MCS-4 y no fué hasta que Intel contratara al italiano Federico Faggin²⁶ que empezó a tomar su forma definitiva. Faggin venía de desarrollar en 1968 en Fairchild Semiconductor una tecnología de compuertas de silicio (SGT o *Silicon Gate Technology*²⁷), técnica que se sigue aplicando hoy en día. Faggin también era responsable del desarrollo de la técnica llamada *Random Logic* que permitió la síntesis de descripciones complejas de lógica en hardware sencillo, como compuertas AND y OR. En el momento en el que se estaba desarrollando el proyecto MCS-4, el responsable del área de diseño de MOS de Intel era Leslie L. Vadász, cuya atención estaba enfocada

²¹Nota sobre *Electronic News*

²²Nota sobre Busicom

²³Nota sobre Ted Hoff

²⁴Nota sobre Stanley Mazor

²⁵Nota sobre Masatoshi Shima

²⁶Nota sobre Federico Faggin

²⁷Nota sobre SGT, o también llamada Self-aligned Gate

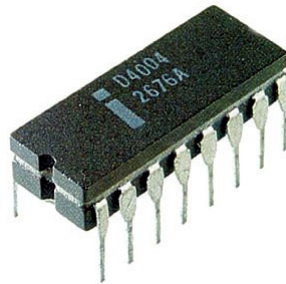


Figura 2.7: Encapsulado del Intel 4004; que es considerado el primer microprocesador comercial en ser introducido en el mercado.

en el mercado de memorias. Fue gracias a esto que le otorgó el liderazgo del proyecto MCS-4 a Faggin, quien fue finalmente el responsable de llevar el proyecto hasta la concepción final del 4004 gracias a la aplicación de las técnicas de diseño y fabricación antes mencionadas. Las primeras unidades del 4004 fueron entregadas a Busicom en Marzo de 1971.

2.2. Evolución hacia las arquitecturas modernas

El concepto moderno de las arquitecturas se basa en las máquinas que almacenan tanto el programa que se ejecuta, como los datos que se manjan en memoria de lectura-escritura (*Stored-Program digital computer*). Esto las diferencia de las primeras implementaciones de los años 40, tales como Colossus²⁸ y la ENIAC.

En esta sección se estudian las distintas clasificaciones de arquitecturas existentes hoy en día. Luego se detallarán diversas arquitecturas existentes y se las enmarcará dentro de dichas clasificaciones.

2.2.1. Clasificación de las arquitecturas según el acceso a la memoria de instrucciones y datos

Las instrucciones que ejecuta el procesador, al igual que los datos del programa en ejecución, deben residir en la memoria. De este hecho se desprende que existen -al menos- dos tipos de memorias: de instrucciones, y de datos. El acceso a estas dos memorias puede realizarse a través del mismo *bus*²⁹ o con *buses* separados. Los nombres que la historia les ha provisto a estos dos posibles enfoques, son *von Neumann* para las arquitecturas de *bus* único y *Harvard* para aquellas de *buses* separados.

²⁸Nota sobre Colossus

²⁹Bus: definir bus

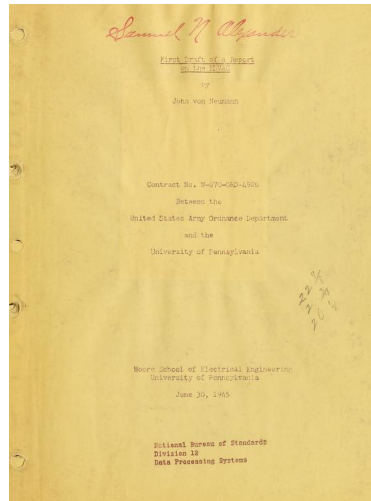


Figura 2.8: Carátula del escrito de von Neumann First Draft of a Report on the EDVAC

Arquitectura von Neumann

También conocida como modelo de von Neumann y arquitectura Princeton, este tipo de arquitectura debe su nombre a John von Neumann³⁰, quién en el año 1945 describió en el documento inconcluso, cuya carátula puede verse en la figura 2.8 y conocido como *First Draft of a Report on the EDVAC*[1], una arquitectura de computadora electrónica digital para ser implementada con válvulas de vacío. Ésta fue la primera publicación dónde se describe el diseño lógico de una computadora que utilizase el concepto de *stored-program* o programa almacenado. Este diseño, tal como fue llamado por el propio von Neumann, *very high speed automatic digital computing system* o sistema automático digital de cómputo de muy alta velocidad, estaba dividido en 6 componentes:

