

Sistemas Operativos



AUTORES

Gunnar Wolf
Esteban Ruiz
Federico Bergero
Erwin Meza Vega

Sistemas Operativos

1a ed. - Iniciativa Latinoamericana de Libros de Texto Abiertos (LATIn), 2014. 248 pag.

Primera Edición: Marzo 2014

Iniciativa Latinoamericana de Libros de Texto Abiertos (LATIn)

<http://www.proyectolatin.org/>



Los textos de este libro se distribuyen bajo una licencia Reconocimiento-CompartirIgual 3.0 Unported (CC BY-SA 3.0) http://creativecommons.org/licenses/by-sa/3.0/deed.es_ES

Esta licencia permite:

Compartir: copiar y redistribuir el material en cualquier medio o formato.

Adaptar: remezclar, transformar y crear a partir del material para cualquier finalidad.

Siempre que se cumplan las siguientes condiciones:



Reconocimiento. Debe reconocer adecuadamente la autoría, proporcionar un enlace a la licencia e indicar si se han realizado cambios. Puede hacerlo de cualquier manera razonable, pero no de una manera que sugiera que tiene el apoyo del licenciador o lo recibe por el uso que hace.



CompartirIgual — Si remezcla, transforma o crea a partir del material, deberá difundir sus contribuciones bajo la misma licencia que el original.

Las figuras e ilustraciones que aparecen en el libro son de autoría de los respectivos autores. De aquellas figuras o ilustraciones que no son realizadas por los autores, se coloca la referencia respectiva.



Este texto forma parte de la Iniciativa Latinoamericana de Libros de Texto abiertos (LATIn), proyecto financiado por la Unión Europea en el marco de su Programa ALFA III EuropeAid.

El Proyecto LATIn está conformado por: Escuela Superior Politécnica del Litoral, Ecuador (ESPOL); Universidad Autónoma de Aguascalientes, México (UAA), Universidad Católica de San Pablo, Perú (UCSP); Universidade Presbiteriana Mackenzie, Brasil(UPM); Universidad de la República, Uruguay (UdelaR); Universidad Nacional de Rosario, Argentina(UR); Universidad Central de Venezuela, Venezuela (UCV), Universidad Austral de Chile, Chile (UACH), Universidad del Cauca, Colombia (UNICAUCA), Katholieke Universiteit Leuven, Bélgica (KUL), Universidad de Alcalá, España (UAH), Université Paul Sabatier, Francia (UPS).

Índice general

1	Presentación	9
	Presentación	9
1.1	Acerca del libro	9
2	Introducción	15
2.1	¿Qué es un sistema operativo?	15
2.1.1	¿Por qué estudiar los sistemas operativos?	15
2.2	Funciones y objetivos de los sistemas operativos	16
2.3	Evolución de los sistemas operativos	16
2.3.1	Proceso por lotes (<i>batch processing</i>)	16
2.3.2	Sistemas en lotes con dispositivos de carga (<i>spool</i>)	17
2.3.3	Sistemas multiprogramados	17
2.3.4	Sistemas de tiempo compartido	18
2.4	Y del lado de las computadoras personales	18
2.4.1	Primeros sistemas para entusiastas	19
2.4.2	La revolución de los 8 bits	19
2.4.3	La computadora para fines “serios”: La familia PC	20
2.4.4	El impacto del entorno gráfico (WIMP)	20
2.4.5	Convergencia de los dos grandes mercados	22
2.5	Organización de los sistemas operativos	22
2.6	Otros recursos	24
3	Relación con el hardware	25
3.1	Introducción	25
3.2	Unidad de Procesamiento	25
3.2.1	Jerarquía de almacenamiento	25
3.2.2	Interrupciones y excepciones	28
3.3	Terminales	29
3.4	Dispositivos de almacenamiento	30
3.5	Relojes y temporizadores	30
3.6	Canales y puentes	31
3.6.1	Contención	32
3.6.2	Acceso directo a memoria (DMA)	32

3.7	Interfaz del Sistema Operativo: llamadas al sistema	33
3.7.1	Llamadas al sistema, arquitecturas y APIs	34
3.8	Abstracciones comunes	35
3.8.1	Sistemas tipo Windows	35
3.8.2	Sistemas tipo Unix	35
3.9	Cuando dos cabezas piensan mejor que una	36
3.9.1	Multiprocesamiento	36
3.9.2	Cómputo distribuído	39
3.9.3	Amdahl y Gustafson: ¿qué esperar del paralelismo?	39
3.10	Otros recursos	42
4	Administración de procesos	45
4.1	Concepto y estados de un proceso	45
4.1.1	Estados de un proceso	45
4.1.2	Información asociada a un proceso	46
4.2	Procesos e hilos	46
4.2.1	Los hilos y el sistema operativo	47
4.2.2	Patrones de trabajo con hilos	47
4.3	Concurrencia	49
4.3.1	Introducción	49
4.3.2	Problema: el jardín ornamental	50
4.3.3	Mecanismos de sincronización	57
4.3.4	Problema productor-consumidor	64
4.3.5	Bloqueos mutuos e inanición	66
4.3.6	Problema lectores-escritores	66
4.3.7	La cena de los filósofos	68
4.3.8	Los fumadores compulsivos	70
4.3.9	Otros mecanismos	72
4.4	Bloqueos mutuos	76
4.4.1	Prevención de bloqueos	77
4.4.2	Evasión de bloqueos	80
4.4.3	Detección y recuperación de bloqueos	83
4.4.4	Algoritmo del aveSTRUZ	86
4.5	Otros recursos	88
5	Planificación de procesos	89
5.1	Tipos de planificación	89
5.1.1	Tipos de proceso	90
5.1.2	Midiendo la respuesta	91
5.2	Algoritmos de planificación	93
5.2.1	Objetivos de la planificación	94
5.2.2	Primero llegado, primero servido (<i>FCFS</i>)	95
5.2.3	Ronda (<i>Round Robin</i>)	95
5.2.4	El proceso más corto a continuación (<i>SPN</i>)	97
5.2.5	Ronda egoísta (<i>SRR</i>)	99
5.2.6	Retroalimentación multinivel (<i>FB</i>)	100
5.2.7	Lotería	102

5.2.8	Esquemas híbridos	102
5.2.9	Resumiendo	104
5.3	Planificación de hilos	106
5.3.1	Los hilos POSIX (<i>pthreads</i>)	107
5.4	Planificación de multiprocesadores	108
5.4.1	Afinidad a procesador	108
5.4.2	Balanceo de cargas	109
5.4.3	Colas de procesos: ¿Una o varias?	109
5.4.4	Procesadores con soporte a <i>hilos hardware (hyperthreading)</i>	110
5.5	Tiempo real	111
5.5.1	Tiempo real duro y suave	111
5.5.2	Sistema operativo interrumpible (<i>prevenible</i>)	112
5.5.3	Inversión de prioridades	113
5.6	Otros recursos	113
6	Administración de memoria	115
6.1	Funciones y operaciones del administrador de memoria	115
6.1.1	Espacio de direccionamiento	115
6.1.2	Hardware: de la unidad de manejo de memoria (MMU)	116
6.1.3	La memoria <i>caché</i>	117
6.1.4	El espacio en memoria de un proceso	118
6.1.5	Resolución de direcciones	119
6.2	Asignación de memoria contigua	120
6.2.1	Partición de la memoria	120
6.3	Segmentación	122
6.3.1	Permisos	123
6.3.2	Intercambio parcial	123
6.3.3	Ejemplificando	124
6.4	Paginación	125
6.4.1	Tamaño de la página	126
6.4.2	Almacenamiento de la tabla de páginas	126
6.4.3	Memoria compartida	128
6.5	Memoria virtual	130
6.5.1	Paginación sobre demanda	131
6.5.2	Rendimiento	131
6.5.3	Reemplazo de páginas	132
6.5.4	Asignación de marcos	136
6.5.5	Hiperpaginación	139
6.6	Consideraciones de seguridad	140
6.6.1	Desbordamientos de buffer (<i>buffer overflows</i>)	140
6.6.2	Ligado estático y dinámico de bibliotecas	145
6.7	Otros recursos	147
7	Organización de archivos	159
7.1	Introducción	159

7.2	Concepto de archivo	160
7.2.1	Operaciones con archivos	160
7.2.2	Tablas de archivos abiertos	162
7.2.3	Acceso concurrente: Bloqueo de archivos	162
7.2.4	Tipos de archivo	163
7.2.5	Estructura de los archivos y métodos de acceso	165
7.2.6	Transferencias orientadas a bloques	167
7.3	Organización de archivos	168
7.3.1	Evolución del concepto de <i>directorio</i>	168
7.3.2	Operaciones con directorios	172
7.3.3	<i>Montaje</i> de directorios	174
7.4	Sistemas de archivos remotos	176
7.4.1	Network File System (NFS)	176
7.4.2	Common Internet File System (CIFS)	178
7.4.3	Sistemas de archivos distribuidos: Andrew File System (AFS)	178
7.5	Otros recursos	179
8	Sistemas de archivos	181
8.1	Plasmando la estructura en el dispositivo	181
8.1.1	Conceptos para la organización	181
8.1.2	Diferentes sistemas de archivos	182
8.1.3	El volumen	183
8.1.4	El directorio y los i-nodos	184
8.1.5	Compresión y <i>desduplicación</i>	188
8.2	Esquemas de asignación de espacio	191
8.2.1	Asignación contigua	191
8.2.2	Asignación ligada	191
8.2.3	Asignación indexada	193
8.2.4	Las tablas en FAT	195
8.3	Fallos y recuperación	197
8.3.1	Datos y metadatos	198
8.3.2	Verificación de la integridad	199
8.3.3	Actualizaciones suaves (<i>soft updates</i>)	199
8.3.4	Sistemas de archivo con bitácora (<i>journaling file systems</i>)	200
8.3.5	Sistemas de archivos estructurados en bitácora (<i>log-structured file systems</i>)	201
8.4	Otros recursos	202
A	Software libre y licenciamiento	205
A.1	Software libre	205
A.1.1	<i>Free as in Freedom</i> : El proyecto GNU	205
A.1.2	El software libre antes de GNU	206
A.1.3	El <i>software propietario</i> como anomalía histórica	206
A.1.4	Esquemas libres de licenciamiento	207
A.2	Obras culturales libres	209
A.2.1	La familia de licencias <i>Creative Commons</i>	209
A.3	El licenciamiento empleado para la presente obra	211
A.4	Otros recursos	211

B	Virtualización	213
B.1	Introducción	213
B.2	Emulación	213
B.2.1	Emulando arquitecturas inexistentes	214
B.2.2	De lo abstracto a lo concreto	215
B.2.3	¿Emulación o simulación?	217
B.3	Virtualización asistida por hardware	217
B.3.1	El hipervisor	218
B.3.2	Virtualización asistida por hardware en x86	218
B.4	Paravirtualización	219
B.4.1	Paravirtualización y software libre	220
B.4.2	Paravirtualización de dispositivos	220
B.5	Contenedores, o <i>virtualización a nivel sistema operativo</i>	221
B.6	Otros recursos	223
C	El medio físico y el almacenamiento	225
C.1	El medio físico	225
C.1.1	Discos magnéticos rotativos	225
C.1.2	Almacenamiento en estado sólido	231
C.2	RAID: Más allá de los límites físicos	234
C.2.1	RAID nivel 0: División en <i>franjas</i>	235
C.2.2	RAID nivel 1: Espejo	236
C.2.3	Los niveles 2, 3 y 4 de RAID	236
C.2.4	RAID nivel 5: Paridad dividida por bloques	236
C.2.5	RAID nivel 6: Paridad por redundancia P+Q	238
C.2.6	Niveles combinados de RAID	238
C.3	Manejo avanzado de volúmenes	239
C.3.1	LVM: el Gestor de Volúmenes Lógicos	240
C.3.2	ZFS	241
C.4	Otros recursos	241

1 — Presentación

1.1 Acerca del libro

Este libro busca brindar a estudiantes y docentes de las carreras de *Ingeniería en Computación /, Informática, Ciencias de la Computación* y similares un material completo, general y autocontenido sobre la materia de *Sistemas Operativos*. No se asume conocimiento previo sobre la temática, aunque se utilizarán conceptos de estructuras de datos y algoritmos básicos.

Justificación

Actualmente existe vasta bibliografía sobre Sistemas Operativos, sin embargo la gran mayoría está escrita en inglés, y cuando están disponibles en castellano, su traducción deja mucho que desear, llevando a conceptos confusos y difíciles de comprender. La intención de los autores es que el presente texto provea un material redactado originalmente en castellano, revisado por docentes latinoamericanos utilizando la terminología más adecuada para los alumnos de la región y eliminando muchos de los errores de traducción.

Generalmente el material de cursos de Sistemas Operativos está compuesto por partes de distintos libros, artículos de investigación, recursos en línea, software, ejercitación, etc. Por ello, el alumno debe recurrir a distintas fuentes durante el curso. El presente libro pretende ser de utilidad tanto para alumnos como para docentes como una única publicación autocontenido. Cabe remarcar también que el material bibliográfico generalmente está protegido por derecho de autor, es costoso y en muchos casos de difícil acceso (sobre todo las publicaciones en inglés).

Los contenidos de la bibliografía clásica de Sistemas Operativos están basadas en re-ediciones y compendio de libros de hace varias décadas que incluyen temas obsoletos o desactualizados. Existen también desarrollos y tendencias nuevas en el área que aun no han sido integradas en la bibliografía clásica, y mucho menos a las traducciones. El presente libro pretende también revisar y actualizar los conceptos clásicos de sistemas operativos incluyendo material de publicación reciente.

Este libro se desarrolló dentro del marco del [Proyecto LATIn](#), enfocado a la creación de libros de texto con un esquema de licenciamiento libre, derivados de la creación y colaboración de grupos de trabajo multinacionales, para la región latinoamericana.

Público objetivo

Este libro está apuntado tanto a estudiantes de carreras de informática, computación e ingenierías como a los aficionados de la computadora interesados en conocer un poco más de lo que realmente ocurre dentro de un sistema de cómputo y el rol que cumple el sistema operativo.

Al finalizar el libro se espera que el lector haya adquirido conocimientos y habilidades como:

- Administrar, diseñar y desarrollar un sistema operativo.
- Conociendo el funcionamiento general de los sistemas operativos, poder sacar mejor provecho de la computadora
- Conocer y saber aprovechar no sólo los sistemas, sino las metodologías y principales formas de interacción del software libre

Se asume también que el lector está familiarizado con algún lenguaje de programación de alto nivel, y –al menos a nivel básico– con C. Aunque los ejemplos de código están dados en diversos lenguajes de programación (Bash, Perl, c, PascalFC, Python, Ruby, Ensamblador, entre otros), éstos son tan sencillos que pueden ser fácilmente escritos en el lenguaje de elección del lector sin mayor esfuerzo.

Resultará muy conveniente que tener acceso a una computadora con sistema operativo Linux (GNU) u otro Unix libre.

Estructura temática

El texto comprende los siguientes capítulos:

1. **Introducción** Para comenzar a hablar de sistemas operativos, es necesario en primer término enmarcar *qué es* un sistema operativo y cuáles son sus funciones principales. También es importante detallar algunos puntos que, contrario a la percepción común, *no pueden* considerarse parte de sus funciones.
Este tema se presenta apoyado de la evolución histórica del cómputo, haciendo énfasis en por qué este proceso evolutivo en particular desembocó en los sistemas operativos que existen hoy en día.
2. **Relación con el hardware** Partiendo de que una de las principales tareas del sistema operativo es presentar una *abstracción regular* del hardware a los procesos que se ejecuten, resulta importante presentar cómo éste está estructurado, y cómo el sistema operativo puede comunicarse con él.
Este capítulo aborda la *jerarquía de almacenamiento*, el mecanismo de *interrupciones y excepciones* y el papel que juegan para las *llamadas al sistema*, las características base de diversos tipos de dispositivo del sistema, el concepto de *canales* (o *buses*) de comunicación, el mecanismo de *acceso directo a memoria*, y una introducción a lo que puede ser visto como tema conductor a lo largo de todo el libro: La importancia y complejidad de la *concurrencia*, y su relación con el *parallelismo y multiprocесamiento*.
3. **Administración de procesos** La entidad principal con la que interactúa un sistema operativo (sea para brindarle servicios o para imponerle restricciones) es el proceso. Este capítulo inicia presentando los diferentes estados de los procesos y la relación entre éstas y sus *hermanos menores* (los *hilos*), y los principales modelos empleados para el multiprocесamiento.
Todos los sistemas operativos modernos tienen que enfrentar a la *concurrencia*: La incertidumbre del ordenamiento en el tiempo entre eventos relativos a los diferentes procesos e hilos. La parte medular de este capítulo presenta a las *primitivas de sincronización*: Mutexes, semáforos y monitores. Para ilustrarlas, se emplean los patrones y problemas clásicos que se han seguido a lo largo de su desarrollo histórico. Pero las primitivas pueden sólamente utilizarse entre procesos que *cooperan deliberadamente* entre sí. Un sistema operativo debe implementar protección y separación incluso entre procesos que *compiten* o que sencillamente no saben el uno acerca del otro. Por ello, la última sección de este capítulo aborda los diferentes mecanismos que existen para evitar las situaciones de *bloqueo mutuo*.
4. **Planificación de procesos** Para que varios procesos coexistan en un sistema de cómputo, el primer recurso que el sistema operativo debe *multiplexar* o repartir entre todos ellos es el tiempo de cómputo: El uso del procesador. Este capítulo presenta los

diferentes niveles de *planificador* que forman parte de un sistema operativo, y analiza al *planificador a corto plazo* (también conocido como *despachador*). Se presentan los principales algoritmos, y se ilustra cómo los sistemas operativos modernos van empleando técnicas mixtas de varios de ellos.

Por último, se abordan tres temas brevemente: Los diferentes modelos de planificación de hilos y su relación con los procesos, las particularidades de la planificación en un entorno con multiprocesadores reales, y las necesidades de planificación de *tiempo real*.

- 5. Administración de memoria** Los *programas* sólo se vuelven procesos cuando se les asigna memoria y tiempo de cómputo: Cuando dejan de ser el resultado de una compilación guardada estáticamente para convertirse en una entidad dinámica. Este capítulo presenta en primer lugar la visión *desde dentro* de la memoria por parte de cada uno de los procesos: El espacio de direccionamiento y el acomodo clásico de las regiones de un proceso en la memoria que le es asignada.

Para que los distintos procesos compartan la memoria del sistema, a lo largo de la historia se han presentado diferentes esquemas. Se presentan someramente los esquemas de partición contigua fija y variable, para profundizar posteriormente en los que ofrecen mayor flexibilidad al sistema operativo y se mantienen en uso al día de hoy: La segmentación y la paginación. De esta última, se continúa para presentar la abstracción que ha liberado a los sistemas operativos para *sobrecomprometer* la memoria de forma eficiente y prácticamente transparente: La memoria virtual.

Al manejar la memoria de un proceso surgen puntos importantes a tomar en cuenta en lo relativo a la seguridad en cómputo; la parte final de este capítulo presenta la vulnerabilidad conocida como *desbordamiento de pila* (*buffer overflow*), y algunas estrategias de mitigación que se han implementado con el paso de los años para mitigar su peligrosidad.

- 6. Organización de archivos** De cara al usuario, probablemente la principal abstracción implementada por el sistema operativo es la organización de la información sobre un medio persistente. Hoy en día, la norma es que esta organización se realice en *archivos* estructurados sobre una estructura jerárquica llamada *directorio*. Este capítulo se centra en explicar esta abstracción, sin entrar aún en detalles respecto a cómo se llega a un respaldo físico de la misma.

Estos conceptos parecen tan pervasivos y universales que podría pensarse que no requieren mayor análisis. Sin embargo, resulta importante abordar las diferencias semánticas derivadas del desarrollo histórico de distintos sistemas. En este capítulo se presentan varios conceptos cuya implementación en un medio que asegure la persistencia se describirá en el siguiente capítulo.

Por último, en este capítulo se incluye un breve repaso de distintos tipos de sistemas de archivos en red, enfatizando nuevamente en los cambios semánticos derivados de la distinta historia de cada implementación.

- 7. Sistemas de archivos** Este capítulo presenta la contraparte obligada del anterior: ¿Cómo se estructuran los dispositivos de almacenamiento a largo plazo, a los cuales nos referimos genéricamente como *discos*? ¿Cómo se van plasmando las estructuras mediante las cuales el usuario organiza la información en bloques dentro de un dispositivo? ¿Qué problemas pueden derivar del uso de estos *sistemas de archivos*, y qué métodos para evitarlos o resolverlos se han implementado?

La parte central de este capítulo se centra en un sistema de archivos bastante viejo y simple, pero aún en muy amplio uso en el cómputo moderno: La familia *FAT*.

Los siguientes temas resultan muy importantes para la comprensión y para el desarrollo futuro de la materia, pero dado que son *empleados* por el sistema operativo (y no necesariamente son *parte integral* del mismo), se presentan como apéndices

A. Software libre y licenciamiento Estudiar sistemas operativos cruza necesariamente la temática del *software libre*. Uno de los principios fundamentales del desarrollo histórico es la *libertad de aprender*, esto es, todo software que se diga *libre* debe permitir a sus usuarios comprender sus estructuras básicas, la relación entre ellas, y la lógica general de su programación.

Hoy en día existe una gran cantidad de sistemas operativos libres, tanto de propósito general como enfocados a un nicho. El *movimiento ideológico* del software libre, contrario a cualquier pronóstico que pudiera haberse hecho al iniciarse en 1984, claramente ha cambiado el desarrollo del cómputo. Todos los sistemas operativos que pueden ser estudiados de *primera mano*, constatando la implementación de sus principios son necesariamente (aunque con una definición ligeramente laxa) software libre.

Hacia el año 2000 se fue haciendo claro que estas ideas no pueden aplicarse únicamente al software. Poco a poco fue definiéndose una noción mucho más amplia, la de los *bienes culturales libres*. El presente libro busca ser una contribución a esta última categoría.

El primer apéndice aborda brevemente estos temas, así como los principales modelos de licenciamiento libre utilizados.

B. Virtualización La virtualización es una herramienta muy útil, y cada vez más al alcance de todos, para el aprendizaje de los sistemas operativos. Hay una gran cantidad de recursos para comprender desde los primeros momentos del arranque de la computadora. Empleando imágenes de máquinas virtuales, pueden comprenderse y desmenuzarse los distintos elementos del sistema operativo, e incluso observar el resultado de realizar modificaciones sobre un sistema operativo real. Es, por tanto, una herramienta muy importante para acompañar al aprendizaje de esta materia. La virtualización es también una tecnología que permea cada vez más aspectos del uso profesional del cómputo, y comprenderlo ayudará al lector a elegir las herramientas específicas a emplear.

Pero hablar de *la virtualización* como un todo ignoraría aspectos fundamentales de la riqueza que presenta este campo. Al igual que con los conceptos presentados a lo largo del libro, la virtualización es presentada a partir de su perspectiva histórica, y detallando hacia las distintas modalidades que se han desarrollado al paso del tiempo.

C. El medio físico y el almacenamiento En el capítulo 8 se presenta cómo se *concretiza* la abstracción de archivos y directorios para plasmarlo en un gran arreglo lineal de datos, en una entidad aún abstracta a la cual se sigue haciendo referencia con el nombre genérico de *disco*. Este apéndice se ocupa de los detalles físicos del acomodo de la información en su medio.

Pero un *disco* va mucho más allá de un dispositivo que simplemente vuelca dicho arreglo a un medio persistente. En primer término, los *discos magnéticos rotativos* (el medio dominante de almacenamiento) presentan peculiaridades que los sistemas operativos tuvieron que saber resolver. El desarrollo de la tecnología, sin embargo, fue *arrebatiando* estas áreas del ámbito del sistema operativo, entregándolas a la optimización realizada dentro del *hardware controlador*.

Por otro lado, la tecnología de *almacenamiento en estado sólido* ha llegado a niveles de madurez que en determinados mercados ya la colocan claramente por encima de los discos magnéticos. Esto implica cambios importantes para el modo en que el sistema operativo debe estructurar y modificar la información.

Por último, un *volumen* ya no necesariamente se refiere a un único medio físico. Este apéndice aborda tanto a *RAID*, el primer mecanismo que se popularizó para *agregar* varias unidades para mejorar tanto la capacidad máxima y la confiabilidad de un

volumen, como al manejo avanzado de volúmenes, en que el sistema operativo incorpora la lógica de RAID con la del manejo de sistemas de archivos para lograr mucho mayor flexibilidad.

Licenciamiento

Este libro fue desarrollado como parte del [Proyecto LATIn](#), que busca la creación de libros de texto *libres* para nivel universitario, y enfocado a Latinoamérica.

Cualquier porción de este libro puede ser reproducido y utilizado para todo fin, bajo los términos de la licencia [Creative Commons-Atribución-CompartirIgual](#) (CC-BY-SA) versión 4.0.

Este modelo de licenciamiento se presenta y explica en la sección A.2.1.

2—Introducción

2.1 ¿Qué es un sistema operativo?

El *sistema operativo* es el principal programa que se ejecuta en toda computadora de propósito general.

Hay sistemas operativos de todo tipo, desde muy simples hasta terriblemente complejos, y entre más casos de uso hay para el cómputo en la vida diaria, más variedad habrá en ellos.

No nos referiremos al sistema operativo como lo ve el usuario final, o como lo vende la mercadotecnia — El ambiente gráfico, los programas que se ejecutan en éste, los lenguajes de programación en que están desarrollados y en que más fácilmente se puede desarrollar para ellos, e incluso el conjunto básico de funciones que las bibliotecas base ofrecen son principalmente *clientes* del sistema operativo — Se ejecutan sobre él, y ofrecen su implementación a sus usuarios (incluídos, claro, los desarrolladores). La diferencia en el uso son sólo —y si mucho— *consecuencias* del diseño de un sistema operativo. Más aún, con el mismo sistema operativo —como pueden constatarlo comparando dos distribuciones de Linux, o incluso la forma de trabajo de dos usuarios en la misma computadora— es posible tener *entornos operativos* completamente disímiles.

2.1.1 ¿Por qué estudiar los sistemas operativos?

La importancia de estudiar este tema radica no sólo en comprender los mecanismos que emplean los sistemas operativos para cumplir sus tareas sino en entender estos mecanismos para evitar los errores más comunes al programar, que pueden resultar desde un rendimiento deficiente hasta pérdida de información.

Como desarrolladores, comprender el funcionamiento básico de los sistemas operativos y las principales alternativas que nos ofrecen en muchos de sus puntos, o saber diseñar algoritmos y procesos que se ajusten mejor al sistema operativo en que vayamos a ejecutarlo, puede resultar en una diferencia cualitativa decisiva en nuestros productos.

Como administradores de sistemas, muchas veces podemos enfrentarnos a situaciones de bajo rendimiento, de conflictos entre aplicaciones, demoras en la ejecución, y comprender lo que ocurre *tras bambalinas* resulta fundamental para realizar nuestro trabajo. Los sistemas de archivos resultan un área de especial interés para administradores de sistemas: ¿Cómo comparar las virtudes y desventajas de tantos sistemas existentes? ¿Por qué puede resultarnos conveniente mezclarlos en el mismo servidor? ¿Cómo evitar la corrupción o pérdida de información? Lo que es más, ¿cómo recuperar información de un disco dañado?

En el área de la seguridad en cómputo, la relación resulta obvia: si nos interesa localizar vulnerabilidades que nos permitan elevar nuestro nivel de privilegios, ¿cómo podríamos hacerlo sin comprender cómo se engranan los diversos componentes de un sistema? La cantidad de tareas que debe cubrir un sistema operativo es tremenda, y veremos ejemplos de sitios donde un atacante puede enfocar sus energías. Del mismo modo, para quien busca *defender* un sistema (o una red), resulta fundamental comprender cuáles son los

vectores de ataque más comunes y –nuevamente– la relación entre los componentes involucrados para poder remediar o, mejor, prevenir dichos ataques.

Y claro está, podemos ver al mundo en general, fuera del entorno del cómputo, como una serie de modelos interactuantes. Muchos de los métodos y algoritmos que aquí veremos pueden emplearse fuera del entorno del cómputo, y una vez que comprendamos los problemas de concurrencia, de competencia por recursos, o de protección y separación que han sido resueltos en el campo de los sistemas operativos, podemos extrapolar estas soluciones a otros campos.

El camino por delante es largo, y puede resultar interesante y divertido.

2.2 Funciones y objetivos de los sistemas operativos

El sistema operativo es el único programa que interactúa directamente con el hardware de la computadora. Sus funciones primarias son:

Abstracción Los programas no deben tener que preocuparse de los detalles de acceso a hardware, o de la configuración particular de una computadora. El sistema operativo se encarga de proporcionar una serie de abstracciones para que los programadores puedan enfocarse en resolver las necesidades particulares de sus usuarios. Un ejemplo de tales abstracciones es que la información está organizada en *archivos* y *directorios* (en uno o muchos *dispositivos de almacenamiento*).

Administración de recursos Una sistema de cómputo puede tener a su disposición una gran cantidad de *recursos* (memoria, espacio de almacenamiento, tiempo de procesamiento, etc.), y los diferentes *procesos* que se ejecuten en él *compiten* por ellos. Al gestionar toda la asignación de recursos, el sistema operativo puede implementar políticas que los asigne de forma efectiva y acorde a las necesidades establecidas para dicho sistema.

Aislamiento En un sistema multiusuario y multitarea cada proceso y cada usuario no tendrá que preocuparse por otros que estén usando el mismo sistema — Idealmente, su *experiencia* será la misma que si el sistema estuviera exclusivamente dedicado a su atención (aunque fuera un sistema menos poderoso).

Para implementar correctamente las funciones de aislamiento hace falta que el sistema operativo utilice hardware específico para dicha protección.

2.3 Evolución de los sistemas operativos

No se puede comenzar a abordar el tema de los sistemas operativos sin revisar brevemente su desarrollo histórico. Esto no sólo permitirá comprender por qué fueron apareciendo determinadas características y patrones de diseño que se siguen empleando décadas más tarde, sino (como resulta particularmente bien ejemplificado en el discurso de recepción del premio Turing de Fernando Corbató en 1990, *On building systems that will fail*), adecuar un sistema existente a un entorno cambiante, por mejor diseñado que éste estuviera, lleva casi inevitablemente a abrir espacios de comportamiento no previsto — El espacio más propicio para que florezcan los fallos. Conocer los factores que motivaron a los distintos desarrollos puede ayudar a prever y prevenir problemas.

2.3.1 Proceso por lotes (*batch processing*)

Los antecedentes a lo que hoy se conoce como sistema operativo se pueden encontrarlos en la automatización inicial del procesamiento de diferentes programas, surgida en los primeros centros de cómputo: cuando en los '50 aparecieron los dispositivos perforadores/lectores de tarjetas de papel, el tiempo que una computadora estaba improductiva esperando a que estuviera lista una *tarea* (como se designaba a una ejecución de cada

determinado programa) para poder ejecutarla disminuyó fuertemente ya que los programadores entregaban su lote de tarjetas perforadas (en inglés, batches) a los operadores, quienes las alimentaban a los dispositivos lectores, que lo cargaban en memoria en un tiempo razonable, iniciaban y monitoreaban la ejecución, y producían los resultados.

En esta primer época en que las computadoras se especializaban en tareas de cálculo intensivo y los dispositivos que interactuaban con medios externos eran prácticamente desconocidos, el rol del sistema *monitor* o *de control* era básicamente asistir al operador en la carga de los programas y las bibliotecas requeridas, la notificación de resultados y la contabilidad de recursos empleados para su cobro.

Los sistemas monitores se fueron sofisticando al implementar protecciones que evitaban la corrupción de *otros trabajos* (por ejemplo, lanzar erróneamente la instrucción *leer siguiente tarjeta* causaría que el siguiente trabajo encolado perdiera sus primeros caracteres, corrompiéndolo e impidiendo su ejecución), o que entraran en un ciclo infinito, estableciendo *alarmas (timers)* que interrumpirían la ejecución de un proceso si éste duraba más allá del tiempo estipulado. Estos monitores implicaban la modificación del hardware para contemplar dichas características de seguridad — Y ahí se puede hablar ya de la característica básica de gestión de recursos que identifica a los sistemas operativos.

Cabe añadir que el tiempo de carga y puesta a punto de una tarea seguía representando una parte importante del tiempo que la computadora dedicaba al procesamiento: un lector de cintas rápido procesaba del orden de cientos de caracteres por minuto, y a pesar de la lentitud relativa de las computadoras de los '50 ante los estándares de hoy (se medirían por miles de instrucciones por segundo, KHz, en vez de miles de millones como se hace hoy, GHz), esperar cinco o diez minutos con el sistema completamente detenido por la carga de un programa moderadamente extenso resulta a todas luces un desperdicio.

2.3.2 Sistemas en lotes con dispositivos de carga (*spool*)

Una mejora natural a este último punto fue la invención del *spool*: Un mecanismo de entrada/salida que permitía que una computadora de propósito específico, mucho más económica y limitada, leyera las tarjetas y las fuera convirtiendo a cinta magnética, un medio mucho más rápido, teniéndola lista para que la computadora central la cargara cuando terminara con el trabajo anterior. Del mismo modo, la computadora central guardaba sus resultados en cinta para que equipos especializados la leyieran e imprimieran para el usuario solicitante.

La palabra *spool* (*bobina*) se tomó como acrónimo inverso hacia *Simultaneous Peripheral Operations On-Line, operación simultánea de periféricos en línea*.

2.3.3 Sistemas multiprogramados

A lo largo de su ejecución, un programa normalmente pasa por etapas con muy distintas características: durante un ciclo fuertemente dedicado al cálculo numérico, el sistema opera *limitado por el CPU* (*CPU-bound*), mientras que al leer o escribir resultados a medios externos (incluso a través de *spools*) el límite es impuesto por los dispositivos, esto es, opera *limitado por entrada-salida* (*I-O bound*). La *programación multitareas* o los *sistemas multiprogramados* buscaban maximizar el tiempo de uso efectivo del procesador ejecutando varios procesos al mismo tiempo.

El hardware requerido cambió fuertemente. Si bien se esperaba que cada usuario fuera responsable con el uso de recursos, se hizo necesario que apareciera la infraestructura de protección de recursos: un proceso no debe sobreescribir el espacio de memoria de otro (ni el código ni los datos), mucho menos el espacio del monitor. Esta protección se encuentra en la *Unidad de Manejo de Memoria* (MMU), presente en todas las computadoras de uso genérico desde los '90.

Ciertos dispositivos requieren bloqueo para ofrecer acceso exclusivo/único — Cintas e impresoras, por ejemplo, son de acceso estrictamente secuencial, y si dos usuarios intentaran usarlas al mismo tiempo, el resultado para ambos se corrompería. Para estos dispositivos, el sistema debe implementar otros *spools* y mecanismos de bloqueo.

2.3.4 Sistemas de tiempo compartido

El modo de interactuar con las computadoras se modificó drásticamente durante los '60, al extenderse la multitarea para convertirse en sistemas *interactivos* y *multusuarios*, en buena medida diferenciados de los anteriores por la aparición de las *terminales* (primero teletipos seriales, posteriormente equipos con una pantalla completa como se conocen hasta hoy).

En primer término, la tarea de programación y depuración del código se simplificó fuertemente al poder el programador hacer directamente cambios y someter el programa a la ejecución inmediata. En segundo término, la computadora *nunca más estaría simplemente esperando a que esté listo un programa*: Mientras un programador editaba o compilaba su programa, la computadora seguía calculando lo que otros procesos requerían.

Un cambio fundamental entre el modelo de *multiprogramación* y de *tiempo compartido* es el tipo de control sobre la multitarea: (se verá en detalle en el capítulo 4 (*Administración de procesos*))

Multitarea cooperativa o no apropiativa (*Cooperative multitasking*) La implementaron los sistemas multiprogramados: Cada proceso tenía control del CPU hasta que éste hacía una llamada al sistema (o indicara su *disposición a cooperar* por medio de la llamada `yield`: *ceder el paso*).

Un cálculo largo no era interrumpido por el sistema operativo, en consecuencia un error de programador podía congelar la computadora completa.

Multitarea preventiva o apropiativa (*Premptive multitasking*) En los sistemas de tiempo compartido, el reloj del sistema interrumpe periódicamente a los diversos procesos, transfiriendo *forzosamente* el control nuevamente al sistema operativo. El sistema operativo puede entonces elegir otro proceso para continuar la ejecución.

Además, fueron naciendo de forma natural y paulatina las abstracciones que se conocen hoy en día, como los conceptos de *archivos* y *directorios*, y el código necesario para emplearlos iba siendo enviado a las *bibliotecas de sistema* y, cada vez más (por su centralidad) hacia el núcleo mismo del –ahora sí– sistema operativo.

Un cambio importante entre los sistemas multiprogramados y de tiempo compartido es que la velocidad del cambio entre una tarea y otra es mucho más rápido: si bien en un sistema multiprogramado un *cambio de contexto* podía producirse sólo cuando la tarea cambiaba de un modo de ejecución a otro, en un sistema interactivo, para dar la *ilusión* de uso exclusivo de la computadora, el hardware emitía periódicamente al sistema operativo *interrupciones* (señales) que le indicaban que cambie el *proceso activo* (como ahora se le denomina a una instancia de un programa en ejecución).

Diferentes tipos de proceso pueden tener distinto nivel de importancia — Ya sea porque son más relevantes para el funcionamiento de la computadora misma (procesos de sistema), porque tienen mayor carga de interactividad (por la experiencia del usuario) o por diversas categorías de usuarios (sistemas con *contabilidad* por tipo de atención). Esto requiere la implementación de diversas *prioridades* para cada uno de estos.

2.4 Y del lado de las computadoras personales

Si bien la discusión hasta este momento asume una computadora central con operadores dedicados y múltiples usuarios, en la década de los '70 comenzaron a aparecer las *computadoras personales*, sistemas en un inicio verdaderamente reducidos en prestaciones y

a un nivel de precios que los ponían al alcance, primero, de los aficionados entusiastas y, posteriormente, de cualquiera.

2.4.1 Primeros sistemas para entusiastas



Figura 2.1: La microcomputadora Altair 8800, primer computadora personal con distribución masiva, a la venta a partir de 1975. (Imagen de la Wikipedia: *Altair 8800*)

Las primeras computadoras personales eran distribuidas sin sistemas operativos o lenguajes de programación; la interfaz primaria para programarlas era a través de llaves (*switches*), y para recibir sus resultados, se utilizaban bancos de LEDs. Claro está, esto requería conocimientos especializados, y las computadoras personales eran aún vistas sólo como juguetes caros.

2.4.2 La revolución de los 8 bits

La verdadera revolución apareció cuando, poco tiempo más tarde, comenzaron a venderse computadoras personales con salida de video (típicamente a través de una televisión) y entrada a través de un teclado. Estas computadoras popularizaron el lenguaje de programación BASIC, diseñado para usuarios novatos en los '60, y para permitir a los usuarios gestionar sus recursos (unidades de cinta, pantalla posicional, unidades de disco, impresoras, modem, etc.) llevaban un software mínimo de sistema — Nuevamente, un proto-sistema operativo.



Figura 2.2: La Commodore Pet 2001, en el mercado desde 1977, una de las primeras con intérprete de BASIC. (Imagen de la Wikipedia: *Commodore PET*)

2.4.3 La computadora para fines “serios”: La familia PC

Al aparecer las computadoras personales “serias”, orientadas a la oficina más que al hobby, a principios de los ‘80 (particularmente representadas por la IBM PC, 1981), sus sistemas operativos se comenzaron a diferenciar de los equipos previos al separar el *entorno de desarrollo* en algún lenguaje de programación del *entorno de ejecución*. El rol principal del sistema operativo ante el usuario era administrar los archivos de las diversas aplicaciones a través de una sencilla interfaz de línea de comando, y lanzar las aplicaciones que el usuario seleccionaba.

La PC de IBM fue la primer arquitectura de computadoras personales en desarrollar una amplia familia de *clones*, computadoras compatibles diseñadas para trabajar con el mismo sistema operativo, y que eventualmente capturaron casi el 100% del mercado. Prácticamente todas las computadoras de escritorio y portátiles en el mercado hoy derivan de la arquitectura de la IBM PC.



Figura 2.3: La computadora IBM PC modelo 5150 (1981), iniciadora de la arquitectura predominante en uso hasta el día de hoy. (Imagen de la Wikipedia: *IBM Personal Computer*)

Ante las aplicaciones, el sistema operativo (PC-DOS, en las versiones distribuidas directamente por IBM, o el que se popularizó más, MS-DOS, en los *clones*) ofrecía la ya conocida serie de interfaces y abstracciones para administrar los archivos y la entrada/-salida a través de sus puertos. Cabe destacar que, particularmente en sus primeros años, muchos programas se ejecutaban directamente sobre el hardware, arrancando desde el BIOS y sin emplear el sistema operativo.

2.4.4 El impacto del entorno gráfico (WIMP)

Hacia mediados de los ‘80 comenzaron a aparecer computadoras con interfaces gráficas basadas en el paradigma WIMP (*Windows, Icons, Menus, Pointer*; Ventanas, Iconos, Menúes, Apuntador), que permitían la interacción con varios programas al mismo tiempo. Esto *no necesariamente* significa que sean sistemas multitarea — Por ejemplo, la primer interfaz de MacOS permitía visualizar varias ventanas abiertas simultáneamente, pero sólo el proceso activo se ejecutaba.

Esto comenzó, sin embargo, a plantear inevitablemente las necesidades de concurrencia a los programadores. Los programas ya no tenían acceso directo a la pantalla para manipular a su antojo, sino que a una abstracción (la *ventana*) que podía variar sus medidas, y que requería que toda la salida fuera estrictamente a través de llamadas a bibliotecas de primitivas gráficas que comenzaron a verse como parte integral del sistema



Figura 2.4: Apple Macintosh (1984), popularizó la interfaz usuario gráfica (GUI). (Imagen de la Wikipedia: *Macintosh*)

operativo.

Además, los problemas de protección y separación entre procesos concurrentes comenzaron a hacerse evidentes: los programadores tenían ahora que programar con la conciencia de que compartirían recursos — con el limitante (que no tenían en las máquinas *profesionales*) de no contar con hardware especializado para esta protección. Los procesadores en uso comercial en los ‘80 no manejaban *anillos* o *niveles de ejecución* ni *unidad de administración de memoria* (MMU), por lo que un programa fallado o dañino podía corromper la operación completa del equipo. Y si bien los entornos que más éxito tuvieron (Apple MacOS y Microsoft Windows) no implementaban multitarea real, sí hubo desde el principio sistemas como la Commodore Amiga o la Atari ST que hacían un multitasking *preventivo* verdadero.



Figura 2.5: Commodore Amiga 500 (1987), la computadora más popular de la familia *Amiga*, con amplias capacidades multimedia y multitarea preventiva; una verdadera maravilla para su momento. (Imagen de la Wikipedia: *Amiga*)

Naturalmente, ante el uso común de un entorno de ventanas, los programas que se

ejecutaban sin requerir de la carga del sistema operativo cayeron lentamente en el olvido.

2.4.5 Convergencia de los dos grandes mercados

Conforme fueron apareciendo los CPU con características suficientes en el mercado para ofrecer la protección y aislamiento necesario (particularmente, Intel 80386 y Motorola 68030), la brecha de funcionalidad entre las computadoras personales y las *estaciones de trabajo* y *mainframes* se fue cerrando.

Hacia principios de los 1990, la mayor parte de las computadoras de arquitecturas *alternativas* fueron cediendo a las presiones del mercado, y hacia mediados de la década sólo quedaban dos arquitecturas principales: la derivada de IBM y la derivada de la Apple Macintosh.

Los sistemas operativos primarios para ambas plataformas fueron respondiendo a las nuevas características del hardware: en las IBM, la presencia de Microsoft Windows (originalmente un *entorno operativo* desde su primer edición en 1985, evolucionando hacia un sistema operativo completo ejecutando sobre una base de MS-DOS en 1995) se fue haciendo prevalente hasta ser la norma. Windows pasó de ser un sistema meramente de aplicaciones propias y que operaba únicamente por reemplazo de aplicación activa a ser un sistema de multitarea cooperativa, y finalmente un sistema que requería protección en hardware (80386) e implementaba multitarea preventiva.

A partir del 2003, el núcleo de Windows en más amplio uso fue reemplazado por un desarrollo hecho de inicio como un sistema operativo completo y ya no como una aplicación dependiente de MS-DOS: el núcleo de Nueva Tecnología (Windows NT), que, sin romper compatibilidad con los *APIs* históricos de Windows, ofreció mucho mayor estabilidad.

Por el lado de Apple, la evolución fue muy en paralelo: ante un sistema ya agotado y obsoleto, el MacOS 9, en 2001 anunció una nueva versión de su sistema operativo que fue en realidad un relanzamiento completo: MacOS X es un sistema basado en un núcleo Unix BSD, sobre el *microkernel* Mach.

Y otro importante jugador que entró en escena durante los '90 fue el software libre, por medio de varias implementaciones distintas de sistemas tipo Unix — principalmente, Linux y los *BSD (FreeBSD, NetBSD, OpenBSD). Estos sistemas implementaron, colaborativamente y a escala mundial, software compatibles con las PC y con el que se ejecutaba en las estaciones de trabajo a gran escala, con alta confiabilidad, y cerrando por fin la divergencia del árbol del desarrollo de la computación en *fierros grandes* y *fierros chicos*.

Al día de hoy, la arquitectura derivada de Intel (y la PC) es el claro ganador de este proceso de 35 años, habiendo conquistado casi la totalidad de los casos de uso, incluso las máquinas Apple. Hoy en día, la arquitectura Intel ejecuta desde subportátiles hasta supercomputadoras y centros de datos; el sistema operativo específico varía según el uso, yendo mayoritariamente hacia Windows, con los diferentes Unixes concentrados en los equipos servidores.

En el frente de los dispositivos *embebidos* (las computadoras más pequeñas, desde microcontroladores hasta teléfonos y tabletas), la norma es la arquitectura ARM, también bajo versiones específicas de sistemas operativos Unix y Windows (en ese orden).

2.5 Organización de los sistemas operativos

Para comenzar el estudio de los sistemas operativos, la complejidad del tema requiere que se haga de una forma modular. En este texto no se busca enseñar cómo se usa un determinado sistema operativo, ni siquiera comparar el uso de uno con otro (fuera de hacerlo con fines de explicar diferentes implementaciones).

Al nivel que se estudiará, un sistema operativo es más bien un gran programa, que ejecuta otros programas y les provee un conjunto de interfaces para que puedan aprovechar los recursos de cómputo. Hay dos formas primarias de organización *interna* del sistema operativo: los sistemas monolíticos y los sistemas microkernel. Y si bien no se puede marcar una línea clara a rajatabla que indique en qué clasificación cae cada sistema, no es difícil encontrar líneas bases.

Monolíticos La mayor parte de los sistemas operativos históricamente han sido *monolíticos* — Esto significa que hay un sólo *proceso privilegiado* (justamente el sistema operativo) que opera en modo supervisor, y dentro del cual se encuentran todas las rutinas para las diversas tareas que realiza el sistema operativo.

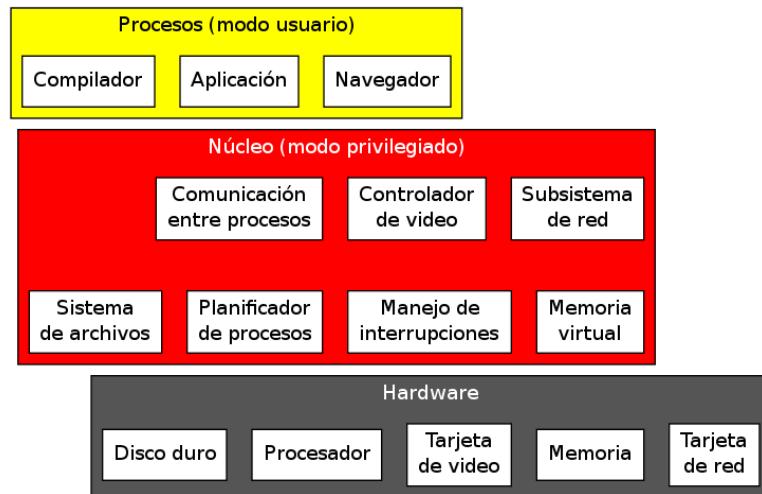


Figura 2.6: Esquematización de los componentes en un sistema monolítico

Microkernel El núcleo del sistema operativo se mantiene en el mínimo posible de funcionalidad, descargando en *procesos especiales sin privilegios* las tareas que implementan el acceso a dispositivos y las diversas políticas de uso del sistema.

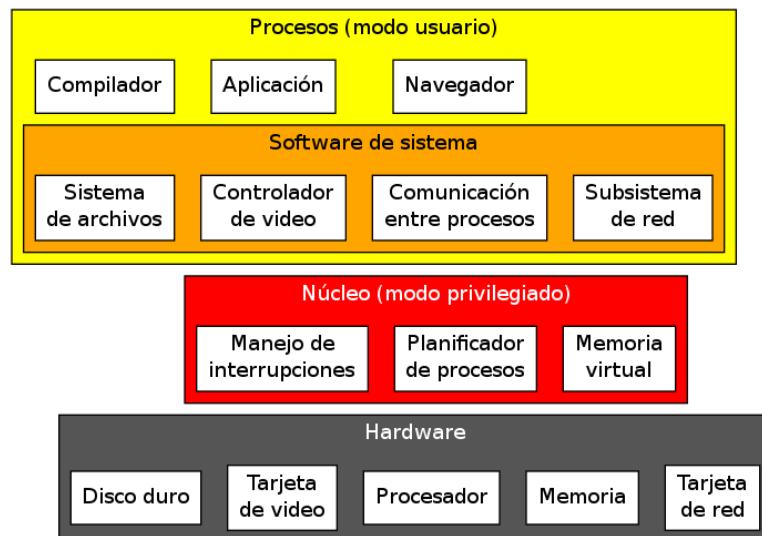


Figura 2.7: Esquematización de los componentes en un sistema microkernel

La principal ventaja de diseñar un sistema siguiendo un esquema monolítico es la simplificación de una gran cantidad de mecanismos de comunicación, que lleva a una

mayor velocidad de ejecución (al requerir menos cambios de contexto para cualquier operación realizada). Además, al manejarse la comunicación directa como paso de estructuras en memoria, el mayor acoplamiento permite más flexibilidad al adecuarse para nuevos requisitos (al no tener que modificar no sólo al núcleo y a los procesos especiales, sino también la interfaz pública entre ellos).

Por otro lado, los sistemas microkernel siguen esquemas lógicos más limpios, permiten implementaciones más elegantes y facilitan la comprensión por separado de cada una de sus piezas. Pueden *auto-repararse* con mayor facilidad, dado que en caso de fallar uno de los componentes (por más que parezca ser de muy bajo nivel), el núcleo puede reiniciarlo o incluso reemplazarlo.

Sistemas con concepciones híbridas No se puede hablar de concepciones únicas ni de verdades absolutas. A lo largo del libro se verán ejemplos de *concepciones híbridas* en este sentido — Sistemas que son mayormente monolíticos pero manejan algunos procesos que parecerían centrales a través de procesos de nivel usuario como los microkernel (por ejemplo, los sistemas de archivos en espacio de usuario, FUSE, en Linux).

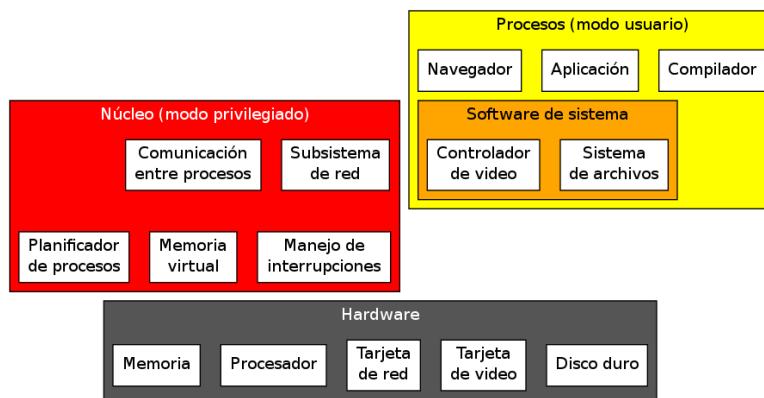


Figura 2.8: Esquematización de los componentes en un sistema híbrido

2.6 Otros recursos

- *On building systems that will fail*
<http://dl.acm.org/citation.cfm?id=1283947>
Fernando J. Corbató (1990); ACM Turing award lectures
- *A Brief History of Computer Operating Systems*
<http://cs.gordon.edu/courses/cs322/lectures/history.html>
R. Bjork (2000); Gordon College
- *Making EPERM friendlier*
<http://lwn.net/Articles/532771/>
Michael Kerrisk (2013); Linux Weekly News: Explica algunas de las limitantes de la semántica POSIX: Falta de granularidad en el reporte de mensajes de error (EPERM), y `errno` global por hilo.
- *Biculturalism*
<http://www.joelonsoftware.com/articles/Biculturalism.html>
Joel Spolsky (2003); Joel on Software

3 — Relación con el hardware

3.1 Introducción

Todos los sistemas de cómputo están compuestos por al menos una unidad de proceso junto con dispositivos que permiten ingresar datos (teclado, mouse, micrófono, etc.) y otros que permiten obtener resultados (pantalla, impresora, parlantes, etc.). Como se vio anteriormente, una de las funciones del sistema operativo es la de abstraer el hardware de la computadora y presentar al usuario una versión unificada y simplificada de los dispositivos. En este capítulo se verá la relación que mantiene el sistema operativo con el hardware, las funciones que cumplen y algunas abstracciones comunes utilizadas en sistemas operativos modernos.

3.2 Unidad de Procesamiento

La *unidad de procesamiento* es la parte fundamental de todo sistema de cómputo. Es la encargada de ejecutar tanto los programas del usuario como el sistema operativo en sí mismo. La funciones del sistema operativo respecto a la unidad de procesamiento son:

Inicialización Luego de ser cargado el sistema operativo debe realizar varias tareas de inicialización como habilitar las interrupciones de hardware y software (excepciones y trampas), configurar el sistema de memoria virtual (paginación, segmentación), etc.

Atender las interrupciones y excepciones Como se verá más adelante, la unidad de procesamiento puede encontrar una situación que no puede resolver por sí misma (una instrucción o dirección inválida, una división por cero, etc.) ante lo cual le pasa el control al sistema operativo para que éste trate o resuelva la situación.

Multiplexación En un sistema multiproceso, el sistema operativo es el encargado de administrar la unidad de procesamiento dando la ilusión a los procesos que están ejecutando de forma exclusiva.

3.2.1 Jerarquía de almacenamiento

Las computadoras que siguen la arquitectura *von Neumann*, esto es, prácticamente la totalidad hoy en día¹ podrían resumir su operación general a alimentar a una *unidad de proceso* (CPU) con los datos e instrucciones almacenados en *memoria*, que pueden incluir llamadas a servicio (y respuestas a eventos) originados en medios externos.

Una computadora von Neumann significa básicamente que es una computadora de *programa almacenado en la memoria primaria* — esto es, se usa el mismo almacenamiento para el programa que está siendo ejecutado y para sus datos, sirviéndose de un *registro*

¹Algunos argumentarán que muchas de las computadoras en uso hoy en día siguen la arquitectura *Harvard modificada*, dado que empleando distintos bancos de memoria caché, un procesador puede tanto referirse a la siguiente instrucción como iniciar una transferencia de memoria primaria. Esta distinción no tiene mayor relevancia para este tema, la referencia se incluye únicamente por no llevar a confusión.

especial para indicar al CPU cuál es la dirección en memoria de la siguiente instrucción a ejecutar.

La arquitectura von Neumann fue planteada, obviamente, sin considerar la posterior diferencia entre la velocidad que adquiriría el CPU y la memoria. En 1977, John Backus presentó al recibir el premio Turing un artículo describiendo el *cuello de botella de von Neumann*. Los procesadores son cada vez más rápidos (se logró un aumento de 1000 veces tanto entre 1975 y 2000 tan sólo en el reloj del sistema), pero la memoria aumentó su velocidad a un ritmo mucho menor —aproximadamente un factor de 50 para la tecnología en un nivel costo-beneficio suficiente para usarse como memoria primaria.

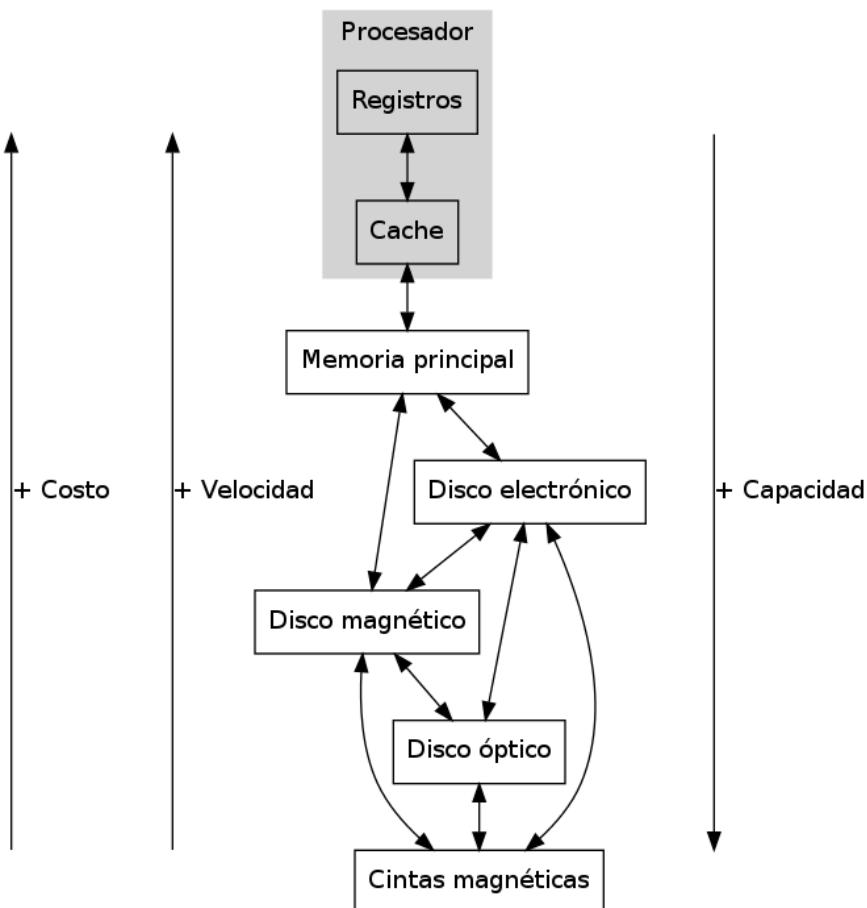


Figura 3.1: Jerarquía de memoria entre diversos medios de almacenamiento.

Una respuesta parcial a este problema es la creación de una jerarquía de almacenamiento, yendo de una pequeña área de memoria mucho más cara hasta un gran espacio de memoria muy económica. En particular, la relación entre las capas superiores está administrada por hardware especializado de modo que su existencia resulta transparente al programador.

Ahora bien, si bien la relación entre estos medios de almacenamiento puede parecer natural, para una computadora tiene una realidad completamente distinta: los registros son parte integral del procesador, y la memoria está a sólo un paso de distancia (el procesador puede referirse a ella directamente, de forma transparente, indicando la dirección desde un programa). El caché no existe para efectos prácticos: el procesador no hace referencia directa a él, sino que es manejado por los controladores de acceso a memoria.

Como se verá, el sistema operativo es el encargado de mantener todas estas jerarquías de memoria consistentes y de realizar las transferencias entre unas y otras.

Cuadro 3.1: Velocidad y gestor de los principales niveles de memoria. (Silberschatz, Galvin, Gagne; p.28)

Nivel	1	2	3	4
Nombre	Registros	Cache	Memoria princ.	Disco
Tamaño	<1KB	<16MB	<64GB	>100GB
Tecnología	Multipuerto, CMOS	SRAM CMOS	CMOS DRAM	Magnética
Acceso (ns)	0.25-0.5	0.5-25	80-250	5,000,000
Transf (MB/s)	20,000-100,000	5,000-10,000	1,000-5,000	20-150
Administra	Compilador	Hardware	Sist. Op.	Sist. op.
Respaldado en	Cache	Memoria princ.	Disco	CD o cinta

Registros

La memoria más rápida de la computadora son los *registros*, ubicados dentro de cada *uno de los núcleos* de cada uno de los CPU. Las arquitecturas tipo RISC (Reduced Instruction Set Computer) sólo contemplan la ejecución de instrucciones entre registros (excepto, claro, las de carga y almacenamiento a memoria primaria).

Los primeros CPU trabajaban con pocos registros, muchos de ellos de propósito específico — trabajaban más bien con una lógica de *registro acumulador*. Por ejemplo, el MOS 6502 (en el cual se basaron las principales computadoras de 8 bits) tenía un acumulador de 8 bits (A), dos registros índice de 8 bits (X e Y), un registro de estado del procesador de 8 bits (P), un apuntador al *stack* de 8 bits (S), y un apuntador al programa de 16 bits (PC). El otro gran procesador de su era, el Zilog Z80, tenía 14 registros (3 de 8 bits y el resto de 16), pero sólo uno era un acumulador de propósito general.

El procesador Intel 8088, en el cual se basó la primera generación de la arquitectura PC, ofrecía cuatro registros de uso *casi* general. En los ochenta comenzaron a producirse los primeros procesadores tipo RISC, muchos de los cuales ofrecían 32 registros, todos ellos de propósito general.

Registros de propósito general	
AH	AL
BH	BL
CH	CL
DH	DL
Registros índices	
SI	Source Index (Índice origen)
DI	Destination Index (Índice Destino)
BP	Base Pointer (Puntero Base)
SP	Stack Pointer (Puntero de Pila)
Registro de Bandera	
- - - - O D I T S Z - A - P - C	Flags (Banderas)
Registros de Segmentos	
CS	Code Segment (Segmento de Código)
DS	Data Segment (Segmento de Datos)
ES	ExtraSegment (Segmento Extra)
SS	Stack Segment (Segmento de Pila)
Registro apuntador de instrucciones	
IP	Instruction Pointer

Figura 3.2: Ejemplo de registros: Intel 8086/8088 (Imagen de la Wikipedia: *Intel 8086 y 8088*)

El compilador² busca realizar muchas operaciones que deben ocurrir reiteradamente, donde la rapidez es fundamental, con sus operadores cargados en los registros. El estado del CPU en un momento dado está determinado por el contenido de los registros. El contenido de la memoria, obviamente, debe estar sincronizado con lo que ocurre dentro de este — pero el estado actual del CPU, lo que está haciendo, las indicaciones respecto a las operaciones recién realizadas que se deben entregar al programa en ejecución están todos representados en los registros. Se debe mantener esto en mente cuando posteriormente se habla de todas las situaciones en que el flujo de ejecución debe ser quitado de un proceso y entregado a otro.

La relación de la computadora y del sistema operativo con la memoria principal será abordada en el capítulo 6.

3.2.2 Interrupciones y excepciones

La ejecución de los procesos podría seguir siempre linealmente, atendiendo a las instrucciones de los programas tal como fueron escritas, pero en el modelo de uso de cómputo actual, eso no serviría de mucho: para que un proceso acepte interacción, su ejecución debe poder responder a los *eventos* que ocurran alrededor del sistema. Y los eventos son manejados a través de las *interrupciones y excepciones* (o *trampas*).

Cuando ocurre algún evento que requiera la atención del sistema operativo, el hardware encargado de procesarlo escribe directamente a una ubicación predeterminada de memoria la naturaleza de la solicitud (el *vector de interrupción*) y, levantando una solicitud de interrupción, *detiene* el proceso que estaba siendo ejecutado. El sistema operativo entonces ejecuta su *rutina de manejo de interrupciones* (típicamente comienza grabando el estado de los registros del CPU y otra información relativa al estado del proceso desplazado) y posteriormente la atiende.

Las interrupciones pueden organizarse por *prioridades*, de modo que una interrupción de menor jerarquía no interrumpa a una más importante — dado que las interrupciones muchas veces indican que hay datos disponibles en algún buffer, el no atenderlas a tiempo podría llevar a la pérdida de datos.

Hay un número limitado de interrupciones definidas para cada arquitectura, mucho más limitado que el número de dispositivos que tiene un equipo de cómputo actual. Las interrupciones son, por tanto, generadas *por el controlador del canal* en que son producidas. Si bien esto resuelve la escasez de interrupciones, dificulta su priorización — con canales de uso tan variado como el USB³, una interrupción puede indicar que hay desde un *teclado* para ser leído hasta un paquete de red esperando a ser procesado — y si bien demorar la atención al primero no llevaría a pérdida notable de información, no atender el paquete de red sí.

El sistema operativo puede elegir ignorar (*enmascarar*) a ciertas interrupciones — pero hay interrupciones que son *no enmascarables*.

Se hace la distinción entre interrupciones y excepciones según su origen: una interrupción es generada por causas externas al sistema (un dispositivo requiere atención), mientras que una excepción es una evento generado por un proceso (una condición en el proceso que requiere la intervención del sistema operativo). Si bien hay distinciones sutiles entre interrupciones, trampas y excepciones, al nivel de discusión que se abordará basta con esta distinción.

Los eventos pueden ser, como ya se mencionó, indicadores de que hay algún dispositivo requiriendo atención, pero pueden también provenir del mismo sistema, como una

²A veces asistido por instrucciones explícitas por parte del programador, pero muchas veces como resultado del análisis del código.

³Algunas arquitecturas, particularmente de sistemas embebidos y por un criterio altamente económico, están estructuradas íntegramente alrededor de un bus USB.

alarma o temporizador (que se emplea para obligar a todo programa a entregar el control en un sistema multitareas) o indicando una condición de error (por ejemplo, una división sobre cero o un error leyendo de disco).

Las funciones del sistema operativo respecto a las interrupciones son:

Administrar el hardware manejar de interrupciones Esto incluye el enmascarado y desenmascarado de las interrupciones, configurar y asignar interrupciones a cada dispositivo, notificar al manejador cuando la interrupción ya ha sido atendida, etc.

Abstraer las interrupciones El sistema operativo oculta a los programas de usuario la existencia de interrupciones de hardware ya que éstas son dependientes de la arquitectura del procesador. En cambio el sistema operativo lo comunica de una forma unificada a través de distintos mecanismos, por ejemplo mensajes o señales o deteniendo el proceso que espera la acción relacionada con una interrupción y continuando su ejecución cuando ésta ocurre.

Punto de entrada al sistema operativo Como se verá más adelante en la sección 3.7, muchos procesadores y sistemas operativos utilizan las interrupciones como medio por el cual un proceso de usuario realiza una llamada al sistema. Por ejemplo, en Linux para arquitecturas x86 el programa de usuario genera la interrupción 0x80 para iniciar una llamada al sistema. En arquitecturas más recientes como x86₆₄, MIPS y ARM esto ha sido reemplazado por una instrucción especial `syscall`.

Atender excepciones y fallas Como se discutió antes, durante la ejecución de un programa pueden ocurrir situaciones anómalas, como por ejemplo, una división sobre cero. Desde el punto de vista del CPU, esto es similar a una interrupción de hardware y debe ser tratada por el sistema operativo. Dependiendo de la causa de la excepción, el sistema operativo tomará acción para resolver en lo posible esta situación. En muchos casos las excepciones resultan en una señal enviada al proceso, y este último es el encargado de tratar la excepción. En otros casos la falla o excepción son irrecuperables (una instrucción inválida o un error de bus) ante la cual el sistema operativo terminará el proceso que la generó. En el capítulo 6 se cubre con mucho mayor detalle un tipo de excepción muy importante que debe tratar el sistema operativo: el fallo de paginación.

3.3 Terminales

Las terminales son dispositivos electrónicos utilizados para ingresar datos y emitir resultados dentro de un sistema de cómputo. Las primeras terminales utilizaban tarjetas perforadas e impresiones en papel. Debido a su limitada velocidad e imposibilidad de “editar” el papel ya impreso, este tipo de terminales fue cediendo terreno ante la aparición sobre principios de los setenta de las terminales de texto con pantalla de video y teclado.

Conceptualmente una terminal de texto es un dispositivo mediante el cual la computadora recibe y envía un flujo de caracteres desde y hacia el usuario respectivamente. Las operaciones más complejas, como edición, borrado y movimiento, en general son tratadas con *secuencias de escape*, esto es, una serie de caracteres simples que tomados en conjunto representan una acción a realizar en la terminal.

Durante la década de los setenta también se desarrollaron terminales gráficas las cuales podían representar imágenes junto con texto. Con la inclusión del ratón o “mouse” estas terminales dieron lugar a lo que hoy se conoce como Interfaz Gráfica de Usuario (*Graphical User Interface* o *GUI*) y a los sistemas de ventana.

En los sistemas operativos modernos es común referirse al *emulador de terminal*, un programa especializado ya sea para tener múltiples instancias de una terminal o para ejecutar una terminal de texto dentro de una interfaz gráfica. Estos programas se denominan de esta forma dado que sólo replican el comportamiento de las terminales (que eran

originalmente equipos independientes), siendo únicamente un programa que recibe la entrada del usuario a través del teclado enviándola al sistema operativo como un flujo de datos, y recibe otro flujo de datos del sistema operativo, presentándolo de forma adecuada al usuario.

3.4 Dispositivos de almacenamiento

El almacenamiento en memoria primaria es *volátil*, esto es, se pierde al interrumpirse el suministro eléctrico. Esto no era muy importante en la época definitoria de los conceptos que se presentan en esta sección, dado que el tiempo total de vida de un conjunto de datos en almacenamiento bajo el control del procesador iba únicamente desde la entrada y hasta el fin de la ejecución del trabajo del usuario. Pero desde la década de los sesenta se popularizó la posibilidad de almacenar en la computadora información *a largo plazo* y con expectativas razonables de permanencia.

De las muchas tecnologías de almacenamiento, la que ha dominado fuertemente durante los últimos 40 años ha sido la de los discos magnéticos⁴. El acceso a disco (miles de veces más lento que el acceso a memoria) no es realizado directamente por el procesador, sino que requiere de la comunicación con controladores externos, con lógica propia, que podrían ser vistos como computadoras independientes de propósito limitado.

El procesador no puede referirse directamente más información que la que forma parte del almacenamiento primario — esto es, de la memoria RAM. En las secciones 3.2.2 (*Interrupciones y excepciones*) y 3.6.2 (*Acceso directo a memoria*), se explica cómo es que se efectúan dichas referencias.

Los dispositivos de almacenamiento (discos, memorias flash, cintas) pueden ser vistos como una región donde la computadora lee y escribe una serie de bytes que preservarán su valor incluso luego de apagada la computadora.

A nivel de hardware el sistema operativo no accede al dispositivo de almacenamiento byte por byte, sino que éstos se agrupan en *bloques* de tamaño fijo. El manejo de estos bloques (administración de bloques libres, lectura y escritura) es una tarea fundamental del sistema operativo, que asimismo se encarga de presentar abstracciones como la de archivos y directorios al usuario. Esto se verá en el capítulo 7.

3.5 Relojes y temporizadores

Todas las computadoras incluyen uno o más relojes y temporizadores que son utilizados para funciones varias como mantener la hora del sistema actualizada, implementar alarmas tanto para los programas de usuario como para el sistema operativo, ejecutar tareas de mantenimiento periódicas, cumplir con requisitos temporales de aplicaciones de tiempo real, etc.

Mantener el tiempo correctamente dentro del sistema operativo es algo crucial. Permite establecer un orden cronológico entre los eventos que ocurren dentro del sistema, por ejemplo la creación de un archivo y de otro o el tiempo consumido en la ejecución de un proceso.

Por otro lado si el sistema operativo utiliza una política de planificación de procesos preventiva (capítulo 5), como la *Ronda (Round Robin)*, éste debe interrumpir al proceso en ejecución luego de cierta cantidad de unidades de tiempo. Esto se implementa haciendo que el temporizador de la computadora genere interrupciones periódicamente, lo cual luego invocará al planificador de procesos.

⁴Se verán en la sección C.1.2 detalles acerca de las tecnologías de almacenamiento *en estado sólido*, que pueden poner fin a esta larga dominación.

3.6 Canales y puentes

Los distintos componentes de un sistema de cómputo se comunican a través de los diferentes *canales* (generalmente se hace referencia a ellos por su nombre en inglés: *buses*). Al nivel más básico, los canales son líneas de comunicación entre el procesador y los demás componentes del chipset⁵, a los cuales a su vez se conectan los diferentes dispositivos del sistema — desde aquellos que requieren mayor velocidad, como la misma memoria, hasta los puertos más sencillos.

Un chipset provee distintos buses, con un agrupamiento lógico según la velocidad requerida por sus componentes y otras características que determinan su topología.

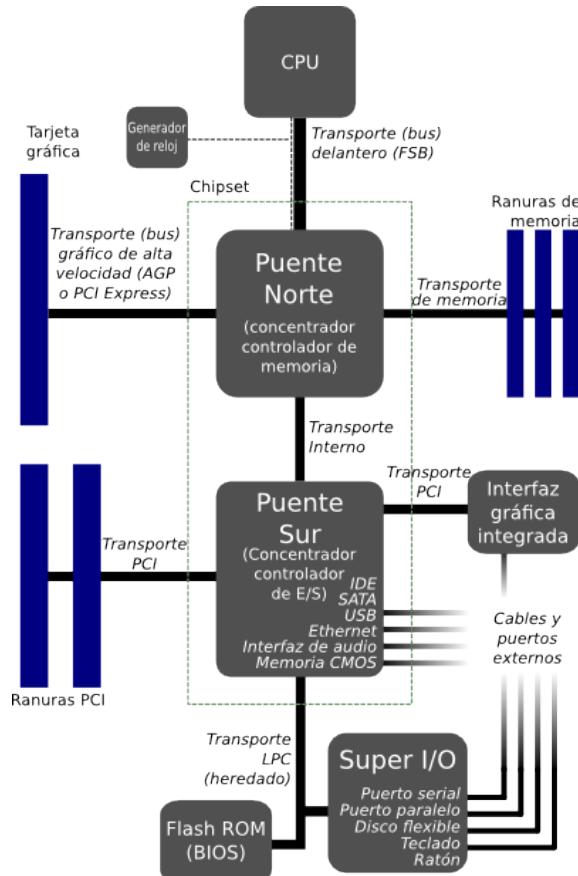


Figura 3.3: Diagrama de la comunicación entre componentes de un sistema de cómputo basado en puente norte y puente sur

Hoy en día, el acomodo más frecuente⁶ de estos buses es a través de una separación en dos chips: el *puente norte* (*Northbridge*), conectado directamente al CPU, encargado de gestionar los buses de más alta velocidad y que, además, son fundamentales para el más básico inicio de la operación del sistema: la memoria y el reloj. La comunicación con algunas tarjetas de video se incorpora al puente norte a través del canal dedicado *AGP* (*Advanced Graphics Port*, *Puerto Gráfico Avanzado*).