- CA: *central arithmetic* o aritmética central
- CC: *central control* o control central
- M: *memory* o memoria
- I: *input* o entrada
- O: *output* o salida
- R: *external memory* o memoria externa

La memoria almacena tanto números (datos) como órdenes (instrucciones). La aritmética central podía realizar sumas, restas, multiplicaciones, divisiones y raíces cuadradas.

³⁰Nota sobre el cambio

Otras operaciones matemáticas, como logaritmos y funciones trigonométricas se debían realizar utilizando tablas de búsqueda e interpolaciones, posiblemente bicuadráticas. Una decisión de diseño establecida en el documento estableció que multiplicaciones y divisiones podían realizarse mediante tablas logarítmicas, pero para poder mantener dichas tablas pequeñas debería utilizarse interpolaciones, lo cuál a su vez necesita de multiplicaciones, aunque de menor precisión.

Los números se representaban mediante notación binaria. Estimó que 27 dígitos binarios³¹ deberían ser suficientes, pero redondó a 30 dígitos, más uno de signo y otro para diferenciar los números de las órdenes, resultando en palabras de 32 dígitos binarios. La aritmética utilizada era complemento a dos para simplificar la operación de resta. Para la multiplicación y la división, propuso ubicar el punto binario luego del bit de signo, lo que implica que el dominio de los operandos y resultados está en el rango -1 a 1, y por lo tanto, los datos y resultados de los programas deberían ser escalados acordeamente.

En cuanto al diseño de los circuitos, estableció que deberían utilizarse válvulas de vacío, dejando de lado los relés, dada a la mayor velocidad provista por las primeras. Otras sugerencias involucraban mantener al sistema de cómputo lo más sencillo posible, evitando cualquier tipo de optimización de performance mediante la superposición de operaciones. Las operaciones aritméticas debían realizarse de a un dígito binario a la vez. Estimó el tiempo de la suma de dos dígitos binarios en un microsegundo, por lo tanto la multiplicación de dos números representados por 30 dígitos binarios debería realizarse en aproximadamente 30^{20} microsegundos; lo cuál era mucho más rápido que el tiempo de cualquier dispositivo de cálculo del momento. El diseño estaba basado en lo que von Neumann llamó “elemento E”, basándose en el modelo biológico de las neuronas, pero implementado de forma digital planteando que podrían ser fabricados mediante una o dos válvulas. En términos modernos, la forma más sencilla del “elemento E” es una compuerta *AND* de dos entradas, con una de sus entradas invertidas, llamada “entrada de inhibición”. Al agregar más entradas a este dispositivo, se establecía un nivel de umbral, el cual al ser superado por la suma de una determinada cantidad de entradas generaba una salida en tanto y en cuanto no se excitara la entrada de inhibición. También estableció que dichos elementos de múltiples entradas podían ser construidos utilizando combinaciones de la versión elemental, pero recomendaba implementarlos por completo utilizando así menos válvulas de vacío con el objetivo de lograr circuitos más sencillos y rápidos, principio que hoy en día se mantiene vigente.

Toda función lógica arbitraria, podía implementarse, entonces, a partir de dichos “elementos E”. Demostró en este trabajo, cómo implementar circuitos que implementaban las funciones aritméticas, así como elementos de memoria y circuitos de control, sin referirse en ningún momento al término de “lógica binaria”.

Los circuitos debían ser sincrónicos, obteniendo la señal de reloj a partir de un circuito oscilador implementado con válvulas y posiblemente controlado mediante cristales osciladores. Los diagramas lógicos incluían los tiempos de demora representados mediante una flecha. Estimó que la velocidad a la que se movía un pulso eléctrico en un cable

³¹En el documento se hace mención a dígitos binarios y no a *bits*, término que fue acuñado en 1948 por Claude Shannon

era de $300\text{mts}/\text{microsegundo}$, por lo tanto no debería generar problemas hasta obtener velocidades de reloj de 100MHz . También menciona pero no desarrolla la necesidad de contar con mecanismos de detección y corrección.