Al puente norte se conecta el *puente sur* (*Southbridge*), que controla el resto de los dispositivos del sistema — normalmente se ven aquí las interfaces de almacenamiento

⁵Los chips que forman parte de un equipo, casi siempre provistos por el mismo fabricante que el procesador mismo.

⁶La separación aquí descrita ha sido característica de las computadoras x86 de los últimos 20 años, aunque la tendencia apunta a que se abandone paulatinamente para dar paso a procesadores que integren en un sólo paquete todos estos componentes. Sin embargo, el acomodo funcional electrónico, al menos hasta el momento, sigue basado en estos puntos.

(SCSI, SATA, IDE), de expansión interna (PCI, PCIe) y de expansión externa (USB, Firewire, puertos *heredados* seriales y paralelos).

3.6.1 Contención

Una de las principales razones de la existencia de tantos *canales* (buses) distintos en un mismo sistema es a la frecuencia acorde a los dispositivos para los cuales está diseñado: la cantidad de datos que tiene que viajar entre el procesador y la memoria a lo largo de la operación del sistema es muy superior que la que tiene que transferirse desde los discos, y a su vez, esta es mucho mayor que la que enviarse a la impresora, o la que se recibe del teclado. Claro está, los demás dispositivos podrían incrementar su frecuencia para participar en un canal más rápido, aunque su costo se incrementaría, dado que harían falta componentes capaces de sostener un reloj varios órdenes de magnitud más rápido.

Pero incluso obviando la diferencia económica: cuando el sistema requiere transferir datos de o hacia varios dispositivos de la misma categoría, es frecuente que ocurra *contención*: puede saturarse el ancho de banda máximo que alcanza uno de los canales y, aún si los dispositivos tienen información lista, tendrán que esperar a que los demás dispositivos desocupen el canal.

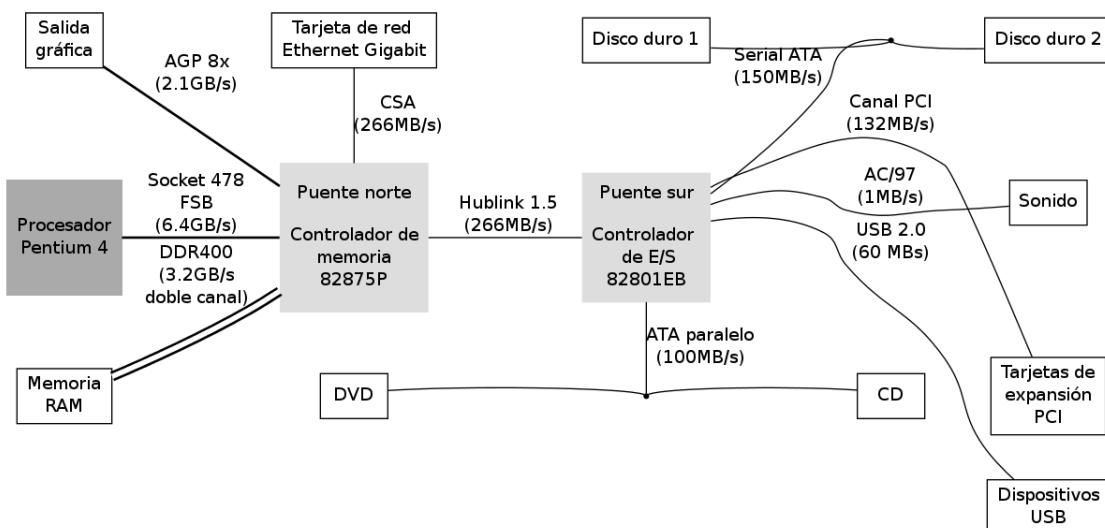


Figura 3.4: Esquema simplificado del chipset Intel 875 (para el procesador Pentium 4) ilustrando la velocidad de cada uno de los canales

En la figura 3.4 se puede ver el diseño general del chipset Intel 875, introducido en el 2003, incluyendo el ancho de banda de cada uno de los canales del sistema. Hay que recordar que hay canales como el USB que permiten la conexión de múltiples dispositivos, los cuales deberán compartir el ancho de banda total permitido por el canal: en la figura se presentan dos discos duros sobre el canal SATA y dos unidades ópticas en el ATA paralelo; el canal USB permite el uso de un máximo de 127 unidades por canal, por lo cual la contención puede ser muy alta.

3.6.2 Acceso directo a memoria (DMA)

La operación de dispositivos de entrada/salida puede ser altamente ineficiente. Cuando un proceso está en una sección *limitada por entrada-salida* (esto es, donde la actividad principal es la transferencia de información entre la memoria principal y cualquier otra área del sistema), si el procesador tiene que encargarse de la transferencia de toda la

información⁷, se crearía un cuello de botella por la cantidad y frecuencia de interrupciones. Hoy en día, para evitar que el sistema se demore cada vez que hay una transferencia grande de datos, todas las computadoras implementan controladores de *acceso directo a memoria* (DMA) en uno o más de sus subsistemas.

El DMA se emplea principalmente al tratar con dispositivos con un gran ancho de banda, como unidades de disco, subsistemas multimedia, tarjetas de red, e incluso para transferir información entre niveles del caché.

Las transferencias DMA se hacen en *bloques* preestablecidos; en vez de que el procesador reciba una interrupción cada vez que hay una palabra lista para ser almacenada en la memoria, el procesador indica al controlador DMA la dirección física base de memoria en la cual operará, la cantidad de datos a transferir, el *sentido* en que se efectuará la operación (del dispositivo a memoria o de memoria al dispositivo), y el *puerto* del dispositivo en cuestión; el controlador DMA efectuará la transferencia solicitada, y sólo una vez terminada ésta (o en caso de encontrar algún error en el proceso) lanzará una interrupción al sistema; el procesador queda libre para realizar otras tareas, sin más limitante que la posible *contención* que tendrá que enfrentar en el bus de acceso a la memoria.

Coherencia de cache

Cuando se realiza una transferencia DMA de un dispositivo a la memoria, puede haber páginas de la memoria en cuestión que estén en alguno de los niveles de la memoria caché; dado que el caché está uno o más niveles por encima de la memoria principal, es posible que la información haya ya cambiado pero el caché retenga la información anterior.

Los sistemas de *caché coherente* implementan mecanismos en hardware que notifican a los controladores de caché que las páginas que alojan están *sucias* y deben ser vueltas a cargar para ser empleadas, los sistemas *no coherentes* requieren que el subsistema de memoria del sistema operativo haga esta operación.

Los procesadores actuales implementan normalmente varios niveles de caché, estando algunos dentro del mismo CPU, por lo que típicamente se encuentran sistemas híbridos, en los que los cachés de nivel 2 son coherentes, pero los de nivel 1 no, y deben ser manejados por software.

3.7 Interfaz del Sistema Operativo: llamadas al sistema

De forma análoga a las interrupciones, se puede hablar de las llamadas al sistema. El sistema operativo protege a un proceso de otro, y previene que un proceso ejecutándose en espacio no privilegiado tenga acceso directo a los dispositivos. Cuando un proceso requiere de alguna acción privilegiada, accede a ellas realizando una *llamada al sistema*. Las llamadas al sistema pueden agruparse, a grandes rasgos, en:

Control de procesos Crear o finalizar un proceso, obtener atributos del proceso, esperar la finalización de un proceso o cierto tiempo, asignar o liberar memoria, etc.

Manipulación de archivos Crear, borrar o renombrar un archivo; abrir o cerrar un archivo existente; modificar sus *metadatos*; leer o escribir de un *descriptor de archivo* abierto, etc.

Manipulación de dispositivos Solicitar o liberar un dispositivo; leer, escribir o reposicionarlo, y otras varias. Muchas de estas llamadas son análogas a las de manipulación de archivos, y varios sistemas operativos las ofrecen como una sola.

Mantenimiento de la información Obtener o modificar la hora del sistema; pedir detalles acerca de procesos o archivos, etc.

Comunicaciones Establecer una comunicación con determinado proceso (local o remoto), aceptar una solicitud de comunicación de otro proceso, intercambiar información

⁷Este modo de operación es también conocido como *entrada/salida programada*.

sobre un canal establecido.

Protección Consultar o modificar la información relativa al acceso de objetos en el disco, otros procesos, o la misma sesión de usuario.

Cada sistema operativo *expone* una serie de llamadas al sistema. Estas son, a su vez, expuestas al programador a través de las *interfaces de aplicación al programador* (API), que se alinean de forma cercana (pero no exacta). Del mismo modo que cada sistema operativo ofrece un conjunto de llamadas al sistema distinto, cada implementación de un lenguaje de programación puede ofrecer un API ligeramente distinto de otros.

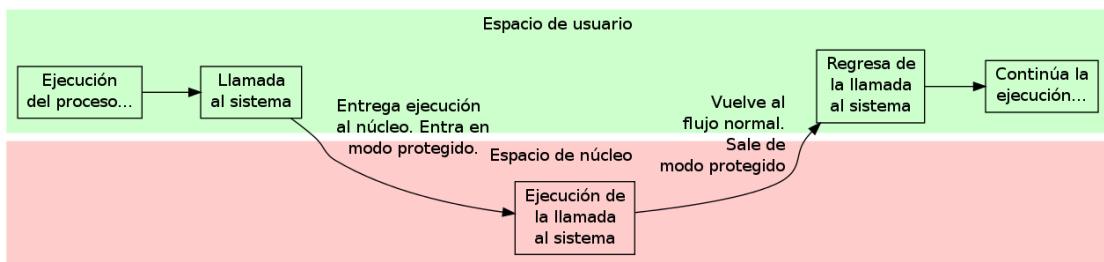


Figura 3.5: Transición del flujo entre espacio usuario y espacio núcleo en una llamada al sistema

3.7.1 Llamadas al sistema, arquitecturas y APIs

Cada familia de sistemas operativos provee distintas llamadas al sistema, y sus lenguajes/bibliotecas implementan distintos APIs. Esto es el que distingue principalmente a uno de otro. Por ejemplo, los sistemas Windows 95 en adelante implementan Win32, Win16 (compatibilidad con Windows previos) y MSDOS; MacOS implementa Cocoa (aplicaciones MacOS X) y Carbon (compatibilidad con aplicaciones de MacOS previos), y Linux y los *BSDs, POSIX (el estándar que define a Unix). El caso de MacOS X es interesante, porque también implementa POSIX, ofreciendo la *semántica* de dos sistemas muy distintos entre sí.

Los lenguajes basados en *máquinas virtuales abstractas*, como Java o la familia .NET, exponen un API con mucha mayor distancia respecto al sistema operativo; la máquina virtual se presenta como un pseudo-sistema operativo intermedio que se ejecuta dentro del real, y esta distinción se hace especialmente notoria cuando se busca conocer los detalles del sistema operativo.

Depuración por trazas (trace)

La mayor parte de los sistemas operativos ofrecen programas que, para fines de depuración, *envuelven* al API del sistema y permiten ver la *traza* de las llamadas al sistema que va realizando un proceso. Algunos ejemplos de estas herramientas son strace en Linux, truss en la mayor parte de los Unixes históricos o ktrace y kdump en los *BSD. A partir de Solaris 10 (2005), Sun incluye una herramienta mucho más profunda y programable para esta tarea llamada dtrace, que al paso del tiempo ha sido *portada*⁸ a otros Unixes (*BSD, MacOS).

La salida de una traza brinda amplio detalle acerca de la actividad realizada por un proceso, y permite comprender a grandes rasgos su interacción con el sistema. El nivel de información que da es, sin embargo, a veces demasiado — eso se puede ver si se considera la siguiente traza, ante uno de los comandos más sencillos: pwd (obtener el directorio actual)

⁸Se denomina *portar* el hacer las adecuaciones necesarias para que una herramienta diseñada para determinado entorno pueda emplearse en otros distintos.

```
$ strace pwd
execve("./bin/pwd", ["pwd"], /* 43 vars */) = 0
brk(0) = 0x8414000
access("/etc/ld.so.nohwcap", F_OK) = -1 ENOENT (No such file or directory)
mmap2(NULL, 8192, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS, -1, 0) = 0xb773d000
access("/etc/ld.so.preload", R_OK) = -1 ENOENT (No such file or directory)
open("/etc/ld.so.cache", O_RDONLY) = 3
fstat64(3, {st_mode=S_IFREG|0644, st_size=78233, ...}) = 0
mmap2(NULL, 78233, PROT_READ, MAP_PRIVATE, 3, 0) = 0xb7729000
close(3) = 0
access("/etc/ld.so.nohwcap", F_OK) = -1 ENOENT (No such file or directory)
open("/lib/i386-linux-gnu/libc.so.6", O_RDONLY) = 3
read(3, "\177ELF\1\1\0\0\0\0\0\0\0\3\0\1\0\0\0po\1\0004\0\0\0...", 512) = 512
fstat64(3, {st_mode=S_IFREG|0755, st_size=1351816, ...}) = 0
mmap2(NULL, 1366328, PROT_READ|PROT_EXEC, MAP_PRIVATE|MAP_DENYWRITE, 3, 0) = 0xb75db000
mprotect(0xb7722000, 4096, PROT_NONE) = 0
mmap2(0xb7723000, 12288, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0x147) = 0xb7723000
mmap2(0xb7726000, 10552, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP_ANONYMOUS, -1, 0) = 0xb7726000
close(3) = 0
mmap2(NULL, 4096, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS, -1, 0) = 0xb75da000
set_thread_area({entry_number:-1 -> 6, base_addr:0xb75da8d0, limit:1048575, seg_32bit:1, contents:0, read_e
mprotect(0xb7723000, 8192, PROT_READ) = 0
mprotect(0xb775c000, 4096, PROT_READ) = 0
munmap(0xb7729000, 78233) = 0
brk(0) = 0x8414000
brk(0x8435000) = 0x8435000
open("/usr/lib/locale/locale-archive", O_RDONLY|O_LARGEFILE) = 3
fstat64(3, {st_mode=S_IFREG|0644, st_size=1534672, ...}) = 0
mmap2(NULL, 1534672, PROT_READ, MAP_PRIVATE, 3, 0) = 0xb7463000
close(3) = 0
getcwd("/home/gwolf/vcs/sistemas_operativos", 4096) = 36
fstat64(1, {st_mode=S_IFCHR|0620, st_rdev=makedev(136, 1), ...}) = 0
mmap2(NULL, 4096, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS, -1, 0) = 0xb773c000
write(1, "/home/gwolf/vcs/sistemas_operati...", 36/home/gwolf/vcs/sistemas_operativos
) = 36
close(1) = 0
munmap(0xb773c000, 4096) = 0
close(2) = 0
exit_group(0) = ?
```

3.8 Abstracciones comunes

Como se mencionó antes una de las funciones del sistema operativo es la de abstraer el hardware de la computadora. Cuál es esa abstracción que ve efectivamente el usuario varía de un sistema operativo a otro. Se verán en esta sección algunas abstracciones utilizadas en varios sistemas dejando las correspondientes a sistemas de archivos para el capítulo 8.

3.8.1 Sistemas tipo Windows

Los sistemas del tipo Windows presentan una abstracción diversa para cada uno de los componentes de la computadora.

Por ejemplo los volúmenes de almacenamiento secundario (discos rígidos, discos compactos, memorias flash, etc.) son relacionados con una letra cada uno, así (en general) C: es el volumen o partición del disco principal, A: y B: se utilizan para discos extraibles.

Una desventaja de esta abstracción es que no queda claro cuáles unidades pertenecen al mismo disco físico y cuáles no.

A los puertos de entrada/salida más utilizados también se les asignan nombres alfanuméricos, por ejemplo el primer puerto paralelo se denomina LPT1 y el segundo puerto serie COM2.

3.8.2 Sistemas tipo Unix

Unix introdujo el concepto de que *todo es un archivo*: en el sistema Unix original, todos los dispositivos podían ser controlados a través de un *archivo especial* que, en vez de

almacenar información, apunta a estructuras en el sistema que controlan a cada dispositivo. Este concepto sobrevive en los sistemas derivados de Unix al día de hoy, aunque varias clases de dispositivo rompen esta lógica. El sistema operativo *Plan9* de Bell Labs mantiene y amplía este concepto e introduce los *espacios de nombres mutables*, que presenta con interfaz de archivo prácticamente cualquier objeto empleado por el sistema.

Las principales estructuras relacionadas de este tipo que existen en un sistema tipo Unix son:

Dispositivos de caracteres Son aquellos en los cuales la información es leída o escrita de a un carácter a la vez y se presentan como *streams* (flujos) de información, ya sea entrante, saliente o mixta. Algunos pueden permitir operaciones adicionales (por ejemplo, rebobinado), pero la manipulación de la información se efectúa de forma secuencial.

Ejemplos: impresora, unidad de cinta, modem.

Dispositivos de bloques Presentan una interfaz de *acceso aleatorio* y entregan o reciben la información en *bloques* de tamaño predeterminado.

El ejemplo más claro de este tipo de dispositivos es una unidad de disco o una de sus particiones.

Archivos especiales Los sistemas Unix actuales incluyen también un gran número de archivos especiales, por medio de los cuales el usuario puede monitorear el estado del sistema (memoria libre, número de procesos, consumo de procesador, etc.), e incluso modificar la configuración del sistema operativo a través de un archivo; por ejemplo, en un sistema Linux, escribir el valor “100” al archivo “/proc/sys/vm/swappiness” hará que el sistema envíe a espacio de intercambio una mayor cantidad de programas de lo que normalmente haría.

3.9 Cuando dos cabezas piensan mejor que una

3.9.1 Multiprocesamiento

El *multiprocesamiento* es todo entorno donde hay más de un procesador (CPU). En un entorno multiprocesado, el conjunto de procesadores se vuelve un recurso más a gestionar por el sistema operativo — y el que haya concurrencia *real* tiene un fuerte impacto en su diseño.

Si bien en el día a día se usan de forma intercambiable⁹, es importante enfatizar en la diferencia fundamental entre el *multiprocesamiento*, que se abordará en esta sección, y la *multiprogramación*, de la cual se hablará en la sección 2.3.3 (*Sistemas multiprogramados*). Un sistema multiprogramado da la *ilusión* de que está ejecutando varios procesos al mismo tiempo, pero en realidad está alternando entre los diversos procesos que compiten por su atención. Un sistema multiprocesador tiene la capacidad de estar atendiendo *simultáneamente* a diversos procesos.

En la figura 3.6, el primer diagrama ilustra una ejecución estrictamente secuencial: cada uno de los procesos que demandan atención del sistema es ejecutado hasta que termina; el segundo muestra cómo se comportaría un sistema multiprogramado, alternando entre los tres procesos, de modo que el usuario vea que los tres avanzan de forma simultánea; el tercero corresponde a un sistema de multitarea pura: cada proceso es ejecutado por un procesador distinto, y avanzan en su ejecución de forma simultánea. El cuarto caso, un esquema híbrido, presenta cómo reaccionaría un equipo con capacidad de atender a dos procesos al mismo tiempo, pero con tres procesos solicitando ser atendidos. Este último esquema es el que más comúnmente se encuentra en equipos de uso general hoy en día.

⁹O poco más que eso, al grado de que rara vez se emplea el término *multiprogramación*, mucho más acorde a los equipos que se emplean día a día.

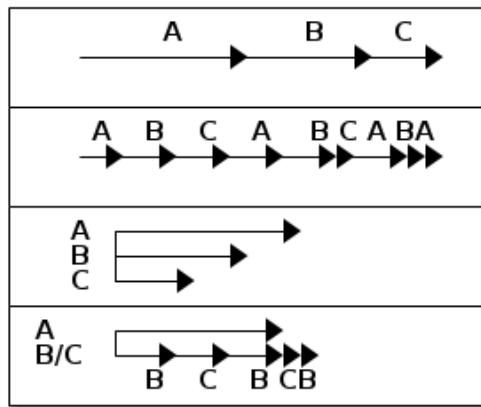


Figura 3.6: Esquema de la ejecución de tres procesos en un sistema secuencial, multiprogramado, multiprocesado, e híbrido

Probablemente el tema que se aborda más recurrentemente a lo largo de este texto será precisamente la complejidad que proviene de la multiprogramación; se la desarrollará particularmente en los capítulos 4 y 5. Valga la nota en este momento únicamente para aclarar la diferencia entre los dos conceptos.

El multiprocesamiento se emplea ampliamente desde los sesenta en los entornos de cómputo de alto rendimiento, pero por muchos años se vio como el área de especialización de muy pocos — las computadoras con más de un procesador eran prohibitivamente caras, y para muchos sistemas, ignorar el problema resultaba una opción válida. Muchos sistemas operativos ni siquiera detectaban la existencia de procesadores adicionales, y en presencia de éstos, ejecutaban en uno sólo.

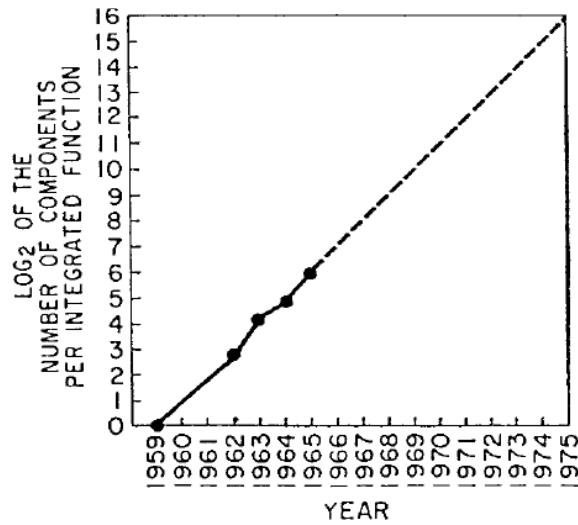


Figura 3.7: La Ley de Moore, en su artículo publicado en 1965, prediciendo la miniaturización por diez años

Esto cambió hacia el 2005. Tras más de 40 años de cumplirse, el modelo conocido como la *Ley de Moore*, enunciando que cada dos años la densidad de transistores por circuito integrado se duplicaría, llevaba a velocidades de CPU que, en el ámbito comercial, excedían los 3GHz, lo cual presentaba ya problemas serios de calentamiento. Además, el diferencial de velocidad con el acceso a memoria era cada vez más alto. Esto motivó a que

las principales compañías productoras de los CPU cambiaran de estrategia, introduciendo chips que son, para propósitos prácticos, paquetes con 2 o más procesadores dentro.

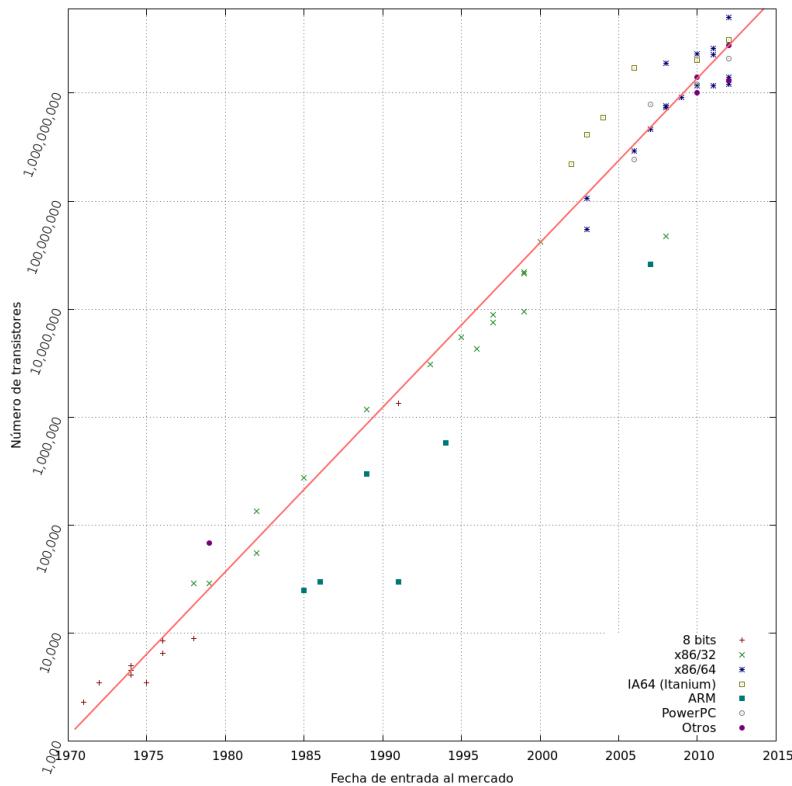


Figura 3.8: La Ley de Moore se sostiene al día de hoy: conteo de transistores por procesador de 1971 al 2012

Con este cambio, el *reloj* de los procesadores se ha mantenido casi sin cambios, cerca de 1GHz, pero el rendimiento de los equipos sigue aumentando. Sin embargo, los programadores de sistemas operativos y programas de aplicación ya no pueden ignorar esta complejidad adicional.

Se denomina *multiprocesamiento simétrico* (típicamente abreviado SMP) a la situación en la que todos los procesadores del sistema son iguales y pueden realizar en el mismo tiempo las mismas operaciones. Todos los procesadores del sistema tienen acceso a la misma memoria (aunque cada uno puede tener su propio *caché*, lo cual obliga a mantener en mente los puntos relacionados con la *coherencia de caché* abordados en la sección anterior).

Existe también el *multiprocesamiento asimétrico*; dependiendo de la implementación, la asimetría puede residir en diferentes puntos. Puede ir desde que los procesadores tengan una *arquitectura* distinta (típicamente dedicada a una tarea específica), en cuyo caso pueden verse como *coprocesadores* o *procesadores coadyuvantes*, casi computadoras independientes contribuyendo sus resultados a un mismo cómputo. Hoy en día, este sería el caso de las tarjetas gráficas 3D, que son computadoras completas con su propia memoria y responsabilidades muy distintas del sistema central.

Es posible tener diferentes procesadores con la misma arquitectura pero funcionando a diferente frecuencia. Esto conlleva una fuerte complejidad adicional, y no se utiliza hoy en día.

Por último, existen los diseños de *Acceso No-Uniforme a Memoria* (*Non-Uniform Memory Access, NUMA*). En este esquema, cada procesador tiene *afinidad* con bancos específicos de memoria — para evitar que los diferentes procesadores estén esperando al mismo tiempo al bus compartido de memoria, cada uno tiene acceso exclusivo a su área. Los sistemas

NUMA pueden ubicarse como en un punto intermedio entre el procesamiento simétrico y el cómputo distribuído, y puede ser visto como un *cómputo distribuído fuertemente acoplado*.

3.9.2 Cómputo distribuído

Se denomina cómputo distribuído a un proceso de cómputo realizado entre computadoras independientes, o, más formalmente, entre procesadores que *no comparten memoria* (almacenamiento primario). Puede verse que un equipo de diseño NUMA está a medio camino entre una computadora multiprocesada y el cómputo distribuído.

Hay diferentes modelos para implementar el cómputo distribuído, siempre basados en la transmisión de datos sobre una *red*. Estos son principalmente:

Cúmulos (clusters) Computadoras conectadas por una red local (de alta velocidad), ejecutando cada una su propia instancia de sistema operativo. Pueden estar orientadas al *alto rendimiento, alta disponibilidad* o al *balanceo de cargas* (o a una combinación de estas). Típicamente son equipos homogéneos, y dedicados a la tarea en cuestión.

Mallas (Grids) Computadoras distribuídas geográficamente y conectadas a través de una red de comunicaciones. Las computadoras participantes pueden ser heterogéneas (en capacidades y hasta en arquitectura); la comunicación tiene que adecuarse a enlaces de mucha menor velocidad que en el caso de un cluster, e incluso presentar la elasticidad para permitir las conexiones y desconexiones de nodos en el transcurso del cómputo.

Cómputo en la nube Un caso específico de cómputo distribuído con partición de recursos (al estilo del modelo cliente-servidor); este modelo de servicio está fuertemente orientado a la *tercerización* de servicios específicos. A diferencia del modelo cliente-servidor tradicional, en un entorno de cómputo en la nube lo más común es que tanto el cliente como el servidor sean procesos que van integrando la información, posiblemente por muchos pasos, y que sólo eventualmente llegarán a un usuario final. La implementación de cada uno de los servicios empleados deja de ser relevante, para volverse un servicio *opaco*. Algunos conceptos relacionados son:

Servicios Web Mecanismo de descripción de funcionalidad, así como de solicitud y recepción de resultados, basado en el estándar HTTP y contenido XML.

Software como servicio El proveedor ofrece una *aplicación completa y cerrada* sobre la red, exponiendo únicamente su interfaz (API) de consultas.

Plataforma como servicio El proveedor ofrece la *abstracción* de un entorno específico de desarrollo de modo que un equipo de programadores pueda *desplegar* una aplicación desarrollada sobre dicha plataforma tecnológica. Puede ser visto como un conjunto de piezas de infraestructura sobre un servidor administrado centralmente.

Infraestructura como servicio El proveedor ofrece computadoras completas (en hardware real o máquinas virtuales); la principal ventaja de esta modalidad es que los usuarios, si bien retienen la capacidad plena de administración sobre sus *granjas*, tienen mucho mayor flexibilidad para aumentar o reducir el consumo de recursos (y por tanto, el pago) según la demanda que alcancen.

El tema del cómputo en la nube va muy de la mano de la virtualización, que se abordará en el apéndice B.

3.9.3 Amdahl y Gustafson: ¿qué esperar del paralelismo?

Al programar una aplicación de forma que aproveche al paralelismo (esto es, diseñarla para que realice en distintos procesadores o nodos sus *porciones paralelizables*) ¿cuál es el incremento al rendimiento que se puede esperar?

En 1967, Gene Amdahl presentó un artículo¹⁰ en que indica los límites máximos en que resultará la programación multiprocesada ante determinado programa: parte de la observación de que aproximadamente 40% del tiempo de ejecución de un programa se dedicaba a administración y mantenimiento de los datos, esto es, a tareas secuenciales.

Si únicamente el 60% del tiempo de procesamiento es susceptible, pues, de ser paralelizado, el rendimiento general del sistema no se incrementará en una proporción directa con el número de procesadores, sino que debe sumársele la porción estrictamente secuencial. Puesto en términos más formales: la ganancia en la velocidad de ejecución de un programa al ejecutarse en un entorno paralelo estará limitado por el tiempo requerido por su fracción secuencial. Esto significa que, si $T(1)$ representa al tiempo de ejecución del programa con un sólo procesador y $T(P)$ al tiempo de ejecución con P procesadores, y si t_s es el tiempo requerido para ejecutar la porción secuencial del programa, y $t_p(P)$ el tiempo que requiere la ejecución de la porción paralelizable, repartida entre P procesadores, se puede hablar de una ganancia g en términos de:

$$g = \frac{T(1)}{T(P)} = \frac{t_s + t_p(1)}{t_s + \frac{t_p(1)}{P}}$$

Esta observación, conocida como la *Ley de Amdahl*, llevó a que por varias décadas el cómputo paralelo fuera relegado al cómputo de propósito específico, para necesidades muy focalizadas en soluciones altamente paralelizables, como el cómputo científico.

En términos del ejemplo presentado en la figura 3.9, se ve un programa que, ejecutado secuencialmente, resulta en $T = 500$. Este programa está dividido en tres secciones secuenciales, de $t = 100$ cada una, y dos secciones paralelizables, totalizando $t = 100$ cada una, como se puede ver al representar una ejecución estrictamente secuencial (??).

Al agregar un procesador adicional (??), se obtiene una ganancia de 1.25x — la ejecución se completa en $T = 400$ (con $g = 1.25$). Las secciones paralelizables sólo toman un tiempo *externo* de 50 cada una, dado que la carga fue repartida entre dos unidades de ejecución. Al ejecutar con cuatro procesadores (??), si bien se sigue notando mejoría, esta apenas lleva a $T = 350$, con $g = 1.4$.

Si el código fuera infinitamente paralelizable, y se ejecutase este programa en una computadora con un número infinito de procesadores, este programa no podría ejecutarse en menos de $T = 300$, lo cual presentaría una ganancia de apenas $g = 1.66$. Esto es, al agregar procesadores adicionales, rápidamente se llegaría a un crecimiento asintótico — el comportamiento descrito por la Ley de Amdahl es frecuentemente visto como una demostración de que el desarrollo de sistemas masivamente paralelos presenta *rendimientos decrecientes*.

Si bien el ejemplo que se acaba de presentar resulta poco optimizado, con sólo un 40% de código paralelizable, se puede ver en la gráfica 3.10 que el panorama no cambia tan fuertemente con cargas típicas. Un programa relativamente bien optimizado, con 80% de ejecución paralela, presenta un crecimiento atractivo en la región de hasta 20 procesadores, y se estaciona apenas arriba de una ganancia de 4.5 a partir de los 40.¹¹ Incluso el hipotético 95% llega a un tope en su crecimiento, imposibilitado de alcanzar una ganancia superior a 20.

Dado que el factor económico resulta muy importante para construir computadoras masivamente paralelas,¹² y que se ve claramente que el poder adicional que da cada procesador es cada vez menor, la Ley de Amdahl resultó (como ya se mencionó) en

¹⁰[Validity of the Single Processor Approach to Achieving Large Scale Computing Capabilities](#), Amdahl, 1967

¹¹De un máximo de 5 al que puede alcanzar con un número infinito de procesadores

¹²Dado que los componentes más caros son necesariamente los procesadores

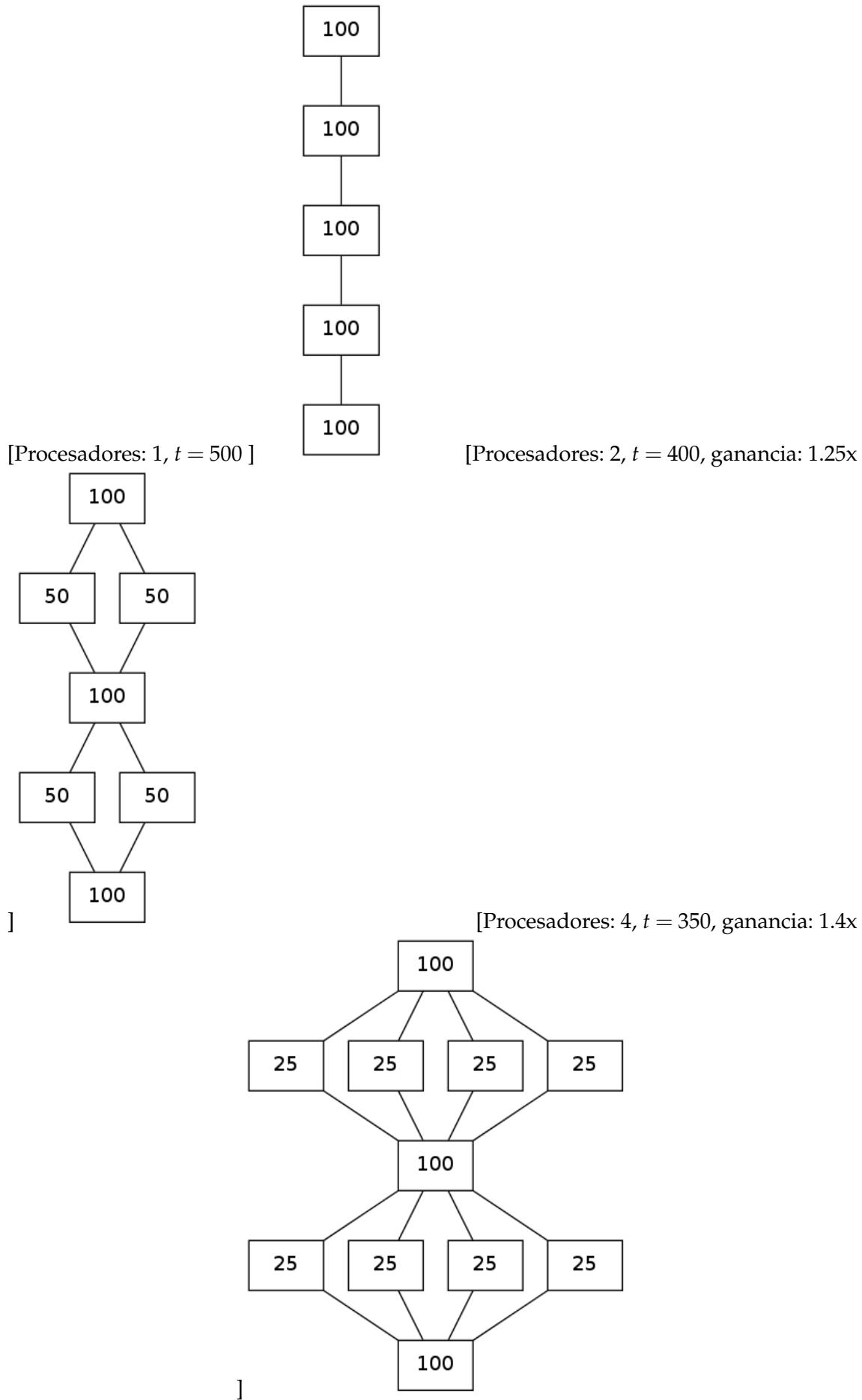


Figura 3.9: Ley de Amdahl: ejecución de un programa con 500 unidades de tiempo total de trabajo con uno, dos y cuatro procesadores.

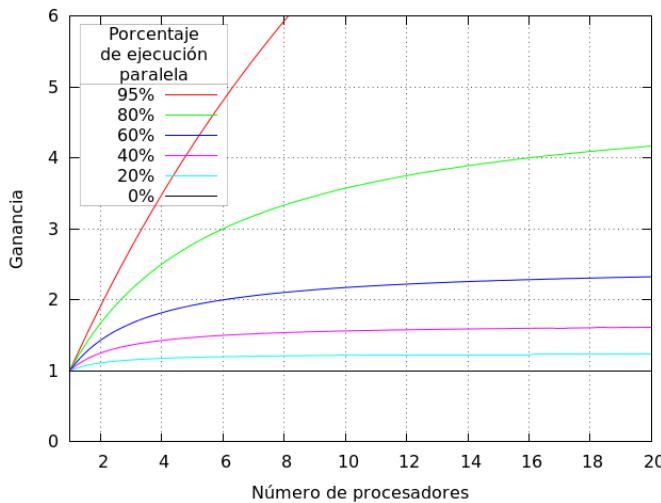


Figura 3.10: Ganancia máxima al paralelizar un programa, según la Ley de Amdahl

varias décadas de mucho mayor atención a la miniaturización y aumento de reloj, y no al multiprocesamiento.

Fue hasta 1988 que John Gustafson publicó una observación a esta ley¹³ que, si bien no la invalida, permite verla bajo una luz completamente diferente y mucho más optimista. Gustafson publicó este artículo corto tras haber obtenido ganancias superiores a 1020 en una supercomputadora con 1024 procesadores — un incremento casi perfectamente lineal al número de procesadores. Sí, respondiendo a una carga altamente optimizada, pero no por eso menos digna de análisis.

El argumento de Gustafson es que al aumentar el número de procesadores, típicamente se verá una modificación *al problema mismo*. Citando de su artículo (traducción propia),

(...)Asumen implícitamente que el tiempo que se ejecuta en paralelo es independiente del número de procesadores, *lo cual virtualmente nunca ocurre de este modo*. Uno no toma un problema de tamaño fijo para ejecutarlo en varios procesadores como no sea para hacer un ejercicio académico; en la práctica, *el tamaño del problema crece con el número de procesadores*. Al obtener procesadores más poderosos, el problema generalmente se expande para aprovechar las facilidades disponibles. Los usuarios tienen control sobre cosas como la resolución de la malla, el número de pasos, la complejidad de los operadores y otros parámetros que usualmente se ajustan para permitir que el programa se ejecute en el tiempo deseado. Por tanto, podría ser más realista que el *tiempo de ejecución*, no el *tamaño del problema*, es constante.

Lo que escribe Gustafson se traduce a que es posible obtener la eficiencia deseada de cualquier cantidad de procesadores *aumentando suficientemente el tamaño del problema*. Al enfrentarse explícitamente con el *bloqueo mental* contra el paralelismo masivo que nació de esta lectura errónea de lo comentado por Amdahl, su artículo sencillo y de apenas más de una cuartilla de extensión cambió la percepción acerca de la utilidad misma del paralelismo masivo.

3.10 Otros recursos

- *Can programming be liberated from the von Neumann style?: a functional style and its algebra of programs*

¹³[Reevaluating Amdahl's Law](#), John L. Gustafson, Communications of the ACM, vol. 31, mayo de 1988

- <https://dl.acm.org/citation.cfm?doid=359576.359579>
John Backus, 1978; Communications of the ACM
- *Cramming more components onto integrated circuits*
<http://cs.utexas.edu/~fussell/courses/cs352h/papers/moore.pdf>
Gordon E. Moore, 1965; Proceedings of the IEEE
 - *An Introduction to the Intel® QuickPath Interconnect*
<http://www.intel.com/content/www/us/en/io/quickpath-technology/quick-path-interconnect.html>
Intel, 2009 (Document Number: 320412-001US)
 - *Intel® Desktop Board D875PBZ Technical Product Specification*
http://downloadmirror.intel.com/15199/eng/D875PBZ_TechProdSpec.pdf
(Intel, 2003)

4 — Administración de procesos

4.1 Concepto y estados de un proceso

En un sistema multiprogramado o de tiempo compartido, un *proceso* es la imagen en memoria de un programa, junto con la información relacionada con el estado de su ejecución.

Un programa es una *entidad pasiva*, una lista de instrucciones; un proceso es una *entidad activa*, que —empleando al programa— define la actuación que tendrá el sistema.

En contraposición con *proceso*, en un sistema por lotes se habla de *tareas*. Una tarea requiere mucha menos estructura, típicamente basta con guardar la información relacionada con la *contabilidad* de los recursos empleados. Una tarea no es interrumpida en el transcurso de su ejecución. Ahora bien, esta distinción no es completamente objetiva — y se pueden encontrar muchos textos que emplean indistintamente una u otra nomenclatura.

Si bien el sistema brinda la *ilusión* de que muchos procesos se están ejecutando al mismo tiempo, la mayor parte de ellos típicamente está esperando para continuar su ejecución — en un momento determinado sólo puede estar ejecutando sus instrucciones un número de procesos igual o menor al número de procesadores que tenga el sistema.

En este capítulo se desarrollan los conceptos relacionados con procesos, hilos, concurrencia y sincronización — Las técnicas y algoritmos que emplea el sistema operativo para determinar cómo y en qué orden hacer los cambios de proceso que nos brindan la ilusión de simultaneidad se abordarán en el capítulo 5.

4.1.1 Estados de un proceso

Un proceso, a lo largo de su vida, alterna entre diferentes *estados* de ejecución. Estos son:

Nuevo Se solicitó al sistema operativo la creación de un proceso, y sus recursos y estructuras están siendo creadas.

Listo Está listo para iniciar o continuar su ejecución pero el sistema no le ha asignado un procesador.

En ejecución El proceso está siendo ejecutado en este momento. Sus instrucciones están siendo procesadas en algún procesador.

Bloqueado En espera de algún evento para poder continuar su ejecución (aun si hubiera un procesador disponible, no podría avanzar).

Zombie El proceso ha finalizado su ejecución, pero el sistema operativo debe realizar ciertas operaciones de limpieza para poder eliminarlo de la lista.¹

Terminado El proceso terminó de ejecutarse; sus estructuras están a la espera de ser *limpiadas* por el sistema operativo

¹Estas operaciones pueden incluir notificar al proceso padre, cerrar las conexiones de red que tenía activas, liberar memoria, etc.

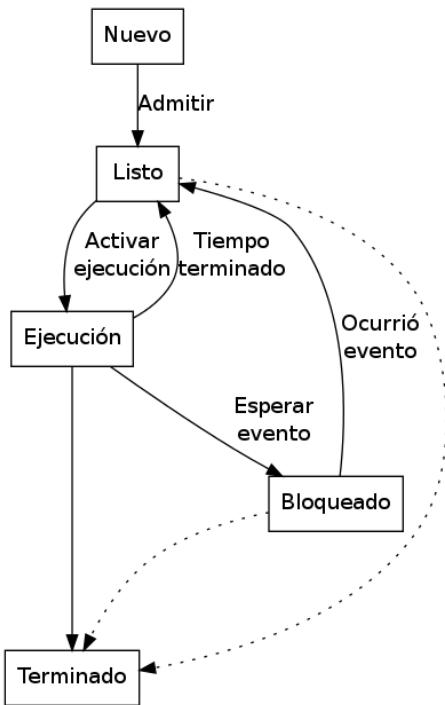


Figura 4.1: Diagrama de transición entre los estados de un proceso

4.1.2 Información asociada a un proceso

La información que debe manipular el sistema operativo relativa a cada uno de los procesos actuales se suele almacenar en una estructura llamada *bloque de control de proceso* (*PCB - Process Control Block*). El PCB incluye campos como:

Estado del proceso El estado actual del proceso.

Contador de programa Cuál es la siguiente instrucción a ser ejecutada por el proceso.

Registros del CPU La información específica del estado del CPU mientras el proceso está en ejecución (debe ser respaldada y restaurada cuando se registra un cambio de estado).

Información de planificación (scheduling) La prioridad del proceso, la *cola* en que está agendado, y demás información que puede ayudar al sistema operativo a planificar los procesos; se profundizará el tema en el capítulo 5.

Información de administración de memoria La información de mapeo de memoria (páginas o segmentos, dependiendo del sistema operativo), incluyendo la pila (*stack*) de llamadas. Se abordará el tema en el capítulo 6.

Información de contabilidad Información de la utilización de recursos que ha tenido este proceso — Puede incluir el tiempo total empleado y otros (*de usuario*, cuando el procesador va avanzando sobre las instrucciones del programa propiamente, *de sistema* cuando el sistema operativo está atendiendo las solicitudes del proceso), uso acumulado de memoria y dispositivos, etc.

Estado de E/S Listado de dispositivos y archivos asignados que el proceso tiene *abiertos* en un momento dado.

4.2 Procesos e hilos

Como se vio, la cantidad de información que el sistema operativo debe manejar acerca de cada proceso es bastante significativa. Si cada vez que el *planificador* elige qué proceso pasar de *Listo* a *En ejecución* debe considerar buena parte de dicha información, la simple

transferencia de todo esto entre la memoria y el procesador podría llevar a un desperdicio *burocrático*² de recursos. Una respuesta a esta problemática fue la de utilizar los *hilos de ejecución*, a veces conocidos como *procesos ligeros* (*LWP, Lightweight processes*).

Cuando se consideran procesos basados en un modelo de hilos, se puede proyectar en sentido inverso que todo proceso es como un sólo hilo de ejecución. Un sistema operativo que no ofreciera soporte expreso a los hilos los planificaría exactamente del mismo modo.

Pero visto desde la perspectiva del proceso hay una gran diferencia: si bien el sistema operativo se encarga de que cada proceso tenga una visión de virtual exclusividad sobre la computadora, todos los hilos de un proceso comparten un sólo espacio de direccionamiento en memoria y los archivos y dispositivos abiertos. Cada uno de los hilos se ejecuta de forma (aparentemente) secuencial y maneja su propio contador de programa y pila (y algunas estructuras adicionales, aunque mucho más ligeras que el PCB).

4.2.1 Los hilos y el sistema operativo

La programación basada en hilos puede hacerse completamente y de forma transparente en espacio de usuario (sin involucrar al sistema operativo). Estos hilos se llaman *hilos de usuario* (*user threads*), y muchos lenguajes de programación los denominan *hilos verdes* (*green threads*). Un caso de uso interesante es en sistemas operativos mínimos (p. ej. para dispositivos embebidos) capaces de ejecutar una máquina virtual de alguno de esos lenguajes: si bien el sistema operativo no maneja multiprocesamiento, a través de los hilos de usuario se crean procesos con multitarea interna.

Los procesos que implementan hilos ganan un poco en el rendimiento gracias a no tener que reemplazar al PCB activo cuando intercalan la ejecución de sus diferentes hilos; pero además de esto, ganan mucho más por la ventaja de compartir espacio de memoria sin tener que establecerlo explícitamente a través de *mecanismos de comunicación entre procesos* (*IPC — Inter Process Communications*). Dependiendo de la plataforma, a veces los hilos de usuario inclusive utilizan multitarea cooperativa para pasar el control dentro de un mismo proceso. Cualquier llamada al sistema *bloqueante* (como obtener datos de un archivo para utilizarlos inmediatamente) interrumpirá la ejecución de *todos* los hilos de ese proceso, dado que el control de ejecución es entregado al sistema operativo quien en este caso no conoce nada sobre los hilos.

Continuando con el desarrollo histórico de este mecanismo, el siguiente paso fue la creación de hilos *informando* al sistema operativo, típicamente denominados *hilos de kernel* (*kernel threads*). Esto se hace a través de bibliotecas de sistema que los implementan de forma estándar para los diferentes sistemas operativos o arquitecturas (p. ej. `pthreads` para POSIX o `Win32_Thread` para Windows). Estas bibliotecas aprovechan la comunicación con el sistema operativo tanto para solicitudes de recursos (p. ej. un proceso basado en hilos puede beneficiarse de una ejecución verdaderamente paralela en sistemas multiprocesador) como para una gestión de recursos más comparable con una situación de multiproceso estándar.

4.2.2 Patrones de trabajo con hilos

Hay tres patrones en los que caen generalmente los modelos de hilos; se puede emplear más de uno de estos patrones en diferentes áreas de nuestra aplicación, e incluso se pueden *anidar* (esto es, se podría tener una *línea de ensamblado* dentro de la cual uno de los pasos sea un *equipo de trabajo*):

Jefe / trabajador Un hilo tiene una tarea distinta de todos los demás: el hilo *jefe* genera o recopila tareas para realizar, las separa y se las entrega a los hilos *trabajadores*.

²Entendiendo *burocrático* como el tiempo que se pierde en asuntos administrativos. Recordar que el tiempo que consume el sistema operativo en administración es tiempo perdido para el uso real, productivo del equipo.

Este modelo es el más común para procesos que implementan servidores (es el modelo clásico del servidor Web *Apache*) y para aplicaciones gráficas (GUIs), en que hay una porción del programa (el hilo jefe) esperando a que ocurran eventos externos. El jefe realiza poco trabajo, se limita a *invocar* a los trabajadores para que hagan el trabajo *de verdad*; como mucho, puede llevar contabilidad de los trabajos realizados.

Típicamente, los hilos trabajadores realizan su operación, posiblemente notifican al jefe de su trabajo, y *finalizan* su ejecución.

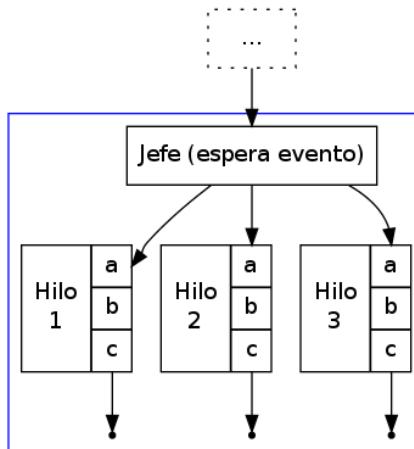


Figura 4.2: Patrón de hilos jefe/trabajador

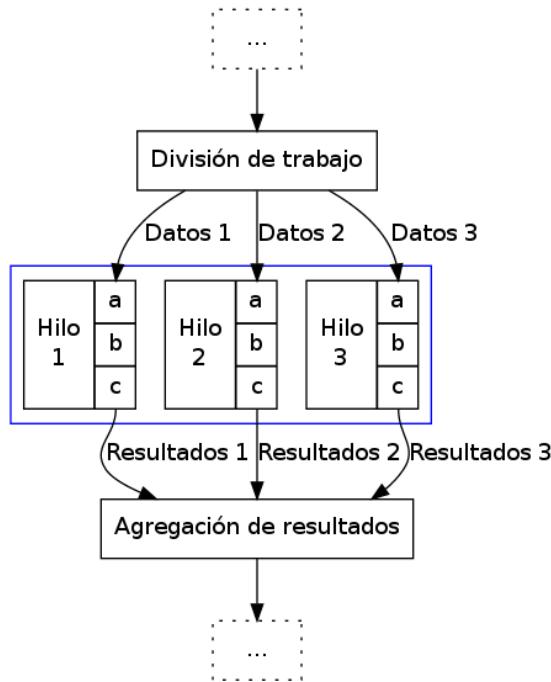
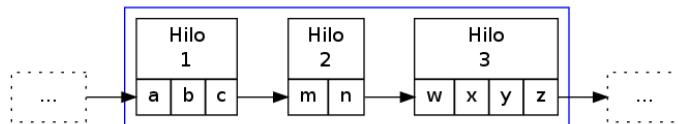
Equipo de trabajo al iniciar la porción multihilos del proceso, se crean muchos hilos idénticos, que realizarán las mismas tareas sobre diferentes datos. Este modelo es muy frecuentemente utilizado para cálculos matemáticos (p. ej.: criptografía, render, álgebra lineal). Puede combinarse con un estilo jefe/trabajador para irle dando al usuario una previsualización del resultado de su cálculo, dado que éste se irá ensamblando progresivamente, pedazo por pedazo.

Su principal diferencia con el patrón *jefe/trabajador* consiste en que el trabajo a realizar por cada uno de los hilos se plantea desde principio, esto es, el paso de *división de trabajo* no es un hilo más, sino que prepara los datos para que éstos sean lanzados en *paralelo*. Estos datos no son resultado de *eventos independientes* (como en el caso anterior), sino partes de un sólo cálculo. Por consecuencia, resulta natural que en este modelo los resultados generados por los diferentes hilos son *agregados* o *totalizados* al terminar su procesamiento. Los hilos no *terminan*, sino que *son sincronizados* y luego continúan la ejecución lineal.

Línea de ensamblado si una tarea larga puede dividirse en pasos sobre bloques de la información total a procesar, cada hilo puede enfocarse a hacer sólo un paso y pasarle los datos a otro hilo conforme vaya terminando. Una de las principales ventajas de este modelo es que nos ayuda a mantener rutinas simples de comprender, y permite que el procesamiento de datos continúe incluso si parte del programa está bloqueado esperando E/S.

Un punto importante a tener en cuenta en una línea de ensamblado es que, si bien los hilos trabajan de forma secuencial, pueden estar ejecutándose paralelamente sobre bloques consecutivos de información, eventos, etc.

Este patrón es claramente distinto de los dos anteriormente presentados; si bien en los anteriores los diferentes hilos (a excepción del hilo *jefe*) eran casi siempre idénticos -aunque operando sobre distintos conjuntos de datos-, en este caso son todos completamente distintos.

Figura 4.3: Patrón de hilos *Equipo de trabajo*Figura 4.4: Patrón de hilos *Línea de ensamblado*

4.3 Concurrencia

4.3.1 Introducción

Desde un punto de vista formal, la *concurrencia* no se refiere a dos o más eventos que ocurren a la vez sino a dos o más eventos cuyo orden es *no determinista*, esto es, eventos acerca de los cuales *no se puede predecir el orden relativo en que ocurrirán*. Si bien dos procesos (o también dos hilos) completamente independientes entre sí ejecutándose simultáneamente son dos procesos concurrentes, la concurrencia se ocupa principalmente de procesos cuya ejecución está vinculada de alguna manera (p. ej.: dos procesos que comparten cierta información o que dependen uno del otro).

Aunque una de las tareas principales de los sistemas operativos es dar a cada proceso la ilusión de que se está ejecutando en una computadora dedicada, de modo que el programador no tenga que pensar en la competencia por recursos, a veces un programa requiere interactuar con otros: parte del procesamiento puede depender de datos obtenidos en fuentes externas, y la cooperación con hilos o procesos externos es fundamental.

Se verá que pueden aparecer muchos problemas cuando se estudia la interacción entre hilos del mismo proceso, la sincronización entre distintos procesos, la asignación de recursos por parte del sistema operativo a procesos simultáneos, o incluso cuando interactúan usuarios de diferentes computadoras de una red — se presentarán distintos conceptos relacionados con la concurrencia utilizando uno de esos escenarios, pero muchos de esos conceptos en realidad son independientes del escenario: más bien nos ocupa la relación entre procesos que deben compartir recursos o deben sincronizar sus tareas.

Para presentar los problemas y conceptos relacionados con la concurrencia suelen

utilizarse algunos problemas clásicos, que presentan casos particulares muy simplificados, y puede encontrárseles relación con distintas cuestiones que un programador enfrentará en la vida real. Cada ejemplo presenta uno o más conceptos. Se recomienda comprender bien el ejemplo, el problema y la solución y desmenuzar buscando los casos límite como ejercicio antes de pasar al siguiente caso. También podría ser útil imaginar en qué circunstancia un sistema operativo se encontraría en una situación similar.

Para cada problema se mostrará *una forma* de resolverlo aunque en general hay más de una solución válida. Algunos de estos problemas fueron originalmente planteados como justificación del desarrollo de las estructuras de control presentadas o de nuevos paradigmas de concurrencia y muchos son aun objeto de debate e investigación.

Para profundizar más en este tema, se recomienda el libro «[The little book of semaphores](#)» de Allen Downey (2008). En este libro (de libre descarga) se encuentran muchos ejemplos que ilustran el uso de semáforos no sólo para resolver problemas de sincronización, sino como un mecanismo simple de comunicación entre procesos. También se desarrollan distintas soluciones para los problemas clásicos (y no tan clásicos).

4.3.2 Problema: el jardín ornamental

Descripción del problema

Un gran jardín ornamental se abre al público para que todos puedan apreciar sus fantásticas rosas, arbustos y plantas acuáticas. Por supuesto, se cobra una modesta suma de dinero a la entrada para lo cual se colocan dos torniquetes, uno en cada una de sus dos entradas. Se desea conocer cuánta gente ha ingresado al jardín así que se instala una computadora conectada a ambos torniquetes: cada torniquete envía una señal cuando una persona ingresa al jardín. Se realiza un modelo simplificado de la situación, así que no se estudiarán los detalles del hardware utilizado. Aquí es importante notar que los dos torniquetes son objetos que existen y se comportan en paralelo e independientemente: los eventos que generan no tienen un orden predecible. Es decir, que cuando se escriba el software no se sabe en qué momento llegará cada visitante ni qué torniquete utilizará.

Se simulará un experimento en el que 20 visitantes ingresan por cada torniquete. Al final de la simulación deberá haber 40 visitantes contados. Una implementación tentativa podría ser la siguiente:³

```
int cuenta;

proceso torniquete1() {
    int i;
    for(i=0;i<20;i++) {
        cuenta = cuenta + 1;
    }
}

proceso torniquete2() {
    int i;
    for(i=0;i<20;i++) {
        cuenta = cuenta + 1;
    }
}

main() {
```

³Se utiliza una versión ficticia del lenguaje C para el ejemplo, evitando entrar en los detalles de sintaxis de un lenguaje concurrente.

```

cuenta = 0;
/* Lanzar ambos procesos concurrentemente*/
concurrentemente { //
    torniquete1();
    torniquete2();
}
/* Esperar a que ambos finalicen */
esperar(torniquete1);
esperar(torniquete2);
printf("Cuenta: %d\n", cuenta);
}

```

Como se ve el problema es muy sencillo. Sin embargo, al intentar ejecutar repetidas veces ese programa muy de vez en cuando el resultado no tiene el valor 40. Si se modifica el programa para utilizar un solo torniquete, cuenta siempre tiene el valor correcto (20).

¿Qué es lo que está ocurriendo? La mayor parte de los lenguajes de programación convierten cada instrucción en una serie más o menos larga de operaciones de máquina (instrucciones ensamblador). De modo, que una instrucción aparentemente simple como cuenta = cuenta + 1 habitualmente implica varias operaciones de más bajo nivel (las instrucciones de ejemplo corresponden a arquitecturas Intel x86):

LEER Leer cuenta desde la memoria (p. ej. mov \$cuenta, %rax).

INC Incrementar el registro (p. ej. add \$1, %rax).

GUARDAR Guardar el resultado nuevamente en memoria (p. ej. mov %rax, \$cuenta).

En un sistema operativo multitarea cuando un proceso agota su porción de tiempo de procesador (quantum) o detiene su ejecución por otra razón, los valores almacenados en registros se preservan (junto con la información sobre el proceso) para poder restaurarlo cuando la ejecución continúa (de esta forma se provee la ilusión de la multitarea en sistemas de un solo núcleo). Así, en el problema del Jardín Ornamental cada torniquete tiene su propia copia de los valores en los registros. Sin embargo, se supone que el resto de la memoria es compartida (en particular, se utiliza ese hecho para llevar la cuenta de personas que ingresan).

Si se considera lo que ocurre cuando dos procesos (p. ej. torniquete1 y torniquete2) ejecutan la instrucción cuenta = cuenta + 1 en un equipo con un solo procesador, puede darse la siguiente secuencia de eventos. Se considera que cuenta está inicialmente en 0.

1. cuenta = 0
2. torniquete1: LEER (resultado: rax de $p_1 = 0$, cuenta = 0)
3. torniquete1: INC (resultado: rax de $p_1 = 1$, cuenta = 0)
4. torniquete1: GUARDAR (resultado: rax de $p_1 = 1$, cuenta = 1)
5. El sistema operativo decide cambiar de tarea, suspende torniquete1 y continúa con torniquete2.
6. torniquete2: LEER (resultado: rax de $p_2 = 1$, cuenta = 1)
7. torniquete2: INC (resultado: rax de $p_2 = 2$, cuenta = 1)
8. torniquete2: GUARDAR (resultado: rax de $p_2 = 2$, cuenta = 2)

Se puede ver que ambos procesos realizaron sus instrucciones para incrementar el contador en 1 y el resultado final fue que la cuenta se incrementó en dos unidades.

Pero, también puede darse la siguiente secuencia de eventos durante la ejecución de estas instrucciones:

1. cuenta = 0
2. torniquete1: LEER (resultado: rax de $p_1 = 0$, cuenta = 0)
3. torniquete1: INC (resultado: rax de $p_1 = 1$, cuenta = 0)

4. El sistema operativo decide cambiar de tarea, suspende torniquete1 y continua con torniquete2.
5. torniquete2: LEER (resultado: rax de $p_2 = 1$, cuenta = 1)
6. torniquete2: INC (resultado: rax de $p_2 = 2$, cuenta = 1)
7. torniquete2: GUARDAR (resultado: rax de $p_2 = 2$, cuenta = 2)
8. El sistema operativo decide cambiar de tarea, suspende torniquete2 y continua con torniquete1.
9. torniquete1: GUARDAR (resultado: rax de $p_1 = 1$, cuenta = 1)

Nuevamente ambos procesos ejecutaron sus instrucciones para incrementar en 1 el contador. Sin embargo, ¡en este caso cuenta tiene el valor 1!. A este problema también se lo conoce como *problema de las actualizaciones múltiples*.

Esto parece muy específico Si bien este análisis parece muy específico es fácil ver que la misma circunstancia podría darse en un sistema de reserva de vuelos (p. ej.: puede de que dos operadores vean un asiento vacío en su copia local de los asientos y ambos marquen el mismo asiento como ocupado) o con dos procesos que decidan cambiar simultáneamente datos en un archivo. Aquí las operaciones ya no son necesariamente operaciones internas de la máquina.

¿Pero no es muy poco probable? Por otro lado, uno podría pensar (con cierta cuota de razón) que la secuencia de eventos propuesta es muy poco probable: usualmente un sistema operativo ejecuta miles de instrucciones antes de cambiar de un proceso a otro. De hecho, en la práctica este problema es muy frecuentemente ignorado y los programas funcionan muy bien la mayoría de las veces. Esto permite ver una característica importante de los programas concurrentes: es muy usual que un programa funcione perfectamente la mayor parte del tiempo, pero de vez en cuando puede fallar. Subsecuentes ejecuciones con los mismos argumentos producen nuevamente el resultado correcto. Esto hace que los problemas de concurrencia sean muy difíciles de detectar y más aun de corregir. Es importante (y mucho más efectivo) realizar un buen diseño inicial de un programa concurrente en lugar de intentar arreglarlo cuando se detecta alguna falla. También es interesante notar que dependiendo del sistema, puede ser que alguna de las instrucciones sea muy lenta, en el caso de un sistema de reserva de asientos de aviones, las operaciones pueden durar un tiempo importante (p. ej.: desde que el operador muestra los asientos disponibles hasta que el cliente elige el asiento) haciendo mucho más probable que ocurra una secuencia no deseada.

Vale la pena preocuparse? A modo de ilustración de la gravedad del problema, estos son algunos valores para el resultado final de la variable cuenta cuando se ejecuta el programa anterior en Pascal-FC⁴: 25 29 31 20 21 26 27 18 31 35. Notesé que incluso uno de los valores es menor que 20 (que es lo mínimo que cuenta cada torniquete). Es un ejercicio interesante pensar qué secuencia de eventos podría producir tal valor y cuál es el mínimo valor posible.

Pero tenemos muchos núcleos Otra cuestión que puede parecer artificiosa es que en el ejemplo hay un solo procesador o núcleo. Sin embargo, tener más de un procesador no sólo no soluciona el problema sino que lo empeora: ahora las operaciones de lectura o escritura pueden ejecutarse directamente en paralelo y aparecen nuevos problemas de coherencia de caché. En la siguiente discusión muchas veces se presupone que hay un solo procesador, sin que eso invalide la discusión para equipos multiprocesadores.

⁴<http://www-users.cs.york.ac.uk/burns/pf.html>

Algunos conceptos de concurrencia

Antes de abordar posibles soluciones al problema presentado, se presentan las definiciones de algunos conceptos importantes.

Operación atómica Operación que requiere la garantía de que se ejecutará como una sola unidad de ejecución, o fallará completamente, sin resultados o estados parciales observables por otros procesos o el entorno. Esto no necesariamente implica que el sistema no retirará el flujo de ejecución en medio de la operación, sino que *el efecto de que se le retire el flujo* no llevará a comportamiento inconsistente.

Condición de carrera (En inglés, *Race condition*) Categoría de errores de programación que involucra a dos procesos que fallan al comunicarse su estado mutuo, llevando a resultados inconsistentes. Es uno de los problemas más frecuentes y difíciles de depurar, y ocurre típicamente por no considerar la *no atomicidad* de una operación.

Sección (o región) crítica El área de código que requiere ser protegida de accesos simultáneos, donde se realiza la modificación de datos compartidos.

Recurso compartido Un recurso que puede ser accedido desde más de un proceso. En muchos escenarios esto es un variable en memoria (como *cuenta* en el jardín ornamental), pero podrían ser archivos, periféricos, etc...

Dado que el sistema no tiene forma de saber cuáles instrucciones (o áreas del código) deben funcionar de forma atómica, el programador debe asegurar la atomicidad de forma explícita, mediante la sincronización de los procesos. El sistema no debe permitir la ejecución de parte de esa área en dos procesos de forma simultánea (sólo puede haber un proceso en la sección crítica en un momento dado).

- ¿Y qué tiene que ver esto con el problema del Jardín Ornamental?

En el problema hay claramente un *recurso compartido* que es la *cuenta*, así la sección que modifica la cuenta es una *sección crítica* y la operación *cuenta = cuenta + 1* debe ser una *operación atómica*. La secuencia de eventos que se mostró es una *condición de carrera*: el segundo torniquete presume un estado (*cuenta = 0*) que no es el mismo que conoce el torniquete1 (*cuenta = 1*).

Soluciones posibles (y no tanto)

El planteamiento del problema del *jardín ornamental* busca llevar al lector a ir encontrando, a través de sucesivos refinamientos, los mecanismos principales que se emplean para resolver –en general– los problemas que implican el acceso concurrente a una sección crítica. Se presentan a continuación, pues, los sucesivos *intentos*.

Intento 1: No utilizar multitarea En este sencillo ejemplo una posible solución es utilizar una sola entrada (o torniquete). Esto podría ser una solución en tanto que no haya mucha gente que haga cola para entrar. Sin embargo, en un sistema análogo de reserva de pasajes aéreos no parece tener mucho sentido que todos los pasajeros deban ir a Japón a sacar su pasaje. Por otro lado, ya deberían ser claras las ventajas de la multitarea y el poseer distintos núcleos.

Intento 2: Suspender la multitarea durante la sección crítica Una versión más relajada de la alternativa anterior es suspender la multitarea durante la ejecución de la sección crítica. Así, un torniquete deberá hacer:

```
disable(); /* Suspender temporal las interrupciones */
cuenta = cuenta + 1;
enable(); /* Habilitar nuevamente las interrupciones */
```

Durante el lapso en el que las interrupciones están suspendidas no puede haber un cambio de contexto pues el planificador depende de la interrupción del reloj (salvo que el proceso realice una llamada bloqueante durante la región crítica).

Esta solución puede resultar conveniente para sistemas sencillos, pero en un sistema multiusuario se torna inusuable por varios motivos:

- Permitir que un programa de usuario deshabilite las interrupciones en un sistema operativo de propósito general involucra un gran problema de seguridad: cualquier usuario podría hacer un programa malicioso (o sencillamente erróneo) que deshabilite las interrupciones y suspenda indefinidamente el resto del sistema.
- No funciona para sistemas distribuidos (como el sistema de reserva de pasajes aéreos), ni siquiera para sistemas multinúcleo o multiprocesador, ya que las interrupciones se deshabilitan en un sólo núcleo (si bien también es posible detener a los demás procesadores, representa un costo demasiado alto).
- Expone detalles de hardware y podría provocar mal funcionamiento de algún periférico si el procesamiento de la sección crítica es demasiado largo.
- Intento 3: Utilizar una bandera ::

Utilizar una bandera parece ser una solución muy sencilla: mediante una variable de bandera se indica si hay un proceso en la región crítica:

```
int bandera = 0; /* 0 => región crítica libre, 1 => ocupada */
int cuenta = 0;
/* ... */

/* Torniquete1 */
/* ... */
if (bandera) wait;
/* Aquí bandera=0 */
bandera = 1; /* Inicio de la sección crítica */
cuenta = cuenta + 1;
bandera = 0; /* Fin de la sección crítica */
```

Sin embargo esto no funciona, ahora puede darse la siguiente secuencia de eventos:

1. bandera==0;
2. torniquete2: if (bandera) wait;
3. Nuevo cambio de contexto
4. torniquete1: if (bandera) wait;
5. torniquete1: bandera = 1;
6. torniquete2: bandera = 1;
7. torniquete2: cuenta = cuenta + 1;
8. torniquete1: cuenta = cuenta + 1; /* Ups, no se respetó la región crítica */

Notar que el problema aquí es que la bandera también es un recurso compartido: lo único que ha cambiado es que ahora la sección crítica está en otro lugar. La solución funcionaría si se pudiera garantizar que la secuencia de operaciones se realizara atómicamente:

```
if (bandera) wait;
bandera = 1
```

Intento 4: Manejar la bandera con instrucciones atómicas Algunas arquitecturas de computadoras permiten realizar determinadas operaciones sencillas (como actualizar una bandera) de forma atómica (p. ej.: VAX tiene la instrucción `test_and_set` y el i386 tiene la instrucción `INC`).

Usando esto, la solución es:

```
int bandera; /* 0 => desocupada */

while (++bandera != 1) {
```

```

bandera--; /* Debe generar "INC" */
}
/* Sección crítica */
cuenta = cuenta + 1;

bandera--;

```

Esto funciona correctamente siempre que la operación `++bandera` sea atómica. Sin embargo, hay dos problemas a considerar: un proceso puede permanecer mucho tiempo repitiendo el ciclo:

```

while (++bandera!=1) {
    bandera--;
}

```

De hecho, si el sistema operativo decide darle alta prioridad a este proceso es posible que esté un tiempo infinito en este ciclo, impidiendo que otro proceso decremente la bandera. Y aún cuando el sistema operativo decida cambiar de proceso en el siguiente tic de reloj, es evidente que se podría aprovechar el procesador para hacer algo útil durante ese tiempo y que suspender el proceso de otra manera le da más posibilidad a otros procesos para que cambien la bandera. A esto se lo conoce como *espera activa* o *espera ocupada* (*busy waiting* en inglés) y es una situación que se desea evitar.

El otro problema tiene que ver con el hardware: determinadas arquitecturas no permiten instrucciones que lean y actualicen en una única operación una dirección de memoria (se requiere una operación para leer y otra para escribir). En particular, ninguna arquitectura RISC lo permite (p. ej.: SPARC, RS 6000, ...).

Intento 5: Utilizar turnos Una alternativa para evitar el problema de la actualización múltiple a una bandera es utilizar turnos

```
int turno = 1; /* Inicialmente el turno es del proceso 1 */
```

Ahora el código del proceso 1 contendría algo como:

```

while (turno != 1) {
    esperar(); /* ¿Otro proceso? */
}
/* Sección crítica */
cuenta = cuenta + 1;
turno = 2;

```

Y el del proceso dos:

```

while (turno != 2) {
    esperar();
}
/* Sección crítica */
cuenta = cuenta + 1;
turno = 1;

```

Esto garantiza que no hay dos procesos en sección crítica. Pero nótese que hacer esto equivale a tener un solo torniquete: sólo podrá haber una persona ingresando a la vez... o incluso peor, las personas deberán utilizar alternativamente los torniquetes. Así que si bien esto soluciona el problema de la actualización múltiple en realidad es una solución muy restrictiva: un proceso que no está en la sección crítica puede obligar a que otro proceso espere mucho tiempo para ingresar a la sección crítica. De aquí en más se buscarán soluciones en las que no ocurra esto.

Intento 6: Indicar la intención de entrar a la sección crítica Para paliar los efectos de la solución anterior se puede intentar indicar si el otro proceso también está queriendo entrar en la sección crítica. El código sería:

```
int b1, b2;
/* ... */

/* Proceso 1: */
/* ... */
b1 = 1;
if (b2) {
    esperar();
}
/* Sección crítica */
cuenta = cuenta + 1;
b1 = 0;
/* ... */
/* Proceso 2: */
/* ... */
b2 = 1;
if (b1) {
    esperar();
}
/* Sección crítica */
cuenta = cuenta + 1;
b2 = 0;
/* ... */
```

Nuevamente aquí está garantizado que no puede haber dos procesos en la región crítica, pero este enfoque sufre de un problema grave: ambos procesos pueden bloquearse mutuamente (si el proceso 1 coloca su bandera en 1 y luego se cambia el control al proceso 2 quien también colocará su bandera en 1). ORG-LIST-END-MARKER

Una solución: el Algoritmo de Peterson

La primera solución a esta problema fue propuesta por Dekker en 1957. Sin embargo su explicación es bastante extensa (aunque perfectamente comprensible). Se presentará la solución planteada por Peterson unos cuantos años más tarde: en 1970.