En cuanto a la memoria, que es donde entra la clasificación discutida en esta sección, von Neumann estableció que la memoria es uniforme, conteniendo tanto los números (datos) como las órdenes (instrucciones). Citando su trabajo:

“El dispositivo requiere una memoria considerable. A pesar de que pareciese que distintas partes de la memoria tiene que realizar funciones que difieren en su naturaleza y considerablemente en su propósito, resulta tentador tratar a toda la memoria con un sólo órgano, y hacer que sus partes sean tan intercambiables como sea posible para todas las funciones que deban realizar” [1, sección, 2.5]

“Las órdenes recibidas por CC vienen de M, es decir, el mismo lugar donde se almacenan los datos numéricos” [1, sección, 14.0].

Él concluyó que la memoria sería el subsistema más grande y abarcativo de la máquina. Basándose en diversos problemas matemáticos, incluyendo la resolución de ecuaciones diferenciales parciales y ordinarias, ordenamientos y experimentos probabilísticos, estimó que se necesitaba espacio para almacenar 8192 palabras de 32 dígitos binarios para los datos, y algunos cientos de palabras para almacenar las órdenes. Para la implementación de memoria, propuso dos tipos de memoria rápida, *delay line memory* e *iconoscope*. Con estas implementaciones planteó que la memoria debía ser direccionable por palabras y para sortear los problemas de demora asociados a la lectura de estas memorias, organizó las mismas en 256 conjuntos de 1024 dígitos binarios, o sea 32 palabras, logrando así direccionar los conjuntos con 8 dígitos binarios y las palabras con 5 utilizando, entonces, 13 dígitos binarios en total para el direccionamiento completo.

En su trabajo también estableció el formato de las órdenes, al cual llamó “código”. Los tipos de órdenes incluyeron las operaciones aritméticas básicas, así como el movimiento de palabras entre CA y M (análogas a las instrucciones *load* y *store* de hoy en día que veremos más adelante), una orden que elegía entre dos números basado en el signo del resultado de una operación previa (análoga a una instrucción *branch*), órdenes para controlar la entrada y salida de datos y para indicarle a la CC que debía tomar instrucciones desde otra sección de M (análoga a un *jump*). Determinó, asimismo, la cantidad de dígitos binarios necesarios para necesarios para los distintos tipos de órdenes, sugirió lo que llamó “órdenes inmediatas”, donde la siguiente palabra se trata del operando (lo cual también tiene una anagolía con ciertos conjuntos de instrucciones actuales) y dejó planteado si era deseable dejar dígitos sin especificar para propósitos futuros y mayor direccionamiento de memoria.

Arquitectura Harvard

Este tipo de arquitectura, ve sus orígenes en un desarrollo propuesto por el Dr. Howard Aiken³² en la universidad de Harvard en el año 1937 basado en el trabajo llamado

³²Nota sobre Aiken

*Analytical Engine*³³ de Charles Babage³⁴. El desarrollo fue llevado al cabo por IBM³⁵ y se trataba de una máquina de cálculo llamada *Automatic Sequence Controlled Calculator - Harvard Mark I* entregado a la universidad el 24 de agosto de 1944. A diferencia de la máquina propuesta (posteriormente) por von Neumann, ésta era una máquina electromecánica construida con llaves, relés, engranajes y embreagues. Constaba de 765000 componentes electromecánicos, cientos de millas de cableado, un peso de 4.5 Ton. y un consumo de potencia de 3.7 kW. Para el ingreso de datos, contaba con 60 conjuntos de 24 llaves rotativas que permitían ingresar números decimales. Podía almacenar hasta 72 números de hasta 23 dígitos decimales. En un segundo era capaz de realizar 3 sumas o restas. La multiplicación tardaba 6 segundos y una división 13.5 segundos. El cálculo de funciones logarítmicas o trigonométricas utilizaba más de un minuto. Las instrucciones eran ingresadas a través de cintas perforadas de 24 canales (en la figura 2.9 se puede ver una cinta con instrucciones de esta máquina). Otra cinta podía ser utilizada para el ingreso de datos, pero utilizando un formato distinto al de las instrucciones; estas no podían ser ejecutadas desde los registros de almacenamiento. Esta separación entre datos e instrucciones es lo que se conoce como arquitectura Harvard y que la diferencia de la arquitectura von Neumann. El formato de las instrucciones definía tres campos en ocho canales separados. Cada acumulador, cada conjunto de switches, y los registros asociados a la entrada y salida y las unidades aritméticas tenían asignado un índice único. Estos números eran representados en formato binario en la cinta perforada de control, siendo el primer campo el índice del resultado de la operación, el segundo el origen del dato para la operación y el tercero el código que representa la operación a realizar. Una particularidad de esta primera versión es que no tenía instrucción para el salto condicional en el programa. Por lo tanto, los programas complejos debían ser largos y los “loops” se implementaban uniendo físicamente el final de la cinta perforada del programa con el principio.