La solución está basada en una combinación de los intentos anteriores: utilizar banderas para indicar qué proceso puede entrar, pero además usa un turno para *desempatar* en caso de que ambos procesos busquen entrar a la vez. En cierto modo es un algoritmo *amable*: si un proceso detecta que el otro proceso fue el primero en actualizar el turno, entonces lo deja pasar:

```
int b2, b2;
int quien;

/* Proceso 1: */
...
b1=1;
quien=2;
if ( b2 && (quien==2) ) {
    esperar();
}
```

```

/* Sección crítica */
cuenta = cuenta + 1;
b1=0;

/* Proceso 2: */
...
b2=1;
quien=1;
if ( b1 && quien==1) {
    esperar();
}
/* Sección crítica */
cuenta = cuenta + 1;
b1=0;

```

Cabe apuntar las siguientes notas sobre la solución de Peterson:

Espera activa La solución presentada mantiene todavía el problema de la espera activa (también llamados *spinlocks*): un proceso puede consumir mucho tiempo de procesador sólo para esperar que otro proceso cambie una bandera, lo cual en un sistema con manejo de *prioridades*, puede resultar dañino para el desempeño global. Una forma de mitigar los efectos es forzar (o sugerir) cambios de contexto en esos puntos a través de una *primitiva* del lenguaje o del sistema operativo (p. ej.: `sleep` o `yield`), pero debe resultar claro que de ninguna forma es una solución general. Por esta razón los sistemas operativos o lenguajes suelen proveer alguna abstracción para soportar explícitamente operaciones atómicas o implementar una solución más elegante al problema. Se verán algunas de esas abstracciones más adelante.

Para mayores detalles acerca de las razones, ventajas y desventajas del uso de spinlocks en sistemas operativos reales, referirse a [Spin Locks & Other Forms of Mutual Exclusion](#) (Theodore P. Baker 2010)

Solución para más procesos El algoritmo de Peterson sirve únicamente cuando hay dos procesos que compiten para acceder a una región crítica. ¿Qué se puede hacer si hay más de dos entradas al jardín, o si hay más de dos puntos de venta de pasajes aereos? La solución a este problema más general fue propuesta por Dijkstra en 1968 y posteriormente Eisenberg y McGuire en 1972 y Lamport en 1974 presentaron distintas soluciones.

La más ampliamente utilizada y sencilla de entender es la propuesta por Lamport, también conocida como el *algoritmo de la panadería* por su semejanza con el sistema de turnos utilizado para atender a los clientes en una panadería.

Solución para equipos multiprocesadores Esta solución (y también la de Lamport y todos los autores mencionadas hasta ahora) falla en equipos multiprocesadores, pues aparecen problemas de coherencia de caché. Se necesitan precauciones especiales en equipos con más de un procesador.

4.3.3 Mecanismos de sincronización

En la presente sección se enumeran los principales mecanismos que pueden emplearse para programar considerando a la concurrencia: Candados, semáforos y variables de condición.

Regiones de exclusión mutua: candados o mutexes

Una de las alternativas que suele ofrecer un lenguaje concurrente o sistema operativo para evitar la *espera activa* a la que obliga el algoritmo de Peterson (o similares) se llama *mutex* o *candado* (*lock*).

La palabra *mutex* nace de la frecuencia con que se habla de las *regiones de exclusión mutua* (en inglés, *mutual exclusion*). Es un mecanismo que asegura que cierta región del código será ejecutada como si fuera atómica.

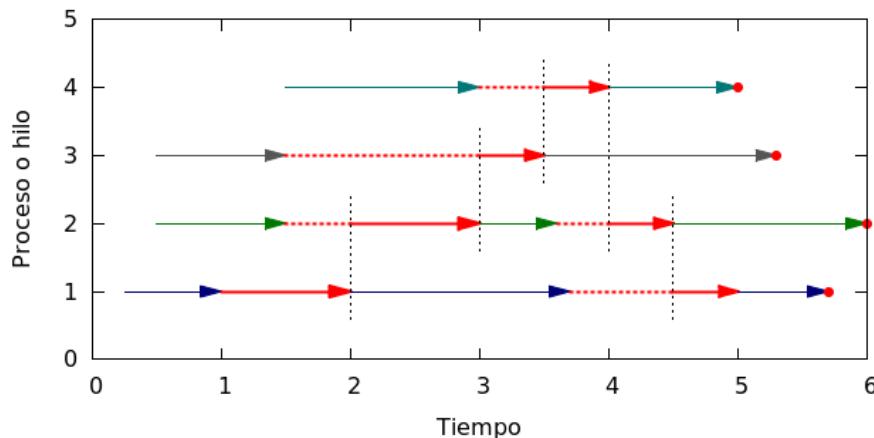


Figura 4.5: Sincronización: La exclusión de las *secciones críticas* entre a varios procesos se protegen por medio de regiones de *exclusión mutua*

Hay que tener en cuenta que un mutex *no implica* que el código no se va a interrumpir mientras se está dentro de esta región — Eso sería muy peligroso, dado que permitiría que el sistema operativo pierda el control del planificador, volviendo (para propósitos prácticos) a un esquema de multitarea cooperativa. El mutex es un *mecanismo de prevención*, que mantiene en espera a cualquier hilo o proceso que quiera entrar a la *sección crítica* protegida por el mutex, reteniéndolo antes de entrar a ésta hasta que el proceso que la está ejecutando salga de ella. Si no hay ningún hilo o proceso en dicha sección crítica (o cuando un hilo sale de ella), uno sólo de los que esperan podrá ingresar.

Como se vio en el ejemplo anterior, para que un mutex sea efectivo tiene que ser implementado a través de una *primitiva* a un nivel inferior⁵, implicando al planificador.

El problema de la actualización múltiple que surge en el caso de la venta de pasajes aereos podría reescribirse de la siguiente manera empleando un mutex:

```
my ($proximo_asiento :shared, $capacidad :shared);
$capacidad = 40;

sub asigna_asiento {
    lock($proximo_asiento);
    if ($proximo_asiento < $capacidad) {
        $asignado = $proximo_asiento;
        $proximo_asiento += 1;
        print "Asiento asignado: $asignado\n";
    } else {
        print "No hay asientos disponibles\n";
        return 1;
    }
    return 0;
}
```

⁵¿Qué significa inferior? Las llamadas de sincronización entre hilos deben implementarse por lo menos a nivel del proceso que los contiene; aquellas que se realizan entre procesos independientes, deben implementarse a nivel del sistema operativo. Debe haber un agente *más abajo* en niveles de abstracción, en control *real* del equipo de cómputo, ofreciendo estas operaciones.

Se debe tener en cuenta que en este caso se utiliza una implementación de hilos, esto hace que la solución sea dependiente del lenguaje específico de implementación, en este caso Perl. Al ser `$proximo_asiento` una variable compartida tiene algunas *propiedades* adicionales, en este caso, la de poder operar como un mutex. La implementación en Perl resulta muy *limpia*, dado que evita el uso de un candado explícito — Se podría leer la línea 5 como *exclusión mutua sobre \$proximo_asiento*.

En la implementación de hilos de Perl, la función `lock()` implementa un mutex delimitado por el *ámbito léxico* de su invocación: el área de exclusión mutua abarca desde la línea 5 en que es invocada hasta la 15 en que termina el bloque en que se invocó.

Un área de exclusión mutua debe:

Ser mínima Debe ser *tan corta como sea posible*, para evitar que otros hilos queden bloqueados fuera del área crítica. Si bien en este ejemplo es demasiado simple, si se hiciera cualquier llamada a otra función (o al sistema) estando dentro de un área de exclusión mutua, se detendría la ejecución de todos los demás hilos por mucho más tiempo del necesario.

Ser completa Se debe analizar bien cuál es el área a proteger y no arriesgarse a proteger de menos. En este ejemplo, se podría haber puesto `lock($asignado)` dentro del `if`, dado que sólo dentro de su evaluación positiva se modifica la variable `$proximo_asiento`. Sin embargo, si la ejecución de un hilo se interrumpiera entre las líneas 7 y 8, la condición del `if` se podría evaluar incorrectamente.

Como comparación, una rutina equivalente en Bash (entre procesos independientes y usando los archivos `/tmp/proximo_asiento` y `/etc/capacidad` como un mecanismo para compartir datos) sería:

```
asigna_asiento() {
    lockfile /tmp/asigna_asiento.lock
    PROX=$(cat /tmp/proximo_asiento || echo 0)
    CAP=$(cat /etc/capacidad || echo 40)
    if [ $PROX -lt $CAP ]
        then
            ASIG=$PROX
            echo $((PROX+1)) > /tmp/proximo_asiento
            echo "Asiento asignado: $ASIG"
        else
            echo "No hay asientos disponibles"
            return 1;
        fi
    rm -f /tmp/asigna_asiento.lock
}
```

Cabe mencionar que `lockfile` no es una función implementada en Bash, sino que envuelve a una *llamada al sistema*. El sistema operativo garantiza que la verificación y creación de este *candado* se efectuará de forma atómica.

Un mutex es, pues, una herramienta muy sencilla, y podría verse como la pieza básica para la sincronización entre procesos. Lo fundamental para emplearlos es identificar las regiones críticas del código, y proteger el acceso *con un mecanismo apto de sincronización*, que garantice atomicidad.

Semáforos

La interfaz ofrecida por los mutexes es muy sencilla, pero no permite resolver algunos problemas de sincronización. Edsger Dijkstra, en 1968, propuso los *semáforos*.⁶

Un semáforo es una variable de tipo entero que tiene definida la siguiente interfaz:

Inicialización Se puede inicializar el semáforo a cualquier valor entero, pero después de esto, su valor no puede ya ser leído. Un semáforo es una *estructura abstracta*, y su valor es tomado como *opaco* (invisible) al programador.

Decrementar Cuando un hilo decrementa el semáforo, si el valor es negativo, el hilo se *bloquea* y no puede continuar hasta que *otro hilo* incremente el semáforo. Según la implementación, esta operación puede denominarse *wait*, *down*, *acquire* o incluso *P* (por ser la inicial de *proberen te verlagen, intentar decrementar* en holandés, del planteamiento original en el artículo de Dijkstra).

Incrementar Cuando un hilo incrementa el semáforo, si hay hilos esperando, uno de ellos es *despertado*. Los nombres que recibe esta operación son *signal*, *up*, *release*, *post* o *V* (de *verhogen, incrementar*).

La interfaz de hilos POSIX (pthreads) presenta esas primitivas con la siguiente definición:

```
int sem_init(sem_t *sem, int pshared, unsigned int value);
int sem_post(sem_t *sem);
int sem_wait(sem_t *sem);
int sem_trywait(sem_t *sem);
```

La variable *pshared* indica si el semáforo puede ser compartido entre procesos o únicamente entre hilos. *sem_trywait* extiende la interfaz sugerida por Dijkstra: verifica si el semáforo puede ser decrementado y, en caso de que no, en vez de bloquearse, indica al proceso que no puede continuar. El proceso debe tener la lógica necesaria para no entrar en las secciones críticas (p. ej., intentar otra estrategia) en ese caso.

sem_trywait se sale de la definición clásica de semáforo, por lo que no se considera en esta sección.

Un semáforo permite la implementación de varios patrones, entre los cuales se mencionarán los siguientes:

Señalar Un hilo debe informar a otro que cierta condición está ya cumplida — Por ejemplo, un hilo prepara una conexión en red mientras que otro calcula lo que tiene que enviar. No se puede arriesgar a comenzar a enviar antes de que la conexión esté lista. Se inicializa el semáforo a 0, y:

```
# Antes de lanzar los hilos
senal = Semaphore(0)

def envia_datos():
    calcula_datos()
    senal.acquire()
    envia_por_red()

def preparaConexion():
    creaConexion()
    senal.release()
```

⁶El símil presentado por Dijkstra no es del semáforo vial, con una luz roja y una luz verde (dicho esquema se asemeja al del *mutex*). La referencia es a la del semáforo de tren, que permite el paso estando *arriba*, e indica espera estando *abajo*.

No importa si `preparaConexion()` termina primero — En el momento en que termine, señalará 1 y `enviaDatos()` podrá proceder.

Rendezvous Así se denomina en francés (y ha sido adoptado al inglés) a quedar en una *cita*. Este patrón busca que dos hilos se esperen mutuamente en cierto punto para continuar en conjunto — Por ejemplo, en una aplicación GUI, un hilo prepara la interfaz gráfica y actualiza sus eventos mientras otro efectúa cálculos para mostrar. Se desea presentar al usuario la simulación desde el principio, así que no debe empezar a calcular antes de que el GUI esté listo, pero preparar los datos del cálculo toma tiempo, y no se quiere esperar doblemente. Para esto, se implementan dos semáforos señalizándose mutuamente:

```
guiListo = Semaphore(0)
calculоЛisto = Semaphore(0)

threading.Thread(target=maneja_gui, args=[]).start()
threading.Thread(target=maneja_calculo, args=[]).start()

def maneja_gui():
    inicializa_gui()
    guiListo.release()
    calculоЛisto.acquire()
    recibe_eventos()

def maneja_calculo():
    inicializa_datos()
    calculоЛisto.release()
    guiListo.acquire()
    procesa_calculo()
```

Mutex El uso de un semáforo inicializado a 1 puede implementar fácilmente un mutex. En Python:

```
mutex = Semaphore(1)
# ...Iniciar el estado y lanzar hilos
mutex.acquire()
# Aquí se está en la región de exclusión mutua
x = x + 1
mutex.release()
# Continúa la ejecución paralela
```

Multiplex Permite la entrada de no más de n procesos a la región crítica. Si se lo ve como una generalización de *Mutex*, basta con inicializar al semáforo al número máximo de procesos deseado.

Su construcción es idéntica a la de un mutex, pero es inicializado al número de procesos que se quiere permitir que ejecuten de forma simultánea.

Torniquete Una construcción que por sí sola no hace mucho, pero resulta útil para patrones posteriores. Esta construcción garantiza que un grupo de hilos o procesos *pasa por un punto determinado* de uno en uno (incluso en un ambiente multiprocesador):

```
torniquete = Semaphore(0)
# ...
if alguna_condicion():
    torniquete.release()
```

```
# (...)  
torniquete.acquire()  
torniquete.release()
```

En este caso, se ve primero una *señalización* que hace que todos los procesos esperen frente al torniquete hasta que alguno marque que `alguna_condicion()` se ha cumplido y libere el paso. Posteriormente, los procesos que esperan pasarán ordenadamente por el torniquete.

El torniquete por sí sólo no es tan útil, pero su función se hará clara a continuación.

Apagador Cuando se tiene una situación de *exclusión categórica* (basada en categorías y no en procesos individuales — Varios procesos de la misma categoría pueden entrar a la sección crítica, pero procesos de dos categorías distintas deben tenerlo prohibido), un *apagador* permite evitar la inanición de una de las categorías ante un flujo constante de procesos de la otra.

El apagador usa, como uno de sus componentes, a un torniquete. Para ver una implementación ejemplo de un apagador, referirse a la solución presentado a continuación para el problema lectores-escritores.

Barrera Una barrera es una generalización de *rendezvous* que permite la sincronización entre varios hilos (no sólo dos), y no requiere que el rol de cada uno de los hilos sea distinto.

Esta construcción busca que ninguno de los hilos continúe ejecutando hasta que todos hayan llegado a un punto dado.

Para implementar una barrera, es necesario que ésta guarde algo de información adicional además del semáforo, particularmente, el número de hilos que se han lanzado (para esperarlos a todos). Esta será una variable compartida y, por tanto, requiere de un mutex. La inicialización (que se ejecuta antes de iniciar los hilos) será:

```
require random  
n = random.randint(1,10) # Número de hilos  
cuenta = 0  
mutex = Semaphore(1)  
barrera = Semaphore(0)
```

Ahora, suponiendo que todos los hilos tienen que realizar, por separado, la inicialización de su estado, y ninguno de ellos debe comenzar el procesamiento hasta que todos hayan efectuado su inicialización:

```
inicializa_estado()
```

```
mutex.acquire()  
cuenta = cuenta + 1  
mutex.release()
```

```
if cuenta == n:  
    barrera.release()
```

```
barrera.acquire()  
barrera.release()
```

```
procesamiento()
```

Las barreras son una construcción suficientemente útil como para que sea común encontrarlas “prefabricadas”. En los hilos POSIX (pthreads), por ejemplo, la interfaz básica es:

```

int pthread_barrier_init(pthread_barrier_t *barrier,
                        const pthread_barrierattr_t *restrict attr,
                        unsigned count);
int pthread_barrier_wait(pthread_barrier_t *barrier);
int pthread_barrier_destroy(pthread_barrier_t *barrier);

```

Cola Se emplea una cola cuando se tienen dos *clases de hilos* que deben proceder en pares. Este patrón es a veces referido como *baile de salón*: para que una pareja baile, hace falta que haya un *líder* y un *seguidor*. Cuando llega una persona al salón, verifica si hay uno de la otra clase esperando bailar. En caso de haberlo, bailan, y en caso contrario, espera a que llegue su contraparte.

El código para implementar esto es muy simple:

```

colaLideres = Semaphore(0)
colaSeguidores = Semaphore(0)
# ...
def lider():
    colaSeguidores.release()
    colaLideres.acquire()
    baila()
def seguidor():
    colaLideres.release()
    colaSeguidores.acquire()
    baila()

```

El patrón debe resultar ya familiar: es un *rendezvous*. La distinción es meramente semántica: en el *rendezvous* se necesitan dos hilos explícitamente, aquí se habla de dos clases de hilos.

Sobre este patrón base se pueden refinar muchos comportamientos. Por ejemplo, asegurar que sólo una pareja esté bailando al mismo tiempo, o asegurar que los hilos en espera vayan bailando en el orden en que llegaron. ORG-LIST-END-MARKER

Variables de condición

Las variables de condición presentan una extensión sobre el comportamiento de los mutexes, buscando darles la “inteligencia” de responder ante determinados eventos. Una variable de condición siempre opera *en conjunto con* un mutex, y en algunas implementaciones es necesario indicar cuál será dicho mutex desde la misma inicialización del objeto.⁷

Una variable de condición presenta las siguientes operaciones:

Espera Se le indica una condición y un mutex. El mutex tiene que haber sido ya adquirido.

Esta operación *libera* al mutex, y se bloquea hasta recibir una *notificación* de otro hilo o proceso. Una vez que la notificación es recibida, y antes de devolver la ejecución al hilo, *re-adquiere* el mutex.

Espera medida Tiene una semántica igual a la de la espera, pero recibe un argumento adicional, indicando el tiempo de expiración. Si pasado el tiempo de expiración no ha sido notificado, despierta al hilo regresándole un error (y sin re-adquirir el mutex).

Señaliza Requiere que el mutex ya haya sido adquirido. Despierta (señaliza) a uno o más hilos (algunas implementaciones permiten indicar como argumento a cuántos hilos) de los que están bloqueados en la espera asociada. *No libera* el mutex — Esto

⁷Mientras que otras implementaciones permiten que se declaren por separado, pero siempre que se invoca a una variable de condición, debe indicársele qué mutex estará empleando.

significa que el flujo de ejecución *se mantiene en el invocante*, quien tiene que salir de su sección crítica (entregar el mutex) antes de que otro de los hilos continúe ejecutando.

Señaliza a todos Indica a *todos los hilos* que estén esperando a esta condición.

La interfaz de hilos POSIX (pthreads) presenta la siguiente definición:

```
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;
int pthread_cond_init(pthread_cond_t *cond, pthread_condattr_t *cond_attr);
int pthread_cond_signal(pthread_cond_t *cond);
int pthread_cond_broadcast(pthread_cond_t *cond);
int pthread_cond_wait(pthread_cond_t *cond, pthread_mutex_t *mutex);
int pthread_cond_timedwait(pthread_cond_t *cond, pthread_mutex_t *mutex, const
                           struct timespec *abstime);
int pthread_cond_destroy(pthread_cond_t *cond);
```

4.3.4 Problema productor-consumidor

Planteamiento

En un entorno multihilos es común que haya una división de tareas tipo *línea de ensamblado*, que se puede generalizar a que un grupo de hilos van *produciendo* ciertas estructuras, a ser *consumidas* por otro grupo.

Un ejemplo de este problema puede ser un programa *orientado a eventos*, en que eventos de distinta naturaleza pueden producirse, y causan que se *disparen* los mecanismos que los puedan atender. Los eventos pueden *apilarse* en un buffer que será procesado por los hilos encargados conforme se vayan liberando. Esto impone ciertos requisitos, como:

- Debido a que el buffer es un recurso compartido por los hilos, agregar o retirar un elemento del buffer tiene que ser hecho de forma atómica. Si más de un proceso intentara hacerlo al mismo tiempo, se correría el riesgo de que se corrompan los datos.
- Si un consumidor está listo y el buffer está vacío, debe bloquearse (*no realizar espera activa!*) hasta que un productor genere un elemento.

Si no se tiene en cuenta la sincronización, el código sería tan simple como el siguiente:

```
import threading
buffer = []
threading.Thread(target=productor, args=[]).start()
threading.Thread(target=consumidor, args=[]).start()

def productor():
    while True:
        event = genera_evento()
        buffer.append(event)

def consumidor():
    while True:
        event = buffer.pop()
        procesa(event)
```

Pero el acceso a buffer no está protegido para garantizar la exclusión mutua, y podría quedar en un estado inconsistente si append() y pop() intentan manipular sus estructuras al mismo tiempo. Además, si bien en este ejemplo se asumió que hay un sólo hilo productor y un sólo hilo consumidor, se puede extender el programa para que haya varios hilos en cualquiera de estos roles.

evento no requiere de protección, dado que es una variable local a cada hilo.

Solución

Para resolver este problema se usarán dos semáforos: `mutex`, que protege el acceso a la sección crítica, y `elementos`. El valor almacenado en `elementos` indica, cuando es positivo, cuántos eventos pendientes hay por procesar, y cuando es negativo, cuántos consumidores están listos y esperando un evento.

Una solución a este problema puede ser:

```
import threading
mutex = threading.Semaphore(1)
elementos = threading.Semaphore(0)
buffer = []
threading.Thread(target=productor, args=[]).start()
threading.Thread(target=consumidor, args=[]).start()

def productor():
    while True:
        event = genera_evento()
        mutex.acquire()
        buffer.append(event)
        mutex.release()
        elementos.release()

def consumidor():
    while True:
        elementos.acquire()
        mutex.acquire()
        event = buffer.pop()
        mutex.release()
        event.process()
```

Se puede ver que la misma construcción, un semáforo, es utilizada de forma muy distinta por `mutex` y `elementos`. `mutex` implementa una exclusión mutua clásica, delimitando tanto como es posible (a una sola línea en este caso) al área crítica, y siempre apareando un `acquire()` con un `release()`. `elementos`, en cambio, es empleado como un verdadero semáforo: como una estructura para la sincronización. Sólo los hilos productores incrementan (sueltan) el semáforo, y sólo los consumidores lo decrementan (adquieren). Esto es, ambos semáforos comunican *al planificador* cuándo es posible despertar a algún consumidor.

Si se supone que `genera_evento()` es eficiente y no utiliza espera activa, esta implementación es óptima: deja en manos del planificador toda la espera necesaria, y no desperdicia recursos de cómputo esperando a que el siguiente elemento esté listo.

Como nota al pie, la semántica del módulo `threading` de Python incluye la declaración de contexto `with`. Todo objeto de `threading` que implemente `acquire()` y `release()` puede ser envuelto en un bloque `with`, aumentando la legibilidad y con exactamente la misma semántica; no se utilizó en este ejemplo para ilustrar el uso tradicional de los semáforos para implementar regiones de exclusión mutua, sin embargo y sólo para concluir con el ejemplo, la función `consumidor()` final podría escribirse así, y ser semánticamente idéntica:

```
def consumidor():
    while True:
        elementos.acquire()
```

```

with mutex:
    event = buffer.pop()
    event.process()

```

A pesar de ser más clara, no se empleará en este texto esta notación por dos razones. La primera es para mantener la conciencia de la semántica de las operaciones `acquire()` y `release()`, y la segunda es mantener la consistencia a través de las distintas implementaciones, ya que se encontrarán varios casos en que no es el mismo hilo el que adquiere un semáforo y el que lo libera (esto es, los semáforos no siempre son empleados como mutexes).

4.3.5 Bloqueos mutuos e inanición

Cuando hay concurrencia, además de asegurar la atomicidad de ciertas operaciones, se debe evitar dos problemas que son consecuencia natural de la existencia de la asignación de recursos de forma exclusiva:

Bloqueo mutuo (o *interbloqueo*; en inglés, *deadlock*) Situación que ocurre cuando dos o más procesos poseen determinados recursos, y cada uno queda detenido, a la espera de alguno de los que tiene el otro. El sistema puede seguir operando normalmente, pero ninguno de los procesos involucrados podrán avanzar.

Inanición (en inglés *resource starvation*) Situación en que un proceso no puede avanzar en su ejecución dado que necesita recursos que están (alternativamente) asignados a otros procesos.

El que se presenten estos conceptos aquí no significa que están *exclusivamente* relacionados con este tema: son conceptos que se enfrentan una y otra vez al hablar de asignación exclusiva de recursos — temática recurrente en el campo de los sistemas operativos.

4.3.6 Problema lectores-escritores

Planteamiento

Una estructura de datos puede ser accedida simultáneamente por muchos procesos *lectores*, pero si algún proceso está escribiendo, se debe evitar que cualquier otro lea (dado que podría encontrarse con los datos en un estado inconsistente). Los requisitos de sincronización son

- Cualquier cantidad de lectores puede estar leyendo al mismo tiempo.
- Los escritores deben tener acceso exclusivo a la sección crítica.
- En pos de un comportamiento más justo: se debe evitar que un influjo constante de procesos lectores dejen a un escritor en situación de *inanición*.

Discusión

Este problema puede ser generalizado como una *exclusión mutua categórica*: se debe separar el uso de un recurso según la categoría del proceso. La presencia de un proceso en la sección crítica no lleva a la exclusión de otros, pero sí hay *categorías* de procesos que tienen distintas reglas — Para los escritores sí hace falta una exclusión mutua completa.

Primera aproximación

Un primer acercamiento a este problema permite una resolución libre de bloqueos mutuos, empleando sólo tres estructuras globales: un contador que indica cuántos lectores hay en la sección crítica, un mutex protegiendo a dicho contador, y otro mutex indicando que no hay lectores ni escritores accediendo al buffer (o cuarto). Se implementan los mutexes con semáforos.

```

import threading
lectores = 0

```

```
mutex = threading.Semaphore(1)
cuarto_vacio = threading.Semaphore(1)

def escritor():
    cuarto_vacio.acquire()
    escribe()
    cuarto_vacio.release()

def lector():
    mutex.acquire()
    lectores = lectores + 1
    if lectores == 1:
        cuarto_vacio.acquire()
    mutex.release()

    lee()

    mutex.acquire()
    lectores = lectores - 1
    if lectores == 0:
        cuarto_vacio.release()
    mutex.release()
```

El semáforo `cuarto_vacio` sigue un patrón visto antes llamado *apagador*. El escritor utiliza al apagador como a cualquier mutex: lo utiliza para rodear a su sección crítica. El lector, sin embargo, lo emplea de otra manera. Lo primero que hace es verificar *si la luz está prendida*, esto es, si hay algún otro lector en el cuarto (si `lectores` es igual a 1). Si es el primer lector en entrar, prende la luz adquiriendo `cuarto_vacio` (lo cual evitará que un escritor entre). Cualquier cantidad de lectores puede entrar, actualizando su número. Cuando el último sale (`lectores` es igual a 0), apaga la luz.

El problema con esta implementación es que un flujo constante de lectores puede llevar a la inanición de un escritor, que está pacientemente parado esperando a que alguien apague la luz.

Solución

Para evitar esta condición de inanición, se puede agregar un *torniquete* evitando que lectores adicionales se *cuelen* antes del escritor. Reescribiendo:

```
import threading
lectores = 0
mutex = threading.Semaphore(1)
cuarto_vacio = threading.Semaphore(1)
torniquete = threading.Semaphore(1)

def escritor():
    torniquete.acquire()
    cuarto_vacio.acquire()
    escribe()
    cuarto_vacio.release()
    torniquete.release()

def lector():
```

```

global lectores
torniquete.acquire()
torniquete.release()

mutex.acquire()
lectores = lectores + 1
if lectores == 1:
    cuarto_vacio.acquire()
mutex.release()

lee()

mutex.acquire()
lectores = lectores - 1
if lectores == 0:
    cuarto_vacio.release()
mutex.release()

```

En la implementación de los escritores, esto puede parecer inútil: Únicamente se agregó un mutex redundante alrededor de lo que ya se tenía. Sin embargo, al obligar a que el lector pase por un torniquete antes de actualizar `lectores`, lo cual obligaría a que se mantenga encendida la luz, se obliga a que espere a que el escritor suelte este mutex exterior. Nuevamente se puede ver cómo la misma estructura es tratada de dos diferentes maneras: para el lector es un torniquete, y para el escritor es un mutex.

4.3.7 La cena de los filósofos

Planteamiento

Cinco filósofos se dan cita para comer arroz en una mesa redonda. En la mesa, cada uno de ellos se sienta frente a un plato. A su derecha, tiene un palito chino, y a su izquierda tiene otro.

Los filósofos sólo saben `pensar()` y `comer()`. Cada uno de ellos va a `pensar()` un tiempo arbitrario, hasta que le da hambre. El hambre es mala consejera, por lo que intenta `comer()`. Los requisitos son:

- Sólo un filósofo puede sostener determinado palito a la vez, esto es los palitos son recursos de acceso exclusivo.
- Debe ser imposible que un filósofo muera de inanición estando a la espera de un palito.
- Debe ser imposible que se presente un bloqueo mutuo.
- Debe ser posible que más de un filósofo pueda comer al mismo tiempo.

Discusión

En este caso, el peligro no es, como en el ejemplo anterior, que una estructura de datos sea sobreescrita por ser accedida por dos hilos al mismo tiempo, sino que se presenten situaciones en el curso normal de la operación que lleven a un bloqueo mutuo.

A diferencia del caso antes descripto, ahora se utilizarán los semáforos no como una herramienta para indicar al planificador cuándo despertar a uno de los hilos, sino como una herramienta de comunicación entre los propios hilos.

Primer acercamiento

Se puede representar a los palillos como un arreglo de semáforos, asegurando la exclusión mutua (esto es, sólo un filósofo puede sostener un palillo al mismo tiempo), pero eso no evita el bloqueo mutuo. Por ejemplo, si la solución fuera:

```

import threading
num = 5
palillos = [threading.Semaphore(1) for i in range(num)]
filosofos = [threading.Thread(target=filosofo, args=[i]).start() for i in range(num)]

def filosofo(id):
    while True:
        piensa(id)
        levanta_palillos(id)
        come(id)
        suelta_palillos(id)

def piensa(id):
    # ...
    print "%d - Tengo hambre..." % id

def levanta_palillos(id):
    palillos[(id+1) % num].acquire()
    print "%d - Tengo el palillo derecho" % id
    palillos[id].acquire()
    print "%d - Tengo ambos palillos" % id

def suelta_palillos(id):
    palillos[(id+1) % num].release()
    palillos[id].release()
    print "%d - Sigamos pensando..." % id

def come(id):
    print "%d - ¡A comer!" % id
    # ...

```

Podría pasar que todos los filósofos quieran comer al mismo tiempo, y el planificador dejara suspendidos a todos con el palillo derecho en la mano.

Solución

Ahora, ¿qué pasa si se hace que algunos filósofos sean *zurdos*? Esto es, que levanten primero el palillo izquierdo y luego el derecho:

```

def levanta_palillos(id):
    if (id % 2 == 0): # Zurdo
        palillo1 = palillos[id]
        palillo2 = palillos[(id+1) % num]
    else: # Diestro
        palillo1 = palillos[(id+1) % num]
        palillo2 = palillos[id]
    palillo1.acquire()
    print "%d - Tengo el primer palillo" % id
    palillo2.acquire()
    print "%d - Tengo ambos palillos" % id

```

Al asegurar que dos filósofos contiguos no intenten levantar el mismo palillo, se tiene la certeza de que no se producirán bloqueos mutuos. De hecho, incluso si sólo uno de los filósofos es zurdo, se puede demostrar que no habrá bloqueos:

```

def levanta_palillos(id):
    if id == 0: # Zurdo
        palillos[id].acquire()
        print "%d - Tengo el palillo izquierdo" % id
        palillos[(id+1) % num].acquire()
    else: # Diestro
        palillos[(id+1) % num].acquire()
        print "%d - Tengo el palillo derecho" % id
        palillos[id].acquire()
    print "%d - Tengo ambos palillos" % id

```

Cabe apuntar que ninguno de estos mecanismos asegura que en la mesa no haya *inanición*, sólo que no haya bloqueos mutuos.

4.3.8 Los fumadores compulsivos

Planteamiento

Hay tres fumadores empedernidos y un *agente* que, de tiempo en tiempo, consigue ciertos insumos. Los ingredientes necesarios para fumar son tabaco, papel y cerillos. Cada uno de los fumadores tiene una cantidad infinita de alguno de los ingredientes, pero no les gusta compartir. Afortunadamente, del mismo modo que no comparten, no son acaparadores⁸.

De tiempo en tiempo, el agente consigue una dosis de dos de los ingredientes — Por ejemplo, si deja en la mesa un papel y tabaco, el que trae los cerillos educadamente tomará los ingredientes, se hará un cigarro, y lo fumará.

Suhas Patil (1971) planteó este problema buscando demostrar que hay situaciones que no se pueden resolver con el uso de semáforos. Las condiciones planteadas son

- No se puede modificar el código del agente. Si el agente es un sistema operativo, ¡tiene sentido la restricción de no tenerle que notificar acerca de los flujos a cada uno de los programas que ejecuta!
- El planteamiento original de Patil menciona que no debe emplearse arreglos de semáforos o usar condicionales en el flujo. Esta segunda restricción haría efectivamente irresoluble al problema, por lo que se ignorará.

Primer acercamiento

Al haber tres distintos ingredientes, tiene sentido que se empleen tres distintos semáforos, para señalizar a los fumadores respecto a cada uno de los ingredientes. Un primer acercamiento podría ser:

```

import random
import threading
ingredientes = ['tabaco', 'papel', 'cerillo']
semaforos = {}
semaforo_agente = threading.Semaphore(1)
for i in ingredientes:
    semaforos[i] = threading.Semaphore(0)

threading.Thread(target=agente, args=[]).start()
fumadores = [threading.Thread(target=fumador, args=[i]).start() for i in ingre

```

⁸Esto es, no buscan obtener y conservar los recursos *preventivamente*, sino que los toman sólo cuando satisfacen *por completo* sus necesidades.

```

def agente():
    while True:
        semaforo_agente.acquire()
        mis_ingr = ingredientes[:]
        mis_ingr.remove(random.choice(mis_ingr))
        for i in mis_ingr:
            print "Proveyendo %s" % i
            semaforos[i].release()

def fumador(ingr):
    mis_semaf = []
    for i in semaforos.keys():
        if i != ingr:
            mis_semaf.append(semaforos[i])
    while True:
        for i in mis_semaf:
            i.acquire()
            fuma(ingr)
            semaforo_agente.release()

def fuma(ingr):
    print 'Fumador con %s echando humo...' % ingr

```

El problema en este caso es que, al tener que cubrir un número de ingredientes mayor a uno, utilizar sólo un semáforo ya no funciona: si `agente()` decide proveer papel y cerillos, nada garantiza que no sea el fumador ['cerillo'] el que reciba la primer señal o que fumador ['tabaco'] reciba la segunda — Para que este programa avance hace falta, más que otra cosa, la buena suerte de que las señales sean recibidas por el proceso indicado.

Solución

Una manera de evitar esta situación es la utilización de *intermediarios* encargados de notificar al hilo adecuado. Partiendo de que, respetando la primer restricción impuesta por Patil, no se puede modificar el código del agente, se implementan los intermediarios, se reimplementan los fumadores y se agregan algunas variables globales, de esta manera:

```

que_tengo = {}
semaforos_interm = {}
for i in ingredientes:
    que_tengo[i] = False
    semaforos_interm[i] = threading.Semaphore(0)
interm_mutex = threading.Semaphore(1)
intermediarios = [threading.Thread(target=intermediario, args=[i]).start() for i in ingredientes]

def fumador(ingr):
    while True:
        semaforos_interm[ingr].acquire()
        fuma(ingr)
        semaforo_agente.release()

def intermediario(ingr):
    otros_ingr = ingredientes[:]

```

```

otros_ingr.remove(ingr)
while True:
    semaforos[ingr].acquire()
    interm_mutex.acquire()
    for i in otros_ingr:
        if que_tengo[i]:
            que_tengo[i] = False
            semaforos_interm[i].release()
            break
        que_tengo[i] = True
    interm_mutex.release()

```

Si bien se ve que el código de `fumador()` se simplifica (dado que ya no tiene que efectuar ninguna verificación), `intermediario()` tiene mayor complejidad. El elemento clave de su lógica es que, si bien el `agente()` (el sistema operativo) seguirá enviando una señal por cada ingrediente disponible, los tres intermediarios se sincronizarán empleando al arreglo `que_tengo` (protegido por `interm_mutex`), y de este modo cada hilo (independientemente del orden en que fue invocado) señalizará a los otros intermediarios qué ingredientes hay en la mesa, y una vez que sepa a qué fumador notificar, dejará el estado listo para recibir una nueva notificación.

4.3.9 Otros mecanismos

Más allá de los mecanismos basados en mutexes y semáforos, existen otros que emplean diferentes niveles de *encapsulamiento* para proteger las abstracciones. A continuación se presentan muy brevemente algunos de ellos.

Monitores

El principal problema con los mecanismos anteriormente descritos es que no sólo hace falta encontrar un mecanismo que permita evitar el acceso simultáneo a la sección crítica sin caer en bloqueos mutuos o inanición, sino que hay que *implementarlo correctamente*, empleando una semántica que requiere de bastante entrenamiento para entender correctamente.

Además, al hablar de procesos que compiten por recursos de una forma *hostil*, la implementación basada en semáforos puede resultar insuficiente. A modo de ejemplo, se mostrará por qué en el modelo original de Djikstra (así como en los ejemplos presentados anteriormente) sólo existen las operaciones de incrementar y decrementar, y no se permite verificar el estado (como lo ofrece `sem_trywait()` en `pthreads`):

```

while (sem_trywait(semaforo) != 0) {}
seccion_critica();
sem_post(semaforo);

```

El código presentado es absolutamente válido — Pero cae en una *espera activa* que desperdicia innecesariamente y constantemente tiempo de procesador (y no tiene garantía de tener más éxito que una espera pasiva, como sería el caso con un `sem_wait()`).

Por otro lado, algún programador puede creer que su código ejecutará suficientemente rápido y con suficientemente baja frecuencia para que la probabilidad de que usar la sección crítica le cause problemas sea muy baja. Es frecuente ver ejemplos como el siguiente:

```

/* Cruzamos los dedos... a fin de cuentas, ejecutaremos con baja frecuencia!
seccion_critica();

```

Los perjuicios causados por este programador resultan obvios. Sin embargo, es común ver casos como este.

Los *monitores* son estructuras provistas por el lenguaje o entorno de desarrollo que *encapsulan tanto a los datos como a las funciones que los pueden manipular*, e impiden el acceso directo a las funciones potencialmente peligrosas — En otras palabras, son *tipos de datos abstractos* (ADTs), *clases de objetos*, y *exponen* una serie de *métodos públicos*, además de poseer *métodos privados* que emplean internamente.

Al no presentar al usuario/programador una interfaz que puedan *subvertir*, el monitor mantiene todo el código necesario para asegurar el acceso concurrente a los datos en un sólo lugar.

Un monitor puede implementarse utilizando cualquiera de los mecanismos de sincronización presentados anteriormente — la diferencia radica en que esto se hace *en un solo lugar*. Los programas que quieran emplear el recurso protegido lo hacen incluyendo el código del monitor como módulo / biblioteca, lo cual fomenta la *reutilización de código*.

Como ejemplo, el lenguaje de programación *Java* implementa sincronización vía monitores entre hilos como una propiedad de la declaración de método, y lo implementa directamente en la JVM. Si se declara un método de la siguiente manera:

```
public class SimpleClass {
    // ...
    public synchronized void metodoSeguro() {
        /* Implementación de metodoSeguro() */
        // ...
    }
}
```

Y se inicializa a un `SimpleClass sc = new SimpleClass()`, cuando se llame a `sc.metodoSeguro()`, la máquina virtual verificará si ningún otro proceso está ejecutando `metodoSeguro()`; en caso de que no sea así, le permitirá la ejecución obteniendo el candado, y en caso de sí haberlo, el hilo se bloqueará hasta que el candado sea liberado — Esto es, la propiedad `synchronized` hace que todo acceso al método en cuestión sea protegido por una mutex.

El modelo de sincronización basado en monitores no sólo provee la exclusión mutua. A través de *variables de condición* (VCs) se puede también emplear una semántica parecida (aunque no igual) a la de los semáforos, con los métodos `var.wait()` y `var.signal()`. En el caso de los monitores, `var.wait()` suspende al hilo hasta que otro hilo ejecute `var.signal()`; en caso de no haber ningún proceso esperando, `var.signal()` no tiene ningún efecto (no cambia el estado de `var`, a diferencia de lo que ocurre con los semáforos)

Aquí se presenta, a modo de ilustración, la resolución del problema de la *cena de los filósofos* en C⁹. Esto demuestra, además, que si bien se utiliza semántica de orientación a objetos, no sólo los lenguajes clásicamente relacionados con la programación orientada a objetos permiten emplear monitores.

```
/* Implementacion para cinco filósofos */
#define PENSANDO 1
#define HAMBRIENTO 2
#define COMIENDO 3

pthread_cond_t VC[5]; /* Una VC por filósofo */
pthread_mutex_t M;      /* Mutex para el monitor */
```

⁹Implementación basada en el [ejemplo de Ted Baker](#), sobre la solución propuesta por Tanenbaum

```

int estado[5];           /* Estado de cada filósofo */

void palillos_init () {
    int i;
    pthread_mutex_init(&M, NULL);
    for (i = 0; i < 5; i++) {
        pthread_cond_init(&VC[i], NULL);
        estado[i] = PENSANDO;
    }
}

void toma_palillos (int i) {
    pthread_mutex_lock(&M)
    estado[i] = HAMBRIENTO;
    actualiza(i);
    while (estado[i] == HAMBRIENTO)
        pthread_cond_wait(&VC[i], &M);
    pthread_mutex_unlock(&M);
}

void suelta_palillos (int i) {
    estado[i] = PENSANDO;
    actualiza((i - 1) % 5);
    actualiza((i + 1) % 5);
    pthread_mutex_unlock(&M);
}

void come(int i) {
    printf("El filosofo %d esta comiendo\n", i);
}

void piensa(int i) {
    printf("El filosofo %d esta pensando\n", i);
}

/* No incluir 'actualiza' en los encabezados, */
/* es una función interna/privada */
int actualiza (int i) {
    if ((estado[(i - 1) % 5] != COMIENDO) &&
        (estado[i] == HAMBRIENTO) &&
        (estado[(i + 1) % 5] != COMIENDO)) {
        estado[i] = COMIENDO;
        pthread_cond_signal(&VC[i]);
    }
    return 0;
}

```

Esta implementación evita los bloqueos mutuos señalizando el estado de cada uno de los filósofos en el arreglo de variables `estado[]`.

La lógica base de esta resolución marca en la verificación del estado propia y de los vecinos siempre que hay un cambio de estado: cuando el filósofo i llama a la función `toma_palillos(i)`, esta se limita a adquirir el mutex M , marcar su estado como

HAMBRIENTO, y llamar a la función interna `actualiza(i)`. Del mismo modo, cuando el filósofo `i` termina de comer y llama a `suelta_palillos(i)`, esta función marca su estado como PENSANDO e invoca a `actualiza()` dos veces: una para el vecino izquierdo y una para el vecino derecho.

Es importante recalcar que, dado que esta solución está estructurada como un monitor, ni `actualiza()` ni las variables que determinan el estado del sistema (`VC`, `M`, `estado`) son expuestas a los hilos invocantes.

La función `actualiza(i)` es la que se encarga de verificar (y modificar, de ser el caso) el estado no sólo del filósofo invocante, sino que de sus vecinos. La lógica de `actualiza()` permite resolver este problema abstrayendo (y eliminando) a los molestos palillos: en vez de preocuparse por cada palillo individual, `actualiza()` impide que dos filósofos vecinos estén COMIENDO al mismo tiempo, y dado que es invocada tanto cuando un filósofo `toma_palillos()` como cuando `suelta_palillos()`, otorga el turno al vecino sin que éste tenga que adquirir explícitamente el control.

Estas características permite que la lógica central de cada uno de los filósofos se simplifique a sólo:

```
void *filosofo(void *arg) {
    int self = *(int *) arg;
    for (;;) {
        piensa(self);
        toma_palillos(self);
        come(self);
        suelta_palillos(self);
    }
}

int main() {
    int i;
    pthread_t th[5]; /* IDs de los hilos filósofos */
    pthread_attr_t attr = NULL;
    palillos_init();
    for (i=0; i<5; i++)
        pthread_create(&th[i], attr, filosofo, (int*) &i);
    for (i=0; i<5; i++)
        pthread_join(th[i],NULL);
}
```

Al ser una solución basada en monitor, el código que invoca a `filosofo(i)` no tiene que estar al pendiente del mecanismo de sincronización empleado, puede ser comprendido más fácilmente por un lector casual y no brinda oportunidad para que un mal programador haga mal uso del mecanismo de sincronización.

Memoria transaccional

Un área activa de investigación hoy en día es la de la *memoria transaccional*. La lógica es ofrecer primitivas que protejan a *un conjunto de accesos a memoria* del mismo modo que ocurre con las bases de datos, en las que tras abrir una transacción, se puede realizar una gran cantidad (no ilimitada) de tareas, y al terminar con la tarea, *confirmarlas (commit)* o *rechazarlas (rollback)* atómicamente — y, claro, el sistema operativo indicará éxito o fracaso de forma atómica al conjunto entero.

Esto facilitaría mucho más aún la sincronización: en vez de hablar de *secciones críticas*, se podría reintentar la transacción y sólo preocuparse de revisar si fue exitosa o no — por ejemplo:

```

do {
    begin_transaction();
    var1 = var2 * var3;
    var3 = var2 - var1;
    var2 = var1 / var2;
} while (! commit_transaction());

```

Si en el transcurso de la transacción algún otro proceso modifica alguna de las variables, la transacción se abortará, pero se puede volver a ejecutar.

Claro está, el ejemplo presentado desperdicia recursos de proceso (lleva a cabo los cálculos al tiempo que va modificando las variables), por lo cual sería un mal ejemplo de sección crítica.

Hay numerosas implementaciones en software de este principio (Software Transactional Memory, STM) para los principales lenguajes, aunque el planteamiento ideal sigue apuntando a una implementación en hardware. Hay casos en que, sin embargo, esta técnica puede aún llevar a resultados inconsistentes (particularmente si un proceso lector puede recibir los valores que van cambiando en el tiempo por parte de un segundo proceso), y el costo computacional de dichas operaciones es elevado, sin embargo, es una construcción muy poderosa.

4.4 Bloqueos mutuos

Un bloqueo mutuo puede ejemplificarse con la situación que se presenta cuando cuatro automovilistas llegan al mismo tiempo al cruce de dos avenidas del mismo rango en que no hay un semáforo, cada uno desde otra dirección. Los reglamentos de tránsito señalan que la precedencia la tiene *el automovilista que viene más por la derecha*. En este caso, cada uno de los cuatro debe ceder el paso al que tiene a la derecha — Y ante la ausencia de un criterio humano que rompa el bloqueo, deberían todos mantenerse esperando por siempre.

Un bloqueo mutuo se presenta cuando (*Condiciones de Coffman*) (La Red, p. 185):

1. Los procesos reclaman control exclusivo de los recursos que piden (condición de *exclusión mutua*).
2. Los procesos mantienen los recursos que ya les han sido asignados mientras esperan por recursos adicionales (condición de *espera por*).
3. Los recursos no pueden ser extraídos de los procesos que los tienen hasta su completa utilización (condición de *no apropiatividad*).
4. Existe una cadena circular de procesos en la que cada uno mantiene a uno o más recursos que son requeridos por el siguiente proceso de la cadena (condición de *espera circular*).

Las primeras tres condiciones son *necesarias pero no suficientes* para que se produzca un bloqueo; su presencia puede indicar una situación de riesgo. Sólo cuando se presentan las cuatro se puede hablar de un bloqueo mutuo efectivo.

Otro ejemplo clásico es un sistema con dos unidades de cinta (dispositivos de acceso secuencial y no compatible), en que los procesos *A* y *B* requieren de ambas unidades. Dada la siguiente secuencia:

1. *A* solicita una unidad de cinta y se bloquea.
2. *B* solicita una unidad de cinta y se bloquea.
3. El sistema operativo otorga la unidad 1 a *A* y lo vuelve a poner en ejecución.
4. *A* continúa procesando; termina su periodo de ejecución.
5. El sistema operativo otorga la unidad 2 a *B* y lo vuelve a poner en ejecución.
6. *B* solicita otra unidad de cinta y se bloquea.

7. El sistema operativo no tiene otra unidad de cinta por asignar. Mantiene a *B* bloqueado; otorga el control de vuelta a *A*.
8. *A* solicita otra unidad de cinta y se bloquea
9. El sistema operativo no tiene otra unidad de cinta por asignar. Mantiene bloqueado tanto *A* como *B* y otorga el control de vuelta a otro proceso (o queda en espera). En este caso ni *A* ni *B* serán desbloqueados nunca.

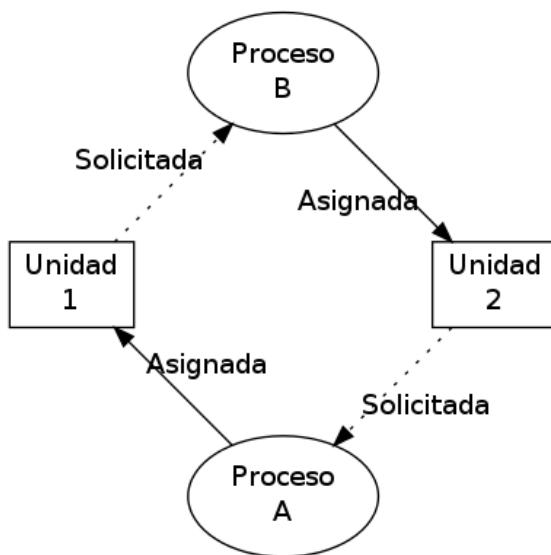


Figura 4.6: Esquema clásico de un bloqueo mutuo simple: Los procesos *A* y *B* esperan mutuamente para el acceso a las unidades de cinta 1 y 2.

Sin una política de prevención o resolución de bloqueos mutuos, no hay modo de que *A* o *B* continúen su ejecución. Se verán algunas estrategias para enfrentar a los bloqueos mutuos.

En el apartado de *Exclusión mutua*, los hilos presentados estaban diseñados para cooperar explícitamente. El rol del sistema operativo va más allá, tiene que implementar políticas que eviten, en la medida de lo posible, dichos bloqueos.

Las políticas tendientes a otorgar los recursos lo antes posible cuando son solicitadas pueden ser vistas como *liberales*, en tanto que las que controlan más la asignación de recursos, *conservadoras*.

Las líneas principales que describen a las estrategias para enfrentar situaciones de bloqueo (La Red, p. 188) son:

Prevención se centra en modelar el comportamiento del sistema para que *elimine toda posibilidad* de que se produzca un bloqueo. Resulta en una utilización subóptima de recursos.

Evasión busca imponer condiciones menos estrictas que en la prevención, para intentar lograr una mejor utilización de los recursos. Si bien no puede evitar *todas las posibilidades* de un bloqueo, cuando éste se produce busca *evitar sus consecuencias*.

Detección y recuperación el sistema *permite* que ocurran los bloqueos, pero busca *determinar si ha ocurrido* y tomar medidas para eliminarlo.

Busca despejar los bloqueos presentados para que el sistema continúe operando sin ellos.

4.4.1 Prevención de bloqueos

Se presentan a continuación algunos algoritmos que implementan la prevención de bloqueos.

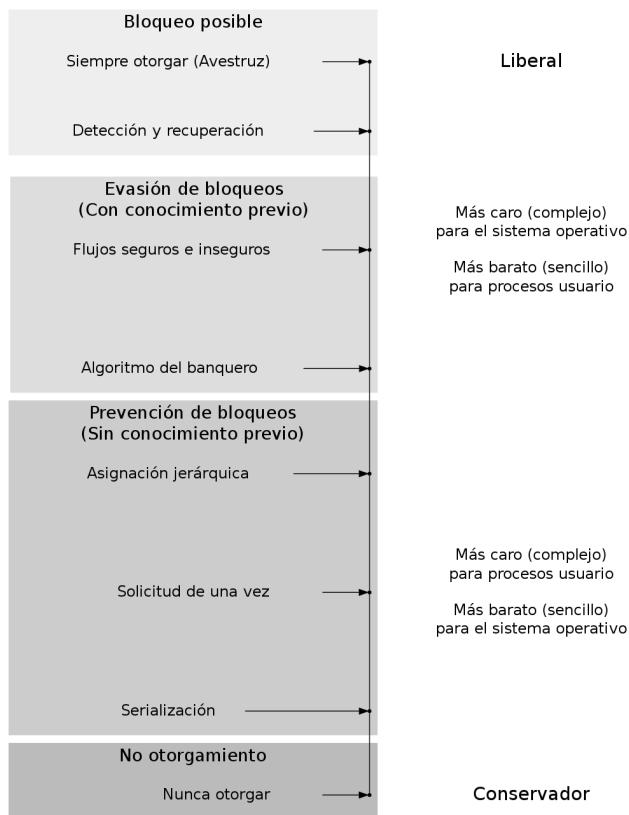


Figura 4.7: Espectro liberal—conservador de esquemas para evitar bloqueos

Serialización

Una manera de evitar bloqueos *por completo* sería el que un sistema operativo jamás asignara recursos a más de un proceso a la vez — Los procesos podrían seguir efectuando cálculos o empleando recursos *no rivales* (que no requieran acceso exclusivo — por ejemplo, empleo de archivos en el disco, sin que exista un acceso directo del proceso al disco), pero sólo uno podría obtener recursos de forma exclusiva al mismo tiempo. Este mecanismo sería la *serialización*, y la situación antes descrita se resolvería de la siguiente manera:

1. A solicita una unidad de cinta y se bloquea
2. B solicita una unidad de cinta y se bloquea
3. El sistema operativo otorga la unidad 1 a A y lo vuelve a poner en ejecución
4. A continúa procesando; termina su periodo de ejecución
5. El sistema operativo mantiene bloqueado a B, dado que A tiene un recurso
6. A solicita otra unidad de cinta y se bloquea
7. El sistema operativo otorga la unidad 2 a A y lo vuelve a poner en ejecución
8. A libera la unidad de cinta 1
9. A libera la unidad de cinta 2 (y con ello, el bloqueo de uso de recursos)
10. El sistema operativo otorga la unidad 1 a B y lo vuelve a poner en ejecución
11. B solicita otra unidad de cinta y se bloquea
12. El sistema operativo otorga la unidad 2 a B y lo vuelve a poner en ejecución
13. B libera la unidad de cinta 1
14. B libera la unidad de cinta 2

Si bien la serialización resuelve la situación aquí mencionada, el mecanismo empleado es subóptimo dado que puede haber hasta $n-1$ procesos esperando a que uno libere los recursos.

Un sistema que implementa una política de asignación de recursos basada en la serialización, si bien no caerá en bloqueos mutuos, sí tiene un peligro fuerte de caer en *inanición*.

Retención y espera (*advance claim*)

Otro ejemplo de política preventiva *menos conservadora* sería la *retención y espera* o *reserva* (*advance claim*): que todos los programas declaran al iniciar su ejecución qué recursos van a requerir. Los recursos son apartados para su uso exclusivo hasta que el proceso termina, pero el sistema operativo puede seguir atendiendo solicitudes *que no rivalicen*: si a los procesos *A* y *B* anteriores se suman procesos *C* y *D*, pero requieren otro tipo de recursos, podrían ejecutarse en paralelo *A*, *C* y *D*, y una vez que *A* termine, podrían continuar ejecutando *B*, *C* y *D*.

El bloqueo resulta ahora imposible por diseño, pero el usuario que inició *B* tiene una percepción de injusticia dado el tiempo que tuvo que esperar para que su solicitud fuera atendida — de hecho, si *A* es un proceso de larga duración (incluso si requiere la unidad de cinta sólo por un breve periodo), esto lleva a que *B* sufra una *inanición* innecesariamente prolongada.

Además, la implementación de este mecanismo preventivo requiere que el programador sepa por anticipado qué recursos requerirá — y esto en la realidad muchas veces es imposible. Si bien podría diseñarse una estrategia de lanzar procesos *representantes* (o *proxy*) solicitando recursos específicos cuando éstos hicieran falta, esto sólo transferiría la situación de bloqueo por recursos a bloqueo por procesos — y un programador poco cuidadoso podría de todos modos desencadenar la misma situación.

Solicitud de una vez (*one-shot*)

Otro mecanismo de prevención de bloqueos sería que los recursos se otorguen exclusivamente a aquellos procesos que *no poseen ningún recurso*. Esta estrategia rompería la condición de Coffman *espera por*, haciendo imposible que se presente un bloqueo.

En su planteamiento inicial, este mecanismo requería que un proceso declarara *una sola vez* qué recursos requeriría, pero posteriormente la estrategia se modificó, permitiendo que un proceso solicite recursos nuevamente, pero únicamente a condición de que previo a hacerlo *renuncien a los recursos* que tenían en ese momento — Claro, pueden volver a incluirlos en la operación de solicitud.

Al haber una *renuncia explícita*, se imposibilita de forma tajante que un conjunto de procesos entre en condición de bloqueo mutuo.

Las principales desventajas de este mecanismo son:

- Requiere cambiar la lógica de programación para tener puntos más definidos de adquisición y liberación de recursos.
- Muchas veces no basta con la *readquisición* de un recurso, sino que es necesario *mantenerlo bloqueado*. Volviendo al ejemplo de las unidades de cinta, un proceso que requiera ir generando un archivo largo no puede arriesgarse a *soltarla*, pues podría ser entregada a otro proceso y corromperse el resultado.

Asignación jerárquica

Otro mecanismo de evasión es la asignación *jerárquica* de recursos. Bajo este mecanismo, se asigna una prioridad o *nivel jerárquico* a cada recurso o clase de recursos.¹⁰ La condición básica es que, una vez que un proceso obtiene un recurso de determinado nivel, sólo puede solicitar recursos adicionales de niveles superiores. En caso de requerir dos dispositivos ubicados al mismo nivel, tiene que hacerse de forma atómica.

¹⁰Incluso varios recursos distintos, o varias clases, pueden compartir prioridad, aunque esto dificultaría la programación. Podría verse a la {solicitud de una vez} como un caso extremo de asignación jerárquica, con una jerarquía plana.

De este modo, si las unidades de cinta tienen asignada la prioridad x , P_1 sólo puede solicitar dos unidades de cinta por medio de *una sola operación*. En caso de también requerir dos unidades de cinta el proceso P_2 al mismo tiempo, al ser atómicas las solicitudes, éstas le serán otorgadas a sólo un de los dos procesos, por lo cual no se presentará bloqueo.

Además, el crear una jerarquía de recursos permitiría ubicar los recursos más escasos o *peleados* en la cima de la jerarquía, reduciendo las situaciones de contención en que varios procesos compiten por dichos recursos — sólo llegarían a solicitarlos aquellos procesos que ya tienen *asegurado* el acceso a los demás recursos que vayan a emplear.

Sin embargo, este ordenamiento es demasiado estricto para muchas situaciones del mundo real. El tener que renunciar a ciertos recursos para adquirir uno de menor prioridad y *volver a competir por ellos*, además de resultar contraintuitivo para un programador, resulta en esperas frustrantes. Este mecanismo llevaría a los procesos a acaparar recursos de baja prioridad, para evitar tener que ceder y re-adquirir recursos más altos, por lo que conduce a una alta inanición.

4.4.2 Evasión de bloqueos

Para la evasión de bloqueos, el sistema partiría de poseer, además de la información descrita en el caso anterior, información acerca de *cuándo* requiere un proceso utilizar cada recurso. De este modo, el planificador puede marcar qué orden de ejecución (esto es, qué *flujos*) entre dos o más procesos son *seguros* y cuáles son *inseguros*.

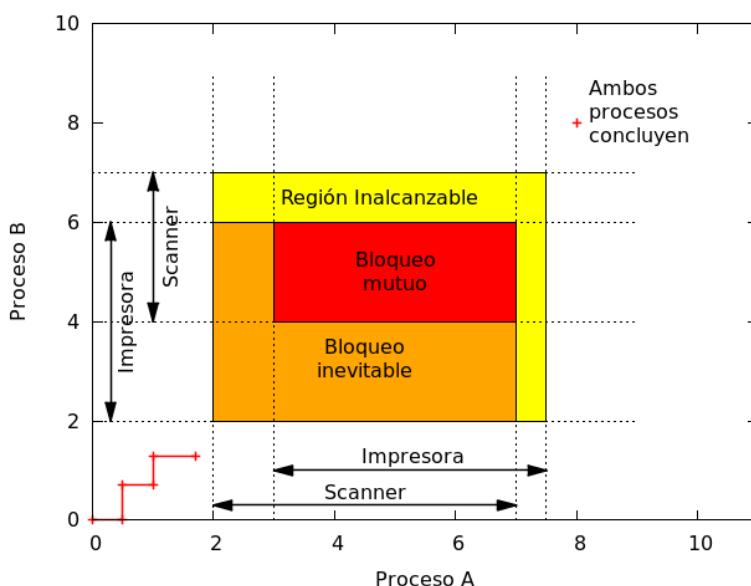


Figura 4.8: Evasión de bloqueos: Los procesos A (horizontal) y B (vertical) requieren del acceso exclusivo a un plotter y una impresora, exponiéndose a bloqueo mutuo.

El análisis de la interacción entre dos procesos se representa como en la figura 4.8; el avance es marcado en sentido horizontal para el proceso A , o vertical para el proceso B ; en un sistema multiprocesador, podría haber avance mutuo, y lo se indicaría en diagonal.

En el ejemplo presentado, el proceso A solicita acceso exclusivo al scanner durante $2 \leq t_A \leq 7$ y a la impresora durante $3 \leq t_A \leq 7,5$, mientras que B solicita acceso exclusivo a la impresora durante $2 \leq t_B \leq 6$ y al scanner durante $4 \leq t_B \leq 7$.

Al saber cuándo reclama y libera un recurso cada proceso, se puede marcar cuál es el área *segura* para la ejecución y cuándo se está aproximando a un área de riesgo.

En el caso mostrado, si bien el bloqueo mutuo sólo se produciría formalmente en

cualquier punto¹¹ en $3 \leq t_A \leq 7$, y $4 \leq t_B \leq 6$ (indicado con el recuadro rojo, *Bloqueo mutuo*).

Pero la existencia del recuadro que indica el bloqueo mutuo podría ser revelada con anterioridad: si el flujo entra en el área marcada como *Bloqueo inminente*, en color naranja (en $3 \leq t_A \leq 7$ y $2 \leq t_B \leq 6$), resulta *ineludible* caer en el bloqueo mutuo.

La región de bloqueo inminente ocurre a partir de que A obtuvo el scanner y B obtuvo la impresora. Si en $t_A = 2,5$ y $t_B = 3$ se cede la ejecución a A por 0.5 unidades, se llegará al punto en que solicita la impresora ($t_A = 3$), y no habrá más remedio que ejecutar B ; al avanzar B 0.5 unidades requerirá al scanner, y se habrá desencadenado el bloqueo mutuo. Un caso análogo ocurre, claro está, si desde el punto de inicio se ejecutara primero B y luego A .

Dadas las anteriores condiciones, y conociendo estos patrones de uso, el sistema operativo evitará entrar en el área de bloqueo inminente: el sistema mantendrá *en espera* a B si $t_B \leq 2$ mientras $2 \leq t_A \leq 6$, y mantendrá a A en espera si $t_A \leq 2$ cuando $2 \leq t_B \leq 6$.

La región marcada como *inalcanzable* en color amarillo, no representa ningún peligro: sólo indica aquellos estados en que resulta imposible entrar. Incluso una vez evadido el bloqueo (por ejemplo, si B fue suspendido en $t_B = 1,8$ y A avanza hasta pasar $t_A = 7$, si el sistema operativo vuelve a dar la ejecución a B , este sólo podrá avanzar hasta $t_B = 2$, punto en que B solicita la impresora. Para que B continúe, es necesario avanzar hasta $t_A > 7,5$ para que B siga avanzando).

Este mecanismo proveería una mejor respuesta que los vistos en el apartado de *prevención de bloqueos*, pero es todavía más difícil de aplicar en situaciones reales. Para poder implementar un sistema con evasión de bloqueos, tendría que ser posible hacer un análisis estático previo del código a ejecutar, y tener un listado total de recursos estático. Estos mecanismos podrían ser efectivos en sistemas de uso especializado, pero no en sistemas operativos (o planificadores) genéricos.

Algoritmo del banquero

Edsger Dijkstra propuso un algoritmo de asignación de recursos orientado a la evasión de bloqueos a ser empleado para el sistema operativo THE (desarrollado entre 1965 y 1968 en la Escuela Superior Técnica de Eindhoven, Technische Hogeschool Eindhoven), un sistema multiprogramado organizado en anillos de privilegios. El nombre de este algoritmo proviene de que busca que el sistema opere cuidando de tener siempre la liquidez (nunca entrar a *estados inseguros*) para satisfacer los préstamos (recursos) solicitados por sus clientes (quienes a su vez tienen una línea de crédito pre-autorizada por el banco).

Este algoritmo permite que el conjunto de recursos solicitado por los procesos en ejecución en el sistema sea mayor a los recursos físicamente disponibles, pero a través de un monitoreo y control en su asignación, logra este nivel de *sobre-compromiso* sin poner en riesgo la operación correcta del sistema.

Este algoritmo debe ejecutarse cada vez que un proceso solicita recursos; el sistema evita caer en situaciones conducentes a un bloqueo mutuo ya sea denegando o posponiendo la solicitud. El requisito particular es que, al iniciar, cada proceso debe *anunciar su reclamo máximo* (llamese `claim()`) al sistema: el número máximo de recursos de cada tipo que va a emplear a lo largo de su ejecución — esto sería implementado como una llamada al sistema. Una vez que un proceso presentó su reclamo máximo de recursos, cualquier llamada subsecuente a `claim()` falla. Claro está, si el proceso anuncia una necesidad mayor al número existente de recursos de algún tipo, también falla dado que el sistema no será capaz de cumplirlo.

Para el algoritmo del banquero:

¹¹En realidad, sólo sería posible *tocar* el margen izquierdo o inferior de este bloque: al caer en bloqueo mutuo, avanzar hacia su área interior resultaría imposible.

Estado matrices de recursos disponibles, reclamos máximos y asignación de recursos a los procesos en un momento dado.

Estado seguro un estado en el cual todos los procesos pueden ejecutar hasta el final sin encontrar un bloqueo mutuo.

Estado inseguro todo estado que no garantice que todos los procesos puedan ejecutar hasta el final sin encontrar un bloqueo mutuo.

Este algoritmo típicamente trabaja basado en diferentes *categorías* de recursos, y los reclamos máximos anunciados por los procesos son por cada una de las categorías.

El estado está compuesto, por clase de recursos y por proceso, por:

Reclamado número de instancias de este recurso que han sido reclamadas.

Asignado número de instancias de este recurso actualmente asignadas a procesos en ejecución.

Solicitado número de instancias de este recurso actualmente pendientes de asignar (solicitudes hechas y no cumplidas).

Además de esto, el sistema mantiene globalmente, por clase de recursos:

Disponibles número total de instancias de este recurso disponibles al sistema.

Libres número de instancias de este recurso que no están actualmente asignadas a ningún proceso.

Cada vez que un proceso solicita recursos, se calcula cuál sería el estado resultante de otorgar dicha solicitud, y se otorga siempre que:

- No haya reclamo por más recursos que los disponibles.
- Ningún proceso solicite (o tenga asignados) recursos por encima de su reclamo.
- La suma de los recursos *asignados* por cada categoría no sea mayor a la cantidad de recursos *disponibles* en el sistema para dicha categoría.

Formalmente, y volviendo a la definición de un estado seguro: un estado *es seguro* cuando hay una secuencia de procesos (denominada *secuencia segura*) tal que:

1. Un proceso *j* puede necesariamente terminar su ejecución, incluso si solicitara todos los recursos que permite su reclamo, dado que hay suficientes recursos libres para satisfacerlo.
2. Un segundo proceso *k* de la secuencia puede terminar si *j* termina y libera todos los recursos que tiene, porque sumado a los recursos disponibles ahora, con aquellos que liberaría *j*, hay suficientes recursos libres para satisfacerlo.
3. El *i*-ésimo proceso puede terminar si todos los procesos anteriores terminan y liberan sus recursos.

En el peor de los casos, esta secuencia segura llevaría a bloquear todas las solicitudes excepto las del único proceso que puede avanzar sin peligro en el orden presentado.

Se presenta un ejemplo simplificando, asumiendo sólo una clase de procesos, e iniciando con 2 instancias libres:

Proceso	Asignado	Reclamando
A	4	6
B	4	11
C	2	7

A puede terminar porque sólo requiere de 2 instancias adicionales para llegar a las 6 que indica en su reclamo. Una vez que termine, liberará sus 6 instancias. Se le asignan entonces las 5 que solicita a C, para llegar a 7. Al terminar éste, habrá 8 disponibles, y asignándole 7 a B se garantiza poder terminar. La secuencia (A, C, B) es una secuencia segura.

Sin embargo, el siguiente estado es inseguro (asumiendo también dos instancias libres):

Proceso	Asignado	Reclamado
A	4	6
B	4	11
C	2	9

A puede terminar, pero no se puede asegurar que *B* o *C* puedan hacerlo, porque incluso una vez terminando *A*, se tendrían sólo 6 instancias no asignadas.

Es necesario apuntar que no hay *garantía* de que continuar a partir de este estado lleve a un bloqueo mutuo, dado que *B* o *C* pueden no incrementar ya su utilización hasta cubrir su reclamo, esto es, puede que lleguen a finalizar sin requerir más recursos, ya sea porque ya los emplearon y liberaron, o porque el *uso efectivo* de recursos requeridos sencillamente resulte menor al del reclamo inicial.

El algoritmo del banquero, en el peor caso, puede tomar $O(n!)$, aunque típicamente ejecuta en $O(n^2)$. Una implementación de este algoritmo podría ser:

```

l = [1, 2, 3, 4, 5]; # Todos los procesos del sistema
s = []; # Secuencia segura
while ! l.empty? do
  p = l.select {|id| asignado[id] - reclamado[id] > libres}.first
  raise Exception, 'Estado inseguro' if p.nil?
  libres += asignado[p]
  l.delete(p)
  s.push(p)
end
puts "La secuencia segura encontrada es: #{s}" % s

```

Hay refinamientos sobre este algoritmo que logran resultados similares, reduciendo su costo de ejecución (se debe recordar que es un procedimiento que puede ser llamado con muy alta frecuencia), como el desarrollado por Habermann (ref: Finkel, p.136).

El algoritmo del banquero es un algoritmo conservador, dado que evita entrar en un estado inseguro a pesar de que dicho estado no lleve con certeza a un bloqueo mutuo. Sin embargo, su política es la más liberal que permite asegurar que no se caerá en bloqueos mutuos, sin conocer el *orden y tiempo* en que cada uno de los procesos requeriría los recursos.

Una desventaja fuerte de todos los mecanismos de evasión de bloqueos es que requieren saber por anticipado los reclamos máximos de cada proceso, lo cual puede no ser conocido en el momento de su ejecución.

4.4.3 Detección y recuperación de bloqueos

La detección de bloqueos es una forma de *reaccionar* ante una situación de bloqueo que ya se produjo y de buscar la mejor manera de salir de ella. La detección de bloqueos podría ser una tarea *periódica*, y si bien no puede prevenir situaciones de bloqueo, puede detectarlas una vez que ya ocurrieron y limitar su impacto.

Manteniendo una lista de recursos asignados y solicitados, el sistema operativo puede saber cuando un conjunto de procesos están esperándose mutuamente en una solicitud por recursos — al analizar estas tablas como grafos dirigidos, se representará:

- Los procesos, con cuadrados.
- Los recursos, con círculos.
 - Puede representarse como un círculo grande a una *clase* o *categoría* de recursos, y como círculos pequeños dentro de éste a una *serie de recursos idénticos* (p. ej. las diversas unidades de cinta)

- Las flechas que van de un recurso a un proceso indican que el recurso *está asignado* al proceso
- Las flechas que van de un proceso a un recurso indican que el proceso *solicita* el recurso

Cabe mencionar en este momento que, cuando se consideran categorías de recursos, el tener la representación visual de un ciclo no implica que haya ocurrido un bloqueo — este sólo se presenta cuando todos los procesos involucrados están en espera mutua.

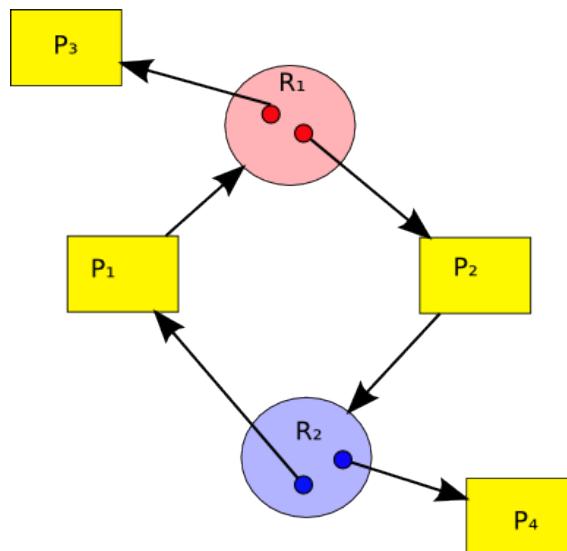


Figura 4.9: Al emplear categorías de recursos, un ciclo *no necesariamente* indica un bloqueo

En la figura 4.9, si bien P_1 y P_2 están esperando que se liberen recursos de tipo R_1 y R_2 , P_3 y P_4 siguen operando normalmente, y es esperable que lleguen a liberar el recurso por el cual están esperando. En el caso ilustrado, dado que el bloqueo se presenta únicamente al ser imposible que un proceso *libere* recursos que *ya le fueron asignados*, tendría que presentarse un caso donde todos los recursos de una misma categoría estuvieran involucrados en una situación de espera circular, como la ilustrada a continuación.

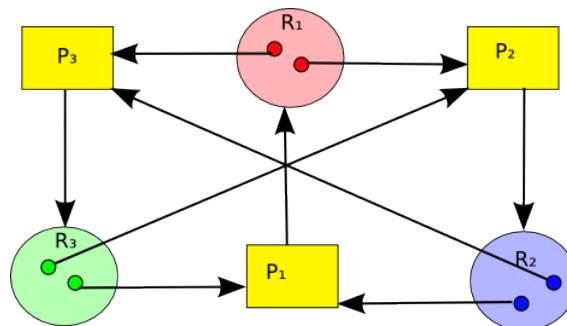


Figura 4.10: Situación en que se presenta espera circular, incluso empleando categorías de recursos

Si se tiene una representación completa de los procesos y recursos en el sistema, la estrategia es *reducir* la gráfica retirando los elementos que no brinden información imprescindible, siguiendo la siguiente lógica (recordar que representan una fotografía del sistema *en un momento dado*):

- Se retiran los procesos que no están solicitando ni tienen asignado ningún recurso.
- Para todos los procesos restantes: si todos los recursos que están solicitando *pueden ser concedidos* (esto es, no están actualmente asignados a otro), se reduce eliminando

del grafo al proceso y a todas las flechas relacionadas con éste.