La Mark I fue sucedida por la Mark II en el año 1948, la Mark III en septiembre de 1949 y la Mark IV en el año 1952, todos estos proyectos fueron de Aiken. La Mark II fue una mejora sobre la Mark I, pero aún basada en componentes electromecánicos. La tercera versión utilizaba en su mayoría componentes electrónicos tales como válvulas de vacío y diodos cristalinos, pero utilizaba cilindros magnéticos para almacenamiento, y relés para la transferencia de datos entre ellos. La cuarta versión fue la primera en ser completamente electrónica, reemplazando los cilindros magnéticos de almacenamiento con memorias de núcleo magnético; el tipo de memoria de acceso aleatorio predominante entre 1955 y 1975.

³³Nota sobre Analytical Engine

³⁴Nota sobre el tipo

³⁵Nota sobre IBM

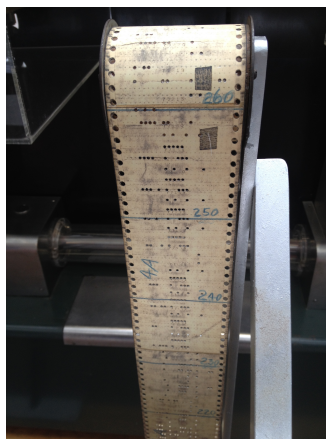


Figura 2.9: Cinta de papel perforado con instrucciones de la Mark I

Harvard vs. von Neumann y Harvard modificada

Como fue establecido al principio de esta sección, lo que diferencia estas dos arquitecturas es el acceso a los datos y a las instrucciones. La arquitectura von Neumann puede tratar a los datos como instrucciones y viceversa, dado que para dicha arquitectura, la memoria es exactamente la misma, con el mismo tipo de acceso y el mismo direccionamiento, contrariamente a lo que es una arquitectura puramente Harvard, cuyas memorias de datos e instrucciones se acceden por diferentes medios y cuyas direcciones pueden estar sobrepuestas (es decir que la misma dirección de memoria no representa lo mismo si se trata de datos o instrucciones) y por lo tanto, los datos y las instrucciones no pueden ser intercambiados en el tratamiento que se realiza dentro de la unidad de cómputo. Dadas estas condiciones, las principales desventajas son que, la arquitectura von Neumann no puede acceder simultáneamente a datos y a instrucciones y la máquina de Harvard no puede ni leer ni escribir en el espacio de memoria de las instrucciones, denegando por lo tanto cualquier posibilidad de código auto-optimizable o inteligente que reescriba sus instrucciones, así como también la carga de los programas a partir de medios de almacenamiento de datos.

Hoy en día, las máquinas puramente Harvard no son más que una especialidad. Típicamente, en las máquinas modernas, se implementan arquitecturas Harvard modificadas. Estas modificaciones pueden ser de diversos tipos:

2.2.2. Clasificación de las arquitecturas según su conjunto de instrucciones

RISC vs CISC

La ISA define el modelo lógico: Tres cosas, los registros, el modelo de la memoria y como se relacionan los registros (operaciones).

2.2.3. Taxonomía de Flynn

Instruction and data streams

IMPORTANTE:

El modelo formal del procesador queda definida por el ISA. El isa no es una tabla de instrucciones. Es el modelo lógico formal que describe la arquitectura.

Hablar de Alan Turing. Problemas Turing Computables -¿Investigar. Teorema de incompletitud de Gödel.

Referencias

- [1] John von Neumann, “First Draft of a Report on the EDVAC”, 1945.

Luciano César Natale