- Si después de esta reducción se eliminan todos los procesos del grafo, entonces no hay interbloqueos y se puede continuar. En caso de permanecer procesos en el grafo, los procesos “irreducibles” constituyen la serie de procesos interbloqueados de la gráfica.

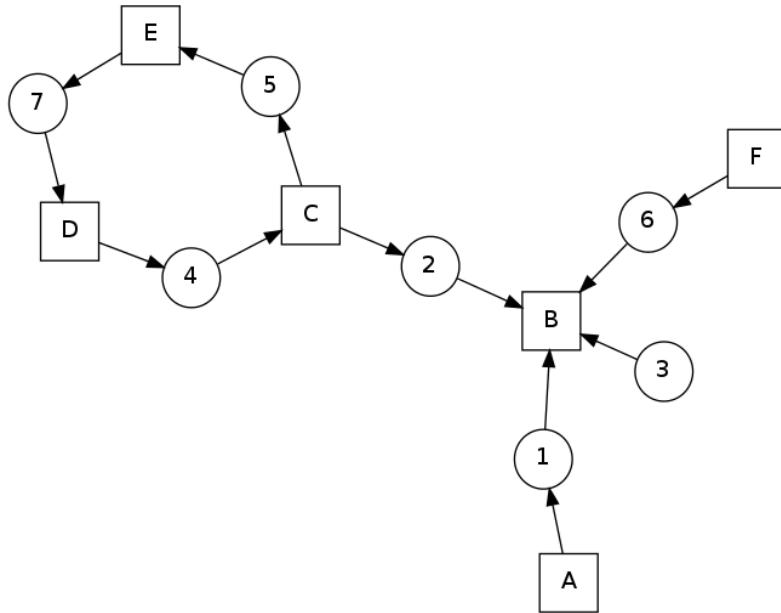


Figura 4.11: Detección de ciclos denotando bloqueos: Grafo de procesos y recursos en un momento dado

En la gráfica *Detección de ciclos denotando bloqueos*, se procede así:

- Se reduce por *B*, dado que actualmente no está esperando a ningún recurso.
- Se reduce por *A* y *F*, dado que los recursos por los cuales están esperando quedan libres en ausencia de *B*.

Y queda un interbloqueo entre *C*, *D* y *E*, en torno a los recursos 4, 5 y 7.

Nótese que *reducir* un proceso del grafo no implica que éste haya *entregado sus recursos*, sino únicamente que, hasta donde se tiene conocimiento, *tiene posibilidad de hacerlo*. Los procesos que están esperando por recursos retenidos por un proceso pueden sufrir inanición aún por un tiempo indeterminado.

Una vez que un bloqueo es diagnosticado, dado que los procesos no podrán terminar por sí mismos (pues están precisamente bloqueados, su ejecución no avanzará más), hay varias estrategias para la recuperación:

- Terminar a todos los procesos bloqueados. Esta es la técnica más sencilla y, de cierto modo, más justa — Todos los procesos implicados en el bloqueo pueden ser relanzados, pero todo el estado del cómputo que han realizado hasta este momento se perderá.
- Retroceder a los procesos implicados hasta el último *punto de control* (*checkpoint*) seguro conocido. Esto es posible únicamente cuando el sistema implementa esta funcionalidad, que tiene un elevado costo adicional. Cuando el estado de uno de los procesos depende de factores externos a éste, es imposible implementar fielmente los *puntos de control*.

Podría parecer que retroceder a un punto previo llevaría indefectiblemente a que se repita la situación — pero los bloqueos mutuos requieren de un orden de ejecución específico para aparecer. Muy probablemente, una ejecución posterior logrará salvar el bloqueo — y en caso contrario, puede repetirse este paso.

- Terminar, uno por uno y no en bloque, a cada uno de los procesos bloqueados. Una vez que se termina uno, se evalúa la situación para verificar si logró romperse la situación de bloqueo, en cuyo caso la ejecución de los restantes continúa sin interrupción.

Para esto, si bien podría elegirse un proceso al azar de entre los bloqueados, típicamente se consideran elementos adicionales como:

- Los procesos que demandan garantías de *tiempo real* son los más sensibles para detener y relanzar
- La menor cantidad de tiempo de procesador consumido hasta el momento. Dado que el proceso probablemente tenga que ser re-lanzado (re-ejecutado), puede ser conveniente *apostarle* a un proceso que haya hecho poco cálculo (para que el tiempo que tenga que invertir para volver al punto actual sea mínimo).
- Mayor tiempo restante estimado. Si se puede estimar cuánto tiempo de procesamiento *queda pendiente*, conviene terminar al proceso que más le falte por hacer.
- Menor número de recursos asignados hasta el momento. Un poco como criterio de justicia, y un poco partiendo de que es un proceso que está haciendo menor uso del sistema.
- Prioridad más baja. Cuando hay un ordenamiento de procesos o usuarios por prioridades, siempre es preferible terminar un proceso de menor prioridad o perteneciente a un usuario poco importante que uno de mayor prioridad.
- En caso de contar con la información necesaria, es siempre mejor interrumpir un proceso que *pueda ser repetido sin pérdida de información* que uno que la cause. Por ejemplo, es preferible interrumpir una compilación que la actualización de una base de datos.

Un punto importante a considerar es cada cuánto debe realizarse la verificación de bloqueos. Podría hacerse:

- Cada vez que un proceso solicite un recurso. pero esto llevaría a un gasto de tiempo en este análisis demasiado frecuente.
- Con una periodicidad fija, pero esto arriesga a que los procesos pasen más tiempo bloqueados.
- Cuando el nivel del uso del CPU baje de cierto porcentaje. Esto indicaría que hay un nivel elevado de procesos en espera.
- Una estrategia combinada.

Por último, si bien los dispositivos aquí mencionados requieren bloqueo exclusivo, otra estrategia es la *apropiación temporal*: tomar un recurso asignado a determinado proceso para otorgárselo *temporalmente* a otro. Esto no siempre es posible, claro, y depende fuertemente de la naturaleza del mismo — pero podría, por ejemplo, interrumpirse un proceso que tiene asignada (pero inactiva) a una impresora para otorgársela temporalmente a otro que tiene un trabajo corto pendiente. Esto último, sin embargo, es tan sensible a detalles de cada clase de recursos que rara vez puede hacerlo el sistema operativo — es normalmente hecho *de acuerdo* entre los procesos competidores, por medio de algún protocolo pre-establecido.

4.4.4 Algoritmo del aveSTRUZ

Una cuarta línea (que, por increíble que parezca, es la más común, empleada en todos los sistemas operativos de propósito general) es el llamado *algoritmo del aveSTRUZ*: ignorar las situaciones de bloqueo (escondiéndose de ellas como aveSTRUZ que esconde la cabeza bajo la tierra), esperando que su ocurrencia sea suficientemente poco frecuente, o si ocurre, que su impacto no afecte al sistema.

Justificando a los avestruces

Hay que comprender que esto ocurre porque las condiciones impuestas por las demás estrategias resultan demasiado onerosas, el conocimiento previo resulta insuficiente, o los bloqueos simplemente pueden presentarse ante recursos externos y no controlados (o conocidos siquiera) por el sistema operativo.

Ignorar la posibilidad de un bloqueo *cuando su probabilidad es suficientemente baja* será preferible para los usuarios (y programadores) ante la disyuntiva de afrontar restricciones para la forma y conveniencia de solicitar recursos.

En este caso, se toma una decisión entre lo *correcto* y lo *conveniente* — Un sistema operativo formalmente no debería permitir la posibilidad de que hubiera bloqueos, pero la inconveniencia presentada al usuario sería inaceptable.

Por último, cabe mencionar algo que a lo largo de todo este apartado mencionamos únicamente de forma casual, evadiendo definiciones claras: ¿qué es un recurso? La realidad es que no está muy bien definido. Se podría, como mencionan los ejemplos, hablar de los clásicos recursos de acceso rival y secuencial: impresoras, cintas, terminales seriales, etc. Sin embargo, también se pueden ver como recursos a otras entidades administradas por el sistema operativo — el espacio disponible de memoria, el tiempo de procesamiento, o incluso estructuras lógicas *creadas y gestionadas* por el sistema operativo, como archivos, semáforos o monitores. Y para esos casos, prácticamente ninguno de los mecanismos aquí analizados resolvería las características de acceso y bloqueo necesarias.

Enfrentando a los avestruces

La realidad del cómputo marca que es el programador de aplicaciones quien debe prever las situaciones de carrera, bloqueo e inanición en su código — El sistema operativo empleará ciertos mecanismos para asegurar la seguridad en general entre los componentes del sistema, pero el resto recae en las manos del programador.

Una posible salida ante la presencia del *algoritmo del avestruz* es adoptar un método *defensivo* de programar. Un ejemplo de esto sería que los programas soliciten un recurso pero, en vez de solicitarlo por medio de una *llamada bloqueante*, hacerlo por medio de una *llamada no bloqueante* y, en caso de fallar ésta, esperar un tiempo aleatorio e intentar nuevamente acceder al recurso un número dado de veces, y, tras n intentos, abortar limpiamente el proceso y notificar al usuario (evitando un bloqueo mutuo circular indefinido).

Por otro lado, hay una gran cantidad de aplicaciones de monitoreo en espacio de usuario. Conociendo el funcionamiento esperable específico de determinados programas es posible construir aplicaciones que los monitoreen *de una forma inteligente* y tomen acciones (ya sea alertar a los administradores o, como se lo revisa en la sección 4.4.3 (*Detección y recuperación de bloqueos*), abortar -y posiblemente reiniciar- la ejecución de aquellos procesos que no puedan recuperarse).

De avestruces, ingenieros y matemáticos

Esta misma contraposición puede leerse, hasta cierto punto en tono de broma, como un síntoma de la tensión que caracteriza a nuestra profesión: la computación nació como ciencia dentro de los departamentos de matemáticas en diversas facultades, sin embargo, al pasar de los años ha ido integrando cada vez más a la ingeniería. Y el campo, una de las áreas más jóvenes pero al mismo tiempo más prolíficas del conocimiento humano, está en una constante discusión y definición: ¿Qué somos? ¿Matemáticos, ingenieros, o... alguna otra cosa?

La asignación de recursos, pues, puede verse desde el punto de vista matemático: es un problema con un planteamiento de origen, y hay varias estrategias distintas (los mecanismos y algoritmos descritos en esta sección). Pueden no ser perfectos, pero el problema no ha *demonstrado* ser intratable. Y un bloqueo es claramente un error — una

situación de excepción, inaceptable. Los matemáticos en nuestro árbol genealógico académico nos llaman a no ignorar este problema, a resolverlo sin importar la complejidad computacional.

Los ingenieros, más aterrizados en el *mundo real*, tienen como parte básica de su formación, sí, el evitar defectos nocivos — pero también contemplan el cálculo de costos, la probabilidad de impacto, los umbrales de tolerancia... Para un ingeniero, si un sistema típico *corre riesgo* de caer en un bloqueo mutuo con una probabilidad $p > 0$, dejando inservibles a dos procesos en un sistema, pero debe también considerar no sólo las fallas en hardware y en los diferentes componentes del sistema operativo, sino que en todos los demás programas que ejecutan en espacio de usuario, y considerando que prevenir el bloqueo conlleva un costo adicional en complejidad para el desarrollo o en rendimiento del sistema (dado que perderá tiempo llevando a cabo las verificaciones ante cualquier nueva solicitud de recursos), no debe sorprender a nadie que los ingenieros se inclinen por adoptar la estrategia del aveSTRUZ — claro está, siempre que no haya opción *razonable*.

4.5 Otros recursos

- *Tutorial de hilos de Perl*
<http://perldoc.perl.org/perlthrtut.html>
John Orwant (1998); The Perl Journal
- *Python higher level threading interface*
<http://docs.python.org/2/library/threading.html>
Python Software Foundation (1990-2014); Python 2.7.6 Documentation
- *Spin Locks & Other Forms of Mutual Exclusion*
<http://www.cs.fsu.edu/~baker/devices/notes/spinlock.html>
Theodore P. Baker (2010); Florida State University
- *Dining philosophers revisited*
<https://dl.acm.org/citation.cfm?id=101091>
Armando R. Gingras (1990), ACM SIGCSE Bulletin

5 — Planificación de procesos

5.1 Tipos de planificación

La *planificación de procesos* se refiere a cómo determina el sistema operativo al orden en que irá cediendo el uso del procesador a los procesos que lo vayan solicitando, y a las políticas que empleará para que el uso que den a dicho tiempo no sea excesivo respecto al uso esperado del sistema.

Existen tres tipos principales de planificación:

A largo plazo Decide qué procesos serán los siguientes en ser iniciados. Este tipo de planificación era el más frecuente en los sistemas de lotes (principalmente aquellos con *spool*) y multiprogramados en lotes; las decisiones eran tomadas principalmente considerando los requisitos pre-declarados de los procesos y los que el sistema tenía libres al terminar algún otro proceso. La planificación a largo plazo puede llevarse a cabo con periodicidad de una vez cada varios segundos, minutos e inclusive horas. En los sistemas de uso interactivo, casi la totalidad de los que se usan hoy en día, este tipo de planificación no se efectúa, dado que es típicamente el usuario quien indica expresamente qué procesos iniciar.

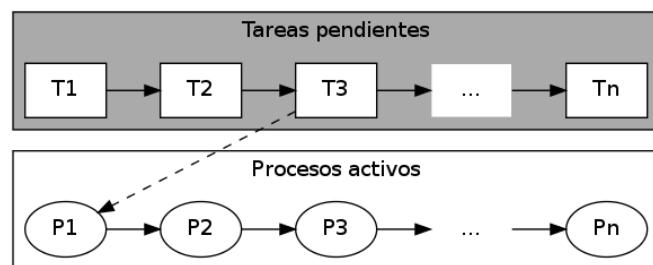


Figura 5.1: Planificador a largo plazo

A mediano plazo Decide cuáles procesos es conveniente *bloquear* en determinado momento, sea por escasez/saturación de algún recurso (como la memoria primaria) o porque están realizando alguna solicitud que no puede satisfacerse momentáneamente; se encarga de tomar decisiones respecto a los procesos conforme entran y salen del estado de *bloqueado* (esto es, típicamente, están a la espera de algún evento externo o de la finalización de transferencia de datos con algún dispositivo).

En algunos textos, al *planificador a mediano plazo* se le llama *agendador (scheduler)*.

A corto plazo Decide cómo compartir *momento a momento* al equipo entre todos los procesos que requieren de sus recursos, especialmente el procesador. La planificación a corto plazo se lleva a cabo decenas de veces por segundo (razón por la cual debe ser código muy simple, eficiente y rápido); es el encargado de planificar *los procesos que están listos para ejecución*.

El *planificador a corto plazo* es también frecuentemente denominado *despachador (dispatcher)*.

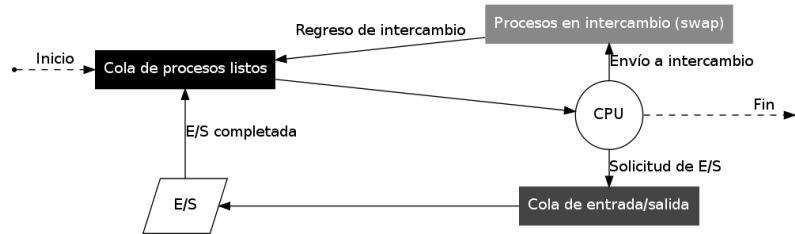


Figura 5.2: Planificador a mediano plazo, o *agendador*

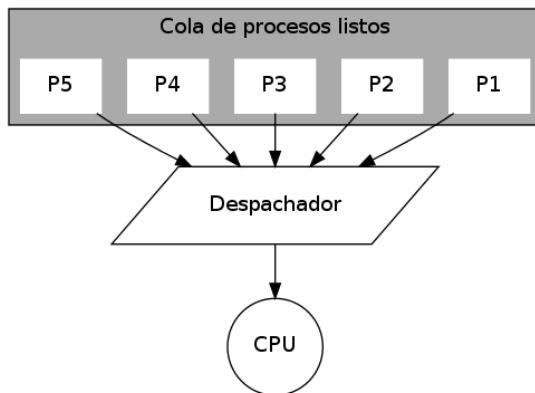


Figura 5.3: Planificador a corto plazo, o *despachador*

Relacionando con los estados de un proceso abordados en la sección 4.1.1, y volviendo al diagrama entonces presentado (reproducido por comodidad de referencia en la figura 5.4), podrían ubicarse a estos tres planificadores en las siguientes transiciones entre estados:

1. El planificador a largo plazo se encarga de *admitir* un nuevo proceso: La transición de *Nuevo* a *Listo*.
2. El planificador a mediano plazo maneja la *activación* y *bloqueo* de un proceso relacionado con *eventos* — Esto es, las transiciones entre *En ejecución* y *Bloqueado*, y entre *Bloqueado* y *Listo*.
3. El planificador a corto plazo decide entre los procesos que están listos para ejecutarse y determina a cuál de ellos *activar*, y detiene a aquellos que *exceden su tiempo* de procesador — Implementa las transiciones entre los estados *Listo* y *En ejecución*.

En esta sección se trata particularmente el planificador a corto plazo, haciendo referencia como mucho a algunos efectos del planificador a mediano plazo.

5.1.1 Tipos de proceso

Como ya se ha visto, los procesos típicamente alternan entre *ráfagas* (periodos, en inglés *bursts*) en que realizan principalmente cómputo interno (están *limitados por CPU*, *CPU-bound*) y otras en que la atención está puesta en transmitir los datos desde o hacia dispositivos externos (están *limitados por entrada-salida*, *I/O-bound*). Dado que cuando un proceso se suspende para realizar entrada-salida deja de estar *listo* (y pasa a estar *bloqueado*), y desaparece de la atención del planificador a corto plazo, en todo momento los procesos que están en ejecución y listos pueden separarse en:

Procesos largos Aquellos que *por mucho tiempo*¹ han estado en *listos* o en ejecución, esto es, procesos que estén en una larga ráfaga limitada por CPU.

Procesos cortos Aquellos que, ya sea que en *este momento*² estén en una ráfaga limitada

¹¿Cuánto es mucho? Dependerá de las políticas generales que se definan para el sistema

²Y también, *este momento* debe ser interpretado con la granularidad acorde al sistema

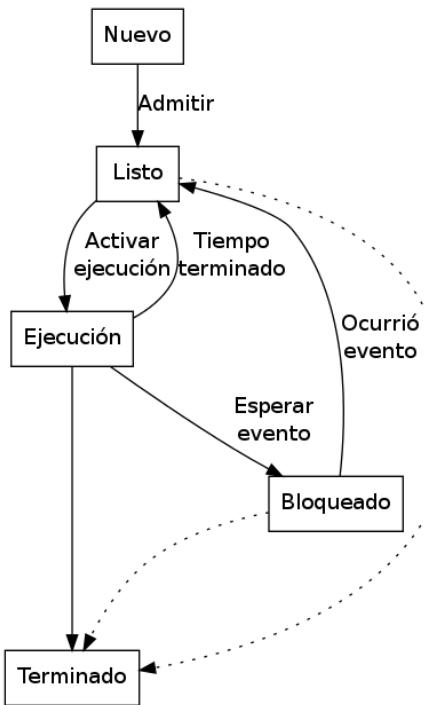


Figura 5.4: Diagrama de transición entre los estados de un proceso

por entrada-salida y requieran atención meramente ocasional del procesador, o tienden a estar bloqueados esperando a eventos (como los procesos interactivos).

Por lo general se busca dar un tratamiento *preferente* a los procesos cortos, en particular a los interactivos. Cuando un usuario está interactuando con un proceso, si no tiene una respuesta *inmediata* a su interacción con el equipo (sea proporcionar comandos, recibir la respuesta a un *teclazo* o mover el puntero en el GUI) su percepción será la de una respuesta degradada.

5.1.2 Midiendo la respuesta

Resulta intuitivo que cada patrón de uso del sistema debe seguir políticas de planificación distintas. Por ejemplo, en el caso de un proceso interactivo, se buscará ubicar al proceso en una *cola* preferente (para obtener un tiempo de respuesta más ágil, para mejorar la percepción del usuario), pero en caso de sufrir demoras, es preferible buscar dar una respuesta *consistente*, aún si la respuesta *promedio* es más lenta. Esto es, si a todas las operaciones sigue una demora de un segundo, el usuario sentirá menos falta de control si en promedio tardan medio segundo, pero ocasionalmente hay picos de cinco.

Para este tema, en vez de emplear unidades temporales formales (p. ej. fracciones de segundo), es común emplear *ticks* y *quantums*. Esto es en buena medida porque, si bien en el campo del cómputo las velocidades de acceso y uso efectivo cambian constantemente, los conceptos y las definiciones permanecen. Además, al ser ambos parámetros ajustables, una misma implementación puede sobrevivir ajustándose a la evolución del hardware.

Tick Una fracción de tiempo durante la cual se puede realizar trabajo útil - Esto es, usar la CPU sin interrupción³. El tiempo correspondiente a un tick está determinado por una señal (interrupción) periódica, emitida por el *temporizador* (timer). La frecuencia con que ocurre esta señal se establece al inicio del sistema. Por ejemplo, una frecuencia de *temporizador* de 100 Hertz implica que éste emitirá una señal cada 10 milisegundos.

³Ignorando las interrupciones causadas por los dispositivos de entrada y salida y otras señales que llegan a la CPU

En Linux (a partir de la versión 2.6.8), un *tick* dura un milisegundo, en Windows, entre 10 y 15 milisegundos.

Quantum El tiempo mínimo que se permitirá a un proceso el uso del procesador. En Windows, dependiendo de la clase de proceso que se trate, un *quantum* durará entre 2 y 12 ticks (esto es, entre 20 y 180 ms), y en Linux, entre 10 y 200 ticks (10 y 200 milisegundos respectivamente).

¿Qué mecanismos o métricas se emplean para medir el comportamiento del sistema bajo determinado planificador? Partiendo de los siguientes conceptos, para un proceso p que requiere de un tiempo t de ejecución:

Tiempo de respuesta (T) Cuánto tiempo total es necesario para completar el trabajo pendiente de un proceso p , incluyendo el tiempo que está inactivo esperando ejecución (pero está en la cola de procesos listos).

Tiempo en espera ($E = T - t$) También referido como *tiempo perdido*. Del tiempo de respuesta total, cuánto tiempo p está listo y esperando ejecutar. Desde la óptica de p , se desearía que $E_p \rightarrow 0$

Proporción de penalización ($P = \frac{T}{t}$) Fracción del tiempo de respuesta durante la cual p estuvo en espera.

Proporción de respuesta ($R = \frac{t}{T}$) Inverso de P . Fracción del tiempo de respuesta durante la cual p pudo ejecutarse.

Para hacer referencia a un grupo de procesos con requisitos similares, todos ellos requiriendo de un mismo tiempo t , se emplea $T(t)$, $E(t) = T(t) - t$, $P(t) = \frac{T(t)}{t}$ y $R(t) = \frac{t}{T(t)}$.

Además de estos tiempos, expresados en relación al tiempo efectivo de los diversos procesos del sistema, es necesario considerar también:

Tiempo núcleo o *kernel* Tiempo que pasa el sistema en espacio de núcleo, incluyendo entre otras funciones⁴ el empleado en decidir e implementar la política de planificación y los cambios de contexto.

Tiempo desocupado (*idle*) Tiempo en que la cola de procesos listos está vacía y no puede realizarse ningún trabajo.

Utilización del CPU Porcentaje del tiempo en que el CPU está realizando *trabajo útil*. Si bien conceptualmente puede ubicarse dicha utilización entre 0 y 100%, en sistemas reales se ha observado (Silberschatz, p.179) que se ubica en un rango entre el 40 y el 90%.

Por ejemplo, si llegan a la cola de procesos listos:

Proceso	Ticks	Llegada
A	7	0
B	3	2
C	12	6
D	4	20

Si el tiempo que toma al sistema efectuar un cambio de contexto es de un *tick*, y la duración de cada *quantum* es de 5 ticks, en un ordenamiento de ronda,⁵ se observaría un resultado como el que ilustra la figura 5.5.

Al considerar al tiempo ocupado por el núcleo como un proceso más, cuyo trabajo en este espacio de tiempo finalizó junto con los demás,⁶ se obtiene por resultado:

⁴Estas funciones incluyen principalmente la atención a interrupciones, el servicio a llamadas al sistema, y cubrir diversas tareas administrativas.

⁵Este mecanismo se presentará en breve, en la sección 5.2.3.

⁶Normalmente no se considera al núcleo al hacer este cálculo, dado que en este ámbito todo el trabajo que hace puede verse como *burocracia* ante los resultados deseados del sistema

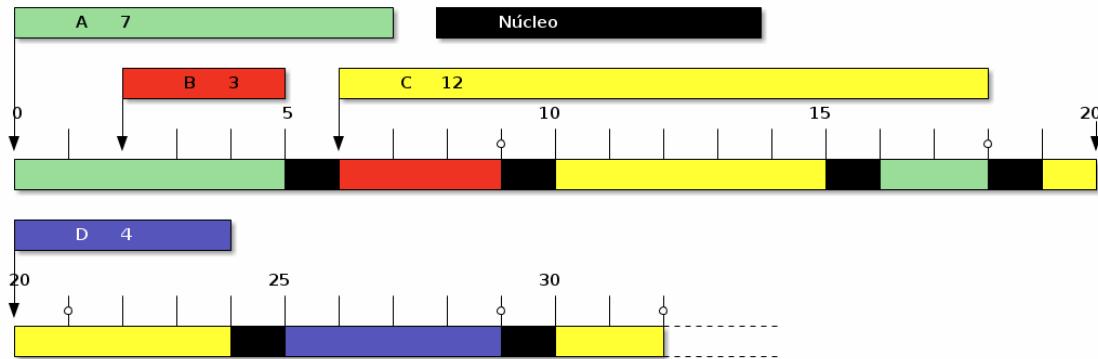


Figura 5.5: Ejecución de cuatro procesos con *quantums* de 5 ticks y cambios de contexto de 2 ticks

Proceso	<i>t</i>	<i>T</i>	<i>E</i>	<i>P</i>	<i>R</i>
A	7	18	11	2.57	0.389
B	3	7	4	2.33	0.429
C	12	26	14	2.17	0.462
D	4	9	5	2.25	0.444
Promedio útil	6.5	15	8.50	2.31	0.433
Núcleo	6	32	26	5.33	0.188
Promedio total	6.4	18.4	12.00	2.88	0.348

Abordando cada proceso, para obtener *T* se parte del momento en que el proceso llegó a la cola, no el punto de inicio de la línea de tiempo. En este caso, dado que el núcleo *siempre* está en ejecución, se asume que inició también en 0.

Respecto al patrón de llegada y salida de los procesos, lo se maneja también basado en una relación. Partiendo de una *frecuencia de llegada* promedio de nuevos procesos a la cola de procesos listos α , y el *tiempo de servicio requerido* promedio β , se define el *valor de saturación* ρ como $\rho = \frac{\alpha}{\beta}$.

Cuando $\rho = 0$, nunca llegan nuevos procesos, por lo cual el sistema estará eventualmente *desocupado*. Cuando $\rho = 1$, los procesos son despachados al mismo ritmo al que van llegando. Cuando $\rho > 1$, el ritmo de llegada de procesos es mayor que la velocidad a la cual la computadora puede darles servicio, con lo cual la cola de procesos listos tenderá a crecer (y la calidad de servicio, la proporción de respuesta R , para cada proceso se decrementará).

5.2 Algoritmos de planificación

El planificador a corto plazo puede ser invocado cuando un proceso se encuentra en algunas de las cuatro siguientes circunstancias:

1. Pasa de estar *ejecutando* a estar *en espera* (por ejemplo, por solicitar una operación de E/S, esperar a la sincronización con otro proceso, etc.)
2. Pasa de estar *ejecutando* a estar *listo* (por ejemplo, al ocurrir la interrupción del temporizador, o de algún evento externo)
3. Deja de estar *en espera* a estar *listo* (por ejemplo, al finalizar la operación de E/S que solicitó)
4. Finaliza su ejecución, y pasa de *ejecutando* a *terminado*

En el primer y cuarto casos, el sistema operativo siempre tomará el control⁷; un sistema que opera bajo *multitarea preventiva* implementará también el segundo y tercer casos, mientras que uno que opera bajo *multitarea cooperativa* no necesariamente reconocerá dichos estados.

Ahora, para los algoritmos a continuación, cabe recordar que se trata únicamente del *despachador*. Un proceso siempre abandonará la cola de procesos listos al requerir de un servicio del sistema.

Para todos los ejemplos a continuación, los tiempos están dados en *ticks*; no es relevante a cuánto *tiempo de reloj* estos equivalen, sino el rendimiento relativo del sistema entero ante una carga dada.

La presente sección está basada fuertemente en el capítulo 2 de *An operating systems vade mecum* (Raphael Finkel, 1988).

5.2.1 Objetivos de la planificación

Los algoritmos que serán presentados a continuación son respuestas que intentan, de diferentes maneras y desde distintos supuestos base, darse a los siguientes objetivos principales (tomando en cuenta que algunos de estos objetivos pueden ser mutuamente contradictorios):

Ser justo Debe tratarse de *igual manera* a todos los procesos que comparten las mismas características⁸, y nunca postergar indefinidamente a uno de ellos.

Maximizar el rendimiento Dar servicio a la mayor parte de procesos por unidad de tiempo.

Ser predecible Un mismo trabajo debe tomar aproximadamente la misma cantidad de tiempo en completarse independientemente de la carga del sistema.

Minimizar la sobrecarga El tiempo que el algoritmo pierda en *burocracia* debe mantenerse al mínimo, dado que éste es tiempo de procesamiento útil perdido

Equilibrar el uso de recursos Favorecer a los procesos que empleen recursos subutilizados, penalizar a los que peleen por un recurso sobreutilizado causando contención en el sistema

Evitar la postergación indefinida Aumentar la prioridad de los procesos más *viejos*, para favorecer que alcancen a obtener algún recurso por el cual estén esperando

Favorecer el uso esperado del sistema En un sistema con usuarios interactivos, maximizar la prioridad de los procesos que sirvan a solicitudes iniciadas por éste (aún a cambio de penalizar a los procesos /de sistema)

Dar preferencia a los procesos que podrían causar bloqueo Si un proceso de baja prioridad está empleando un recurso del sistema por el cual más procesos están esperando, favorecer que éste termine de emplearlo más rápido

Favorecer a los procesos con un comportamiento deseable Si un proceso causa muchas demoras (por ejemplo, atraviesa una ráfaga de entrada/salida que le requiere hacer muchas llamadas a sistema o interrupciones), se le puede penalizar porque degrada el rendimiento global del sistema

Degradarse suavemente Si bien el nivel ideal de utilización del procesador es al 100%, es imposible mantenerse siempre a este nivel. Un algoritmo puede buscar responder con la menor penalización a los procesos preexistentes al momento de exceder este umbral.

⁷En el primer caso, el proceso entrará en el dominio del *planificador a mediano plazo*, mientras que en el cuarto saldrá definitivamente de la lista de ejecución.

⁸Un algoritmo de planificación puede *priorizar* de diferente manera a los procesos según distintos criterios, sin por ello dejar de ser justo, siempre que dé la misma prioridad y respuesta a procesos equivalentes.

5.2.2 Primero llegado, primero servido (FCFS)

El esquema más simple de planificación es el *Primero llegado, primero servido* (*First come, first serve, FCFS*). Este es un mecanismo cooperativo, con la mínima lógica posible: Cada proceso se ejecuta en el orden en que fue llegando, y hasta que suelta el control. El despachador es muy simple, básicamente una cola FIFO.

Para comparar los distintos algoritmos de planificación que serán presentados, se presentará el resultado de cada uno de ellos sobre el siguiente juego de procesos: (Finkel 1988, p.35)

Proceso	Tiempo de Llegada	t	Inicio	Fin	T	E	P
A	0	3	0	3	3	0	1
B	1	5	3	8	7	2	1.4
C	3	2	8	10	7	5	3.5
D	9	5	10	15	6	1	1.2
E	12	5	15	20	8	3	1.6
Promedio		4			6.2	2.2	1.74

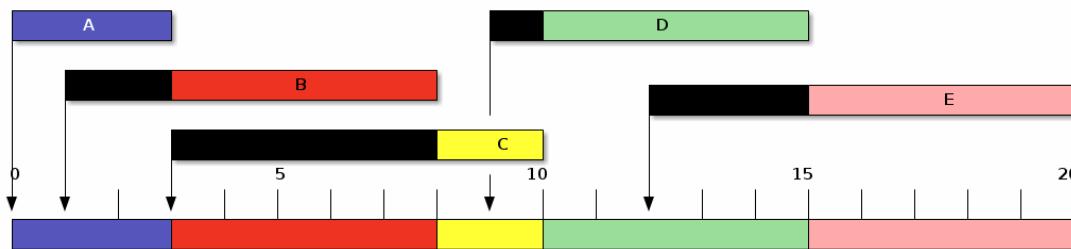


Figura 5.6: Primero llegado, primero servido (FCFS)

Si bien un esquema FCFS reduce al mínimo la *sobrecarga administrativa* (que incluye tanto al tiempo requerido por el planificador para seleccionar al siguiente proceso como el tiempo requerido para el cambio de contexto), el rendimiento percibido por los últimos procesos en llegar (o por procesos cortos llegados en un momento inconveniente) resulta inaceptable.

Este algoritmo dará servicio y salida a todos los procesos siempre que $\rho \leq 1$. En caso de que se sostenga $\rho > 1$, la demora para iniciar la atención de un proceso crecerá cada vez más, cayendo en una cada vez mayor inanición.

FCFS tiene características claramente inadecuadas para trabajo interactivo, sin embargo, al no requerir de hardware de apoyo (como un temporizador) sigue siendo ampliamente empleado.

5.2.3 Ronda (Round Robin)

El esquema *ronda* busca dar una relación de respuesta buena tanto para procesos largos como para los cortos. La principal diferencia entre la ronda y FCFS es que en este caso sí emplea multitarea preventiva: Cada proceso que esté en la lista de procesos listos puede ejecutarse por un sólo *quantum* (q). Si un proceso no ha terminado de ejecutar al final de su *quantum*, será interrumpido y puesto al final de la lista de procesos listos, para que espere a su turno nuevamente. Los procesos que sean *despertados* por los planificadores a mediano o largo plazo se agregarán también al final de esta lista.

Con la misma tabla de procesos presentada en el caso anterior (y, por ahora, ignorando la sobrecarga administrativa provocada por los cambios de contexto) se obtienen los siguientes resultados:

Proceso	Tiempo de Llegada	t	Inicio	Fin	T	E	P
A		0	3	0	6	3	2.0
B		1	5	1	11	10	2.0
C		3	2	4	8	5	2.5
D		9	5	9	18	9	4.0
E		12	5	12	20	8	3.0
Promedio		4			7.6	3.6	1.98

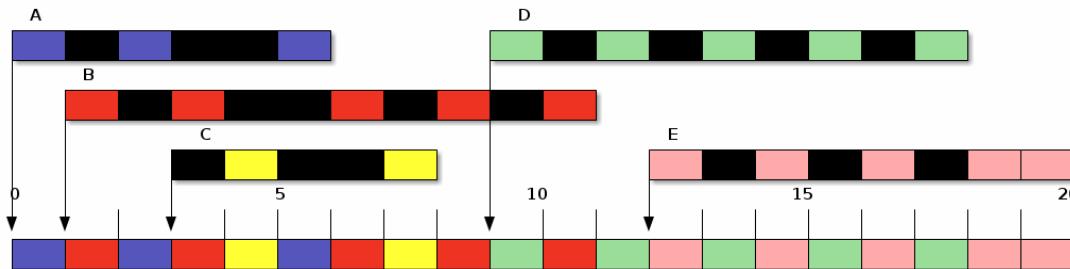


Figura 5.7: Ronda (Round Robin)

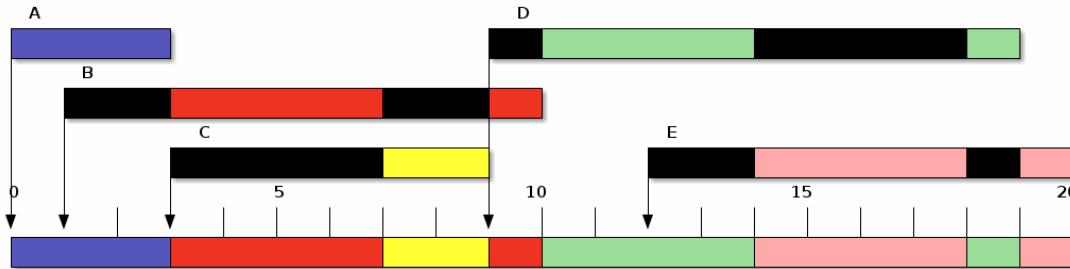
La *ronda* puede ser ajustada modificando la duración de q . Conforme se incrementa q , la ronda tiende a convertirse en FCFS — Si cada *quantum* es arbitrariamente grande, todo proceso terminará su ejecución dentro de su *quantum*; por otro lado, conforme decrece q , se tiene una mayor frecuencia de cambios de contexto; esto llevaría a una mayor ilusión de tener un procesador dedicado por parte de cada uno de los procesos, dado que cada proceso sería incapaz de notar las *ráfagas* de atención que éste le da (avance rápido durante un periodo corto seguido de un periodo sin avance). Claro está, el procesador simulado sería cada vez más lento, dada la fuerte penalización que iría agregando la sobrecarga administrativa.

Finkel (1988, p.35) se refiere a esto como el *principio de la histéresis*: *Hay que resistirse al cambio*. Como ya lo se mencionó, FCFS mantiene al mínimo posible la sobrecarga administrativa, y –aunque sea marginalmente– resulta en mejor rendimiento global.

Si se repite el análisis anterior bajo este mismo mecanismo, pero con un *quantum* de 4 ticks, el resultado es:

Proceso	Tiempo de Llegada	t	Inicio	Fin	T	E	P
A		0	3	0	3	3	0
B		1	5	3	10	9	4
C		3	2	7	9	6	4
D		9	5	10	19	10	5
E		12	5	14	20	8	3
Promedio		4			7.2	3.2	1.88

Si bien aumentar el *quantum* mejora los tiempos promedio de respuesta, aumentarlo hasta convertirlo en un FCFS efectivo degenera en una penalización a los procesos cortos, y puede llevar a la inanición cuando $\rho > 1$. Silberschatz apunta (p.188) a que típicamente el *quantum* debe mantenerse inferior a la duración promedio del 80% de los procesos.

Figura 5.8: Ronda (Round Robin), con $q = 4$

5.2.4 El proceso más corto a continuación (SPN)

(Del inglés, *Shortest Process Next*)

Cuando no se tiene la posibilidad de implementar multitarea preventiva, pero se requiere de un algoritmo más *justo*, contando con información *por anticipado* acerca del tiempo que requieren los procesos que forman la lista, puede elegirse el más corto de los presentes.

Ahora bien, es muy difícil contar con esta información antes de ejecutar el proceso. Es más frecuente buscar *caracterizar* las necesidades del proceso: Ver si durante su historia de ejecución⁹ ha sido un proceso tendiente a manejar ráfagas *limitadas por entrada-salida* o *limitadas por procesador*, y cuál es su tendencia actual.

Para estimar el tiempo que requerirá un proceso p en su próxima invocación, es común emplear el *promedio exponencial* e_p . Se define un *factor atenuante* $0 \leq f \leq 1$, que determinará qué tan reactivo será el promedio obtenido a la última duración; es común que este valor sea cercano a 0.9.

Si el p empleó q *quantums* durante su última invocación,

$$e'_p = fe_p + (1 - f)q$$

Se puede tomar como *semilla* para el e_p inicial un número elegido arbitrariamente, o uno que ilustre el comportamiento actual del sistema (como el promedio del e_p de los procesos actualmente en ejecución). La figura 5.10 presenta la predicción de tiempo requerido que determinado proceso va obteniendo en sus diversas entradas a la cola de ejecución, basado en su comportamiento previo, con distintos factores atenuantes.

Empleando el mismo juego de datos de procesos que se ha venido manejando como resultados de las estimaciones, se obtiene el siguiente resultado:

Proceso	Tiempo de Llegada	t	Inicio	Fin	T	E	P
A	0	3	0	3	3	0	1.0
B	1	5	5	10	9	4	1.8
C	3	2	3	5	2	0	1.0
D	9	5	10	15	6	1	1.2
E	12	5	15	20	8	3	1.6
Promedio		4			5.6	1.6	1.32

Como era de esperarse, SPN favorece a los procesos cortos. Sin embargo, un proceso largo puede esperar mucho tiempo antes de ser atendido, especialmente con valores de ρ

⁹Cabe recordar que todos estos mecanismos se aplican al *planificador a corto plazo*. Cuando un proceso se bloquea esperando una operación de E/S, sigue en ejecución, y la información de contabilidad del mismo sigue alimentándose. SPN se “nutre” precisamente de dicha información de contabilidad.

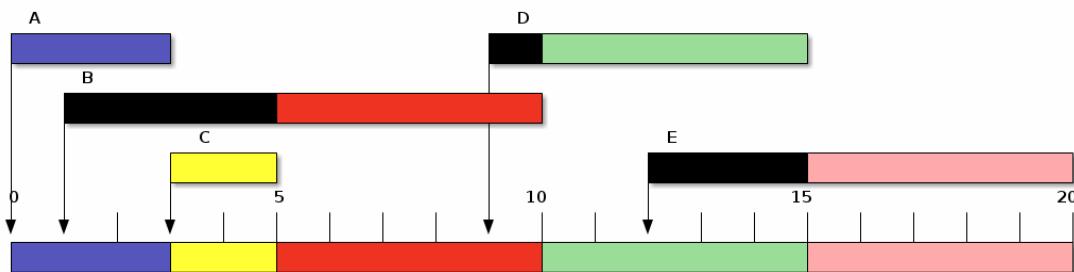


Figura 5.9: El proceso más corto a continuación (SPN)

cercanos o superiores a 1 — Un proceso más largo que el promedio está predispuesto a sufrir inanición.

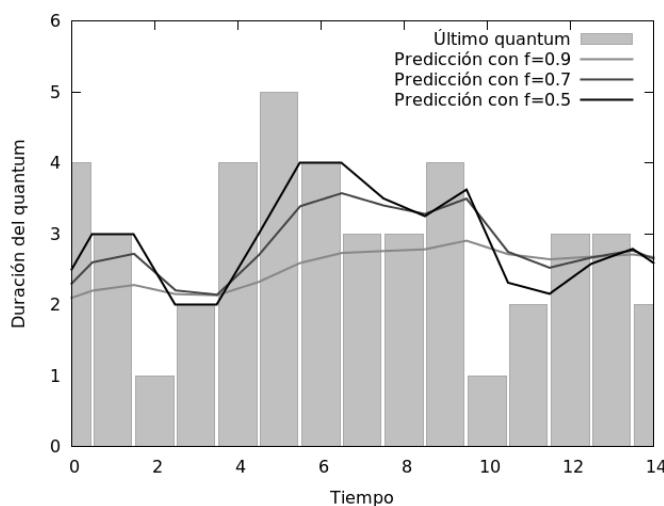


Figura 5.10: Promedio exponencial (predicción de próxima solicitud de tiempo) de un proceso.

En un sistema poco ocupado, en que la cola de procesos listos es corta, SPN generará resultados muy similares a los de FCFS. Sin embargo, puede observarse en el ejemplo que con sólo una permutación en los cinco procesos ejemplos (B y C), los factores de penalización a los procesos ejemplo resultaron muy beneficiados.

SPN preventivo (PSPN)

(*Preemptive Shortest Process Next*)

Finkel (1988, p.44) apunta a que, a pesar de que intuitivamente daría una mayor ganancia combinar las estrategias de SPN con un esquema de multitarea preventiva, el comportamiento obtenido es muy similar para la amplia mayoría de los procesos. Incluso para procesos muy largos, PSPN no los penaliza mucho más allá de lo que lo haría la ronda, y obtiene mejores promedios de forma consistente porque, al despachar primero a los procesos más cortos, mantiene la lista de procesos pendientes corta, lo que lleva naturalmente a menores índices de penalización.

El más penalizado a continuación (HPRN)

(*Highest Penalty Ratio Next*)

En un sistema que no cuenta con multitarea preventiva, las alternativas presentadas hasta ahora resultan invariabilmente injustas: FCFS favorece a los procesos largos, y SPN a los cortos. Un intento de llegar a un algoritmo más balanceado es HPRN.

Todo proceso inicia su paso por la cola de procesos listos con un valor de penalización $P = 1$. Cada vez que es obligado a esperar un tiempo w por otro proceso, P se actualiza como $P = \frac{w+t}{t}$. El proceso que se elige como activo será el que tenga mayor P . Mientras $\rho < 1$, HPRN evitará que incluso los procesos más largos sufran inanición.

En los experimentos realizados por Finkel, HPRN se sitúa siempre en un punto medio entre FCFS y SPN; su principal desventaja se presenta conforme crece la cola de procesos listos, ya que P tiene que calcularse para cada uno de ellos cada vez que el despachador toma una decisión.

5.2.5 Ronda egoísta (SRR)

(*Selfish Round Robin*)

Este método busca favorecer a los procesos que ya han pasado tiempo ejecutando que a los recién llegados. De hecho, los nuevos procesos no son programados directamente para su ejecución, sino que se les forma en la cola de *procesos nuevos*, y se avanza únicamente con la cola de *procesos aceptados*.

Para SRR se emplean los parámetros a y b , ajustables según las necesidades del sistema. a indica el ritmo según el cual se incrementará la prioridad de los procesos de la cola de *procesos nuevos*, y b el ritmo del incremento de prioridad para los *procesos aceptados*. Cuando la prioridad de un proceso nuevo *alcanza* a la prioridad de un proceso aceptado, el nuevo se vuelve aceptado. Si la cola de procesos aceptados queda vacía, se acepta el proceso nuevo con mayor prioridad.

El comportamiento de SRR con los procesos ejemplo es:

Proceso	Tiempo de Llegada	t	Inicio	Fin	T	E	P
A		0	3	0	4	1	1.3
B		1	5	2	10	4	1.8
C		3	2	6	9	6	4
D		9	5	10	15	6	1
E		12	5	15	20	3	1.2
Promedio		4			6.6	2.6	1.79

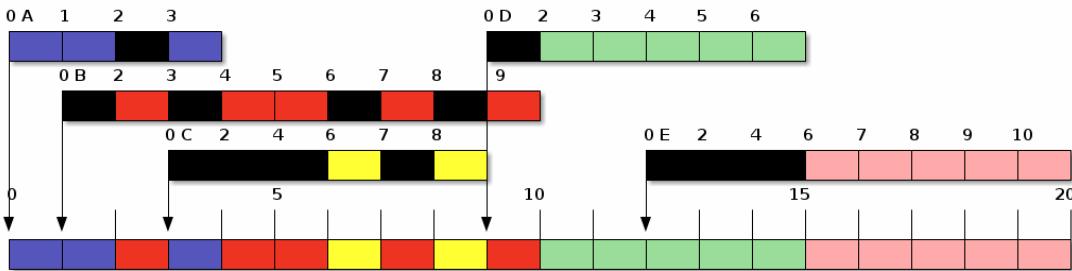


Figura 5.11: Ronda egoísta (SRR) con $a = 2$ y $b = 1$

Mientras $\frac{b}{a} < 1$, la prioridad de un proceso entrante eventualmente alcanzará a la de los procesos aceptados, y comenzará a ejecutarse. Mientras el control va alternando entre dos o más procesos, la prioridad de todos ellos será la misma (esto es, son despachados efectivamente por una simple ronda).

Incluso cuando $\frac{b}{a} \geq 1$, el proceso en ejecución terminará, y B será aceptado. En este caso, este esquema se convierte en FCFS.

Si $\frac{b}{a} = 0$ (esto es, si $b = 0$), los procesos recién llegados serán aceptados inmediatamente, con lo cual se convierte en una ronda. Mientras $0 < \frac{b}{a} < 1$, la ronda será *relativamente egoísta*, dándole entrada a los nuevos procesos incluso si los que llevan mucho tiempo ejecutando son muy largos (y por tanto, su prioridad es muy alta).

5.2.6 Retroalimentación multinivel (FB)

(*Multilevel Feedback*)

El mecanismo descrito en la sección anterior, la *ronda egoísta*, introdujo el concepto de tener no una sino que varias colas de procesos, que recibirán diferente tratamiento. Este mecanismo es muy poderoso, y se emplea en prácticamente todos los planificadores en uso hoy en día. Antes de abordar al esquema de retroalimentación multinivel, conviene presentar cómo opera un sistema con múltiples colas de prioridad.

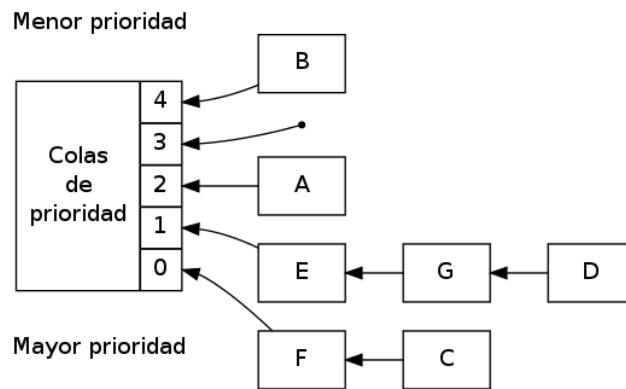


Figura 5.12: Representación de un sistema con cinco colas de prioridad y siete procesos listos

La figura 5.12 ilustra cómo se presentaría una situación bajo esta lógica: El sistema hipotético tiene cinco colas de prioridad, y siete procesos listos para ser puestos en ejecución. Puede haber colas vacías, como en este caso la 3. Dado que la cola de mayor prioridad es la 0, el planificador elegirá únicamente entre los procesos que están *formados* en ella: F o C. Sólo cuando estos procesos terminen (o sean enviados a alguna otra cola), el planificador continuará con aquellos que estén en las siguientes colas.

La *retroalimentación multinivel* basa su operación en más de una cola — Pero en este caso, todas ellas tendrán el mismo tratamiento *general*, distinguiéndose sólo por su nivel de *prioridad*, C_0 a C_n . El despachador elegirá para su ejecución al proceso que esté al frente de la cola de mayor prioridad que tenga algún proceso esperando C_i , y tras un número predeterminado de ejecuciones, lo *degrada* a la cola de prioridad inmediata inferior C_{i+1} .

El mecanismo de retroalimentación multinivel favorece a los procesos cortos, dado que terminarán sus tareas sin haber sido marcados como de prioridades inferiores.

La ejecución del juego de datos con que han sido presentados los algoritmos anteriores bajo este esquema da los siguientes resultados:

Proceso	Tiempo de Llegada	t	Inicio	Fin	T	E	P
A	0	3	0	7	7	4	1.7
B	1	5	1	18	17	12	3.4
C	3	2	3	6	3	1	1.5
D	9	5	9	19	10	5	2.0
E	12	5	12	20	8	3	1.6
Promedio		4			9	5	2.04

Dado que ahora hay que representar la cola en la que está cada uno de los procesos, en la figura 5.13 se presenta sobre cada una de las líneas de proceso la prioridad de la cola en que se encuentra antes del *quantum* a iniciar:

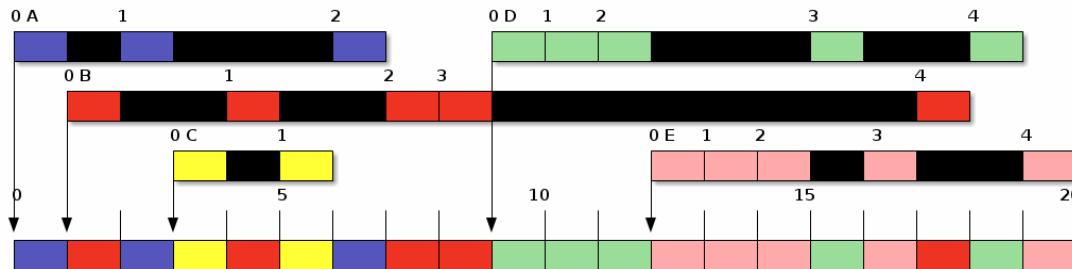


Figura 5.13: Retroalimentación multinivel (FB) básica

Llama la atención que prácticamente todos los números apuntan a que esta es una peor estrategia que las presentadas anteriormente — Los únicos procesos beneficiados en esta ocasión son los recién llegados, que podrán avanzar al principio, mientras los procesos más largos serán castigados y podrán eventualmente (a mayor ρ) enfrentar inanición.

Sin embargo, esta estrategia permite ajustar dos variables: Una es la cantidad de veces que un proceso debe ser ejecutado antes de ser *degradado* a la prioridad inferior, y la otra es la duración del *quantum* asignado a las colas subsecuentes.

Otro fenómeno digno a comentar es el que se presenta a los ticks 8, 10, 11, 13 y 14: El despachador interrumpe la ejecución del proceso activo, para volver a cedérsela. Esto ocurre porque, efectivamente, concluyó su *quantum* — Idealmente, el despachador se dará cuenta de esta situación de inmediato y no iniciará un cambio de contexto *al mismo proceso*. En caso contrario, el trabajo perdido por gasto administrativo se vuelve innecesariamente alto.

El panorama cambia al ajustar estas variables: Si se elige un *quantum* de $2^n q$, donde n es el identificador de cola y q la longitud del *quantum* base, un proceso largo será detenido por un cambio de contexto al llegar a q , $3q$, $7q$, $15q$, etc. lo que llevará al número total de cambios de contexto a $\log(\frac{t(p)}{q})$, lo cual resulta atractivo frente a los $\frac{t(p)}{q}$ cambios de contexto que tendría bajo un esquema de ronda.

Tras de estos ajustes ante el juego de procesos con una retroalimentación multinivel con un incremento exponencial al *quantum* se obtiene como resultado:

Proceso	Tiempo de Llegada	t	Inicio	Fin	T	E	P	
A		0	3	0	4	4	1	1.3
B		1	5	1	10	9	4	1.8
C		3	2	4	8	5	3	2.5
D		9	5	10	18	9	4	1.8
E		12	5	13	20	8	3	1.6
Promedio			4		7	3	1.8	

Los promedios de tiempos de terminación, respuesta, espera y penalización para este conjunto de procesos resultan mejores incluso que los de la ronda.

En este caso, a pesar de que esta estrategia favorece a los procesos recién llegados, al tick 3, 9 y 10, llegan nuevos procesos, pero a pesar de estar en la cola de mayor prioridad, no son puestos en ejecución, dado que llegaron a la mitad del *quantum* (largo) de otro proceso.

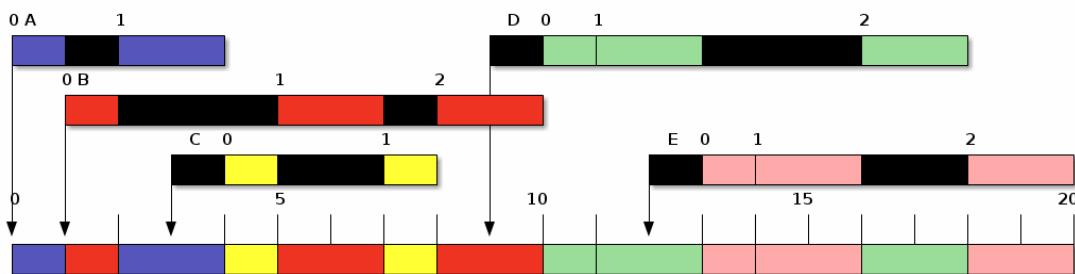


Figura 5.14: Retroalimentación multinivel (FB) con q exponencial

Típicamente se emplean incrementos mucho más suaves, y de crecimiento más controlado, como nq o incluso $q \log(n)$, dado que en caso contrario un proceso muy largo podría causar muy largas inaniciones para el resto del sistema.

Para evitar la inanición, puede considerarse también la retroalimentación en sentido inverso: Si un proceso largo fue *degradado* a la cola C_P y pasa determinado tiempo sin recibir servicio, puede *promoverse* de nuevo a la cola $\$C_{P-1}$ para que no sufra inanición.

Hoy en día, muchos de los principales sistemas operativos operan bajo diferentes versiones de retroalimentación multinivel, y típicamente con hasta decenas de colas.

5.2.7 Lotería

Los mecanismos hasta aquí descritos vienen con largas décadas de desarrollo. Uno de los últimos algoritmos que ha sido ampliamente difundido en unirse a esta lista es el de *planificación por lotería*, publicado por [Carl Waldspurger y William Weihl \(1994\)](#).

Bajo el esquema de la *lotería*, cada proceso tiene un número determinado de boletos, y cada boleto le representa una oportunidad de jugar a la lotería. Cada vez que el planificador tiene que elegir el siguiente proceso a poner en ejecución, elige un número al azar¹⁰, y otorga el siguiente quantum al proceso que tenga el boleto ganador. El boleto ganador *no es retirado*, esto es, la probabilidad de que determinado proceso sea puesto en ejecución no varía entre invocaciones sucesivas del planificador.

Las prioridades pueden representarse en este esquema de forma muy sencilla: Un proceso al que se le quiere dar mayor prioridad simplemente tendrá más boletos; si el proceso *A* tiene 20 boletos y el proceso *B* tiene 60, será tres veces más probable que el siguiente turno toque a *B* que a *A*.

El esquema de planificación por lotería contempla que los procesos puedan cooperar entre sí: Si *B* estuviera esperando un resultado de *A*, podría transferirle sus boletos para aumentar la probabilidad de que sea puesto en ejecución.

A pesar de su simplicidad, el esquema de planificación por lotería resulta justo tanto a procesos cortos como a largos, y presenta una degradación muy suave incluso en entornos de saturación. Claro, al derivar de un proceso aleatorio, resulta imposible presentar una comparación de este mecanismo abordados previamente.

5.2.8 Esquemas híbridos

En líneas generales, los siete algoritmos presentados pueden clasificarse sobre dos discriminadores primarios: Si están pensados para emplearse en multitarea cooperativa o

¹⁰Si bien operar un generador de números aleatorios en estricto sentido sería demasiado caro para un proceso que se ejecuta decenas o cientos de veces por segundo, para *jugar a la lotería* es suficiente emplear un generador débil pseudoaleatorio. El artículo en que este mecanismo fue presentado presenta la implementación del algoritmo Park-Miller, $S' = (A \times S) \bmod (2^{31} - 1)$ con $A = 16807$, implementado en 12 instrucciones de procesador RISC.

preventiva, y si emplean información *intrínseca* a los procesos evaluados o no lo hacen, esto es, si un proceso es tratado de distinta forma dependiendo de su historial de ejecución.

Cuadro 5.1: Caracterización de los mecanismos de planificación a corto plazo

	No considera intrínseca	Considera intrínseca
Cooperativa	Primero llegado primero servido (FCFS)	Proceso más corto (SPN), Proceso más penalizado (HPRN)
Preventiva	Ronda (RR) Lotería	Proceso más corto preventivo (PSPN), Retroalimentación (FB), Ronda egoísta (SRR)

Ahora bien, estas características primarias pueden ser empleadas en conjunto, empleando diferentes algoritmos a diferentes niveles, o cambiándolos según el patrón de uso del sistema, aprovechando de mejor manera sus bondades y logrando evitar sus deficiencias. A continuación, algunos ejemplos de esquemas híbridos.

Algoritmo por cola dentro de FB

Al introducir varias colas, se abre la posibilidad de que cada una de ellas siga un esquema diferente para elegir cuál de sus procesos está a la cabeza. En los ejemplos antes presentados, todas las colas operaban siguiendo una ronda, pero podría contemplarse, por ejemplo, que parte de las colas sean procesadas siguiendo una variación de PSPN que empuje a los procesos más largos a colas que les puedan dar atención con menor número de interrupciones (incluso sin haberlos ejecutado aún).

Podría emplearse un esquema SRR para las colas de menor prioridad, siendo que ya tienen procesos que han esperado mucho tiempo para su ejecución, para –sin que repercutan en el tiempo de respuesta de los procesos cortos que van entrando a las colas superiores– terminen lo antes posible su ejecución.

Métodos dependientes del estado del sistema

Los parámetros de operación pueden variar también dependiendo del estado actual del sistema, e incluso tomando en consideración valores externos al despachador. Algunas ideas al respecto son:

- Si los procesos listos son *en promedio* no muy largos, y el valor de saturación es bajo ($\rho < 1$), optar por los métodos que menos sobrecarga administrativa signifiquen, como FCFS o SPN (o, para evitar los peligros de la multitarea cooperativa, un RR con un *quantum* muy largo). Si el despachador observa que la longitud de la cola excede un valor determinado (o *muestra una tendencia* en ese sentido, al incrementarse ρ), cambiar a un mecanismo que garantice una mejor distribución de la atención, como un RR con *quantum* corto o PSPN.
- Usar un esquema simple de ronda. La duración de un *quantum* puede ser ajustada periódicamente (a cada cambio de contexto, o como un cálculo periódico), para que la duración del siguiente *quantum* dependa de la cantidad de procesos en espera en la lista, $Q = \frac{q}{n}$.

Si hay pocos procesos esperando, cada uno de ellos recibirá un *quantum* más largo, reduciendo la cantidad de cambios de contexto. Si hay muchos, cada uno de ellos tendrá que esperar menos tiempo para comenzar a liberar sus pendientes.

Claro está, la duración de un *quantum* no debe reducirse más allá de cierto valor mínimo, definido según la realidad del sistema en cuestión, dado que podría aumentar demasiado la carga burocrática.

- Despachar los procesos siguiendo una ronda, pero asignarles una duración de *quantum* proporcional a su prioridad externa (fijada por el usuario). Un proceso de mayor prioridad ejecutará *quantums* más largos.
- *Peor servicio a continuación* (*WSN*, *Worst Service Next*). Es una generalización sobre varios de los mecanismos mencionados; su principal diferencia respecto a *HPRN* es que no sólo se considera *penalización* el tiempo que ha pasado esperando en la cola, sino que se considera el número de veces que ha sido interrumpido por el temporizador o su prioridad externa, y se considera (puede ser a favor o en contra) el tiempo que ha tenido que esperar por E/S u otros recursos. El proceso que ha sufrido del *peor servicio* es seleccionado para su ejecución, y si varios empatan, se elige uno en ronda.

La principal desventaja de *WSN* es que, al considerar tantos factores, el tiempo requerido por un lado para recopilar todos estos datos, y por otro lado calcular el peso que darán a cada uno de los procesos implicados, puede impactar en el tiempo global de ejecución. Es posible acudir a *WSN* periódicamente (y no cada vez que el despachador es invocado) para que reordene las colas según criterios generales, y abanzar sobre dichas colas con algoritmos más simples, aunque esto reduce la velocidad de reacción ante cambios de comportamiento.

- Algunas versiones históricas de Unix manejaban un esquema en que la prioridad especificada por el usuario¹¹ era matizada y re-evaluada en el transcurso de su ejecución.

Periódicamente, para cada proceso se calcula una prioridad *interna*, que depende de la prioridad *externa* (especificada por el usuario) y el tiempo consumido recientemente por el proceso. Conforme el proceso recibe mayor tiempo de procesador, esta última cantidad decrece, y aumenta conforme el proceso espera (sea por decisión del despachador o por estar en alguna espera).

Esta prioridad interna depende también del tamaño de la lista de procesos listos para su ejecución: Entre más procesos haya pendientes, más fuerte será la modificación que efectúe.

El despachador ejecutará al proceso que tenga una mayor prioridad después de realizar estos pesos, decidiendo por ronda en caso de haber empates. Claro está, este algoritmo resulta sensiblemente más caro computacionalmente, aunque más justo, que aquellos sobre los cuales construye.

5.2.9 Resumiendo

En esta sección se presentan algunos mecanismos básicos de planificación a corto plazo. Como se indica en la parte final, es muy poco común encontrar a ninguno de estos mecanismos en un *estado puro* — Normalmente se encuentra combinación de ellos, con diferentes parámetros según el nivel en el cual se está ejecutando.

Rendimiento ante diferentes cargas de procesos

Los algoritmos presentados en el transcurso de esta sección fueron presentados y comparados ante una determinada carga de procesos. No se puede asumir, sin embargo, que su comportamiento será igual ante diferentes distribuciones: Un patrón de trabajo donde predominen los procesos cortos y haya unos pocos procesos largos probablemente se verá mucho más penalizado por un esquema SRR (y mucho más favorecido por un

¹¹ La *lindura*, o *niceness* de un proceso, llamada así por establecerse a través del comando `nice` al iniciar su ejecución, o `renice` una vez en ejecución

SPN o PSPN) que uno en el cual predominen los procesos largos.

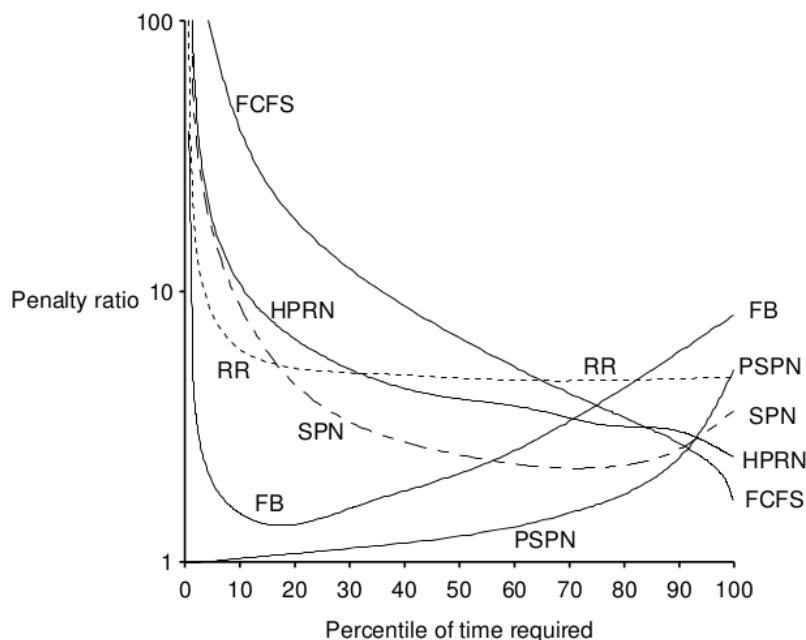


Figura 5.15: Proporción de penalización registrada por cada proceso contra el porcentaje del tiempo que éste requiere (Finkel, p.33)

Raphael Finkel realizó estudios bajo diversas cargas de trabajo, buscando comparar de forma *significativa* estos distintos mecanismos. Finkel simuló el comportamiento que estos algoritmos tendrían ante 50,000 procesos generados de forma aleatoria, siguiendo una distribución exponencial tanto en sus tiempos de llegada como duraciones en ejecución, y manteniendo como parámetro de equilibrio un nivel de saturación $\rho = 0,8$ ($\alpha = 0,8, \beta = 1,0$), obteniendo como resultado las figuras aquí reproducidas (5.15 y 5.16) comparando algunos aspectos importantes de los diferentes despachadores.

Duración mínima del *quantum*

La penalización por cambios de contexto en esquemas preventivos como la *ronda* puede evitarse empleando *quanta* mayores. Pero abordando la contraparte, ¿qué tan corto tiene sentido que sea un *quantum*? Con el hardware y las estructuras requeridas por los sistemas operativos de uso general disponibles hoy en día, un cambio de contexto requiere del orden de 10 microsegundos (Silberschatz, p.187), por lo que incluso con el *quantum* de 10ms (el más corto que manejan tanto Linux como Windows), representa apenas la milésima parte del tiempo efectivo de proceso.

Una estrategia empleada por Control Data Corporation para la CDC6600 (comercializada a partir de 1964, y diseñada por Seymour Cray) fue emplear hardware especializado que permitiera efectivamente *compartir el procesador*: Un sólo procesador tenía 10 juegos de registros, permitiéndole alternar entre 10 procesos con un *quantum* efectivo igual a la velocidad del reloj. A cada paso del reloj, el procesador cambiaba el juego de registros. De este modo, un sólo procesador de muy alta velocidad para su momento (1 MHz) aparecía ante las aplicaciones como 10 procesadores efectivos, cada uno de 100 KHz, reduciendo los costos al implementar sólamente una vez cada una de las unidades funcionales. Puede verse una evolución de esta idea retomada hacia mediados de la década del 2000 en los procesadores que manejan hilos de ejecución.¹².

¹²Aunque la arquitectura de la CDC6600 era *plenamente superescalar*, a diferencia de los procesadores *Hyperthreading*, que será abordada brevemente en la sección 5.4.4, en que para que dos instrucciones se ejecuten simultáneamente deben ser de naturalezas distintas, no requiriendo ambas de la misma *unidad*

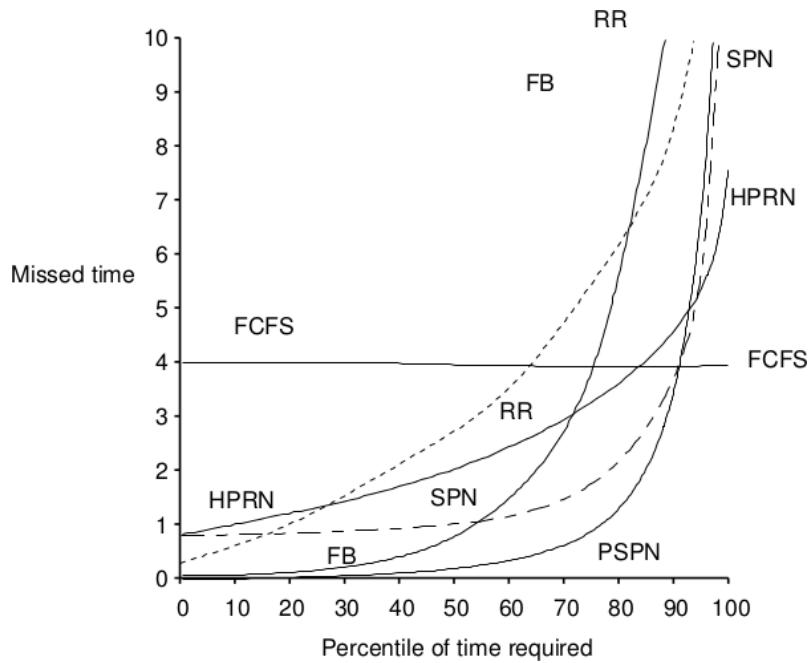


Figura 5.16: Tiempo perdido contra porcentaje de tiempo requerido por proceso (Finkel, p.34)

Esta arquitectura permitía tener multitarea real sin tener que realizar cambios de contexto, sin embargo, al tener un *nivel de concurrencia* fijo establecido en hardware no es tan fácil adecuar a un entorno cambiante, con picos de ocupación.

5.3 Planificación de hilos

Ahora bien, tras centrar toda la presente discusión en los procesos, ¿cómo caben los *hilos* en este panorama? Depende de cómo éstos son *mapeados* a procesos a ojos del planificador.

Como fue expuesto en la sección 4.2.1, hay dos clases principales de hilo: Los *hilos de usuario* o *hilos verdes*, que son completamente gestionados dentro del proceso y sin ayuda del sistema operativo, y los *hilos de núcleo* o *hilos de kernel*, que sí son gestionados por el sistema operativo como si fueran procesos. Partiendo de esto, existen tres modelos principales de mapeo:

Muchos a uno Muchos hilos son agrupados en un sólo proceso. Los *hilos verdes* entran en este supuesto: Para el sistema operativo, hay un sólo proceso; mientras tiene la ejecución, éste se encarga de repartir el tiempo entre sus hilos.

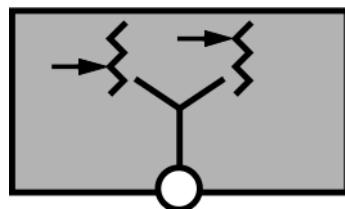


Figura 5.17: Mapeo de hilos *muchos a uno* (Imagen: Beth Plale; ver otros recursos)

Bajo este modelo, si bien el código escrito es más portable entre diferentes sistemas operativos, los hilos no aprovechan *realmente* al paralelismo, y todos los hilos pueden

funcional del CPU. El procesador de la CDC6600 no manejaba *pipelines*, sino que cada ejecución empleaba al CPU completo

tener que bloquearse cuando uno sólo de ellos realiza una llamada *bloqueante* al sistema.

Uno a uno Cada hilo es ejecutado como un *proceso ligero* (*lightweight process* o *LWP*); podría dar la impresión de que este esquema desperdicia la principal característica de los hilos, que es una mayor sencillez y rapidez de inicialización que los procesos, sin embargo, la información de estado requerida para crear un LWP es mucho menor que la de un proceso regular, y mantiene como ventaja que los hilos continúan compartiendo su memoria, descriptores de archivos y demás estructuras.

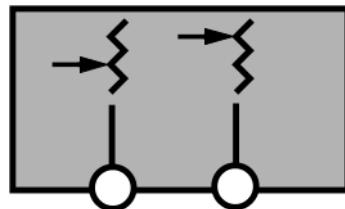


Figura 5.18: Mapeo de hilos *uno a uno* (Imagen: Beth Plale; ver otros recursos)

Este mecanismo permite a los hilos aprovechar las ventajas del paralelismo, pudiendo ejecutarse cada hilo en un procesador distinto, y como única condición para su existencia, el sistema operativo debe poder implementar los LWP.

Muchos a muchos Este mecanismo permite que existan hilos de ambos modelos: Permite la existencia de *hilos unidos* (*bound threads*), en que cada hilo corresponde a un (y sólo un) LWP, y de *hilos no unidos* (*unbound threads*), de los cuales *uno o más* estarán mapeados a cada LWP.

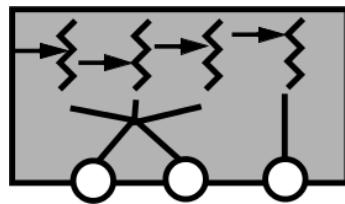


Figura 5.19: Mapeo de hilos *muchos a muchos* (Imagen: Beth Plale; ver otros recursos)

El esquema *muchos a muchos* proporciona las principales características de ambos esquemas; en caso de ejecutarse en un sistema que no soporte más que el modelo *uno a muchos*, el sistema puede caer en éste como *modo degradado*.

No se detalla en el presente texto respecto a los primeros — Cada marco de desarrollo o máquina virtual que emplee *hilos de usuario* actuará cual sistema operativo ante ellos, probablemente con alguno de los mecanismos ilustrados anteriormente.

5.3.1 Los hilos POSIX (pthreads)

La clasificación recién presentada de modelos de mapeo entre hilos y procesos se refleja aproximadamente en la categorización de los hilos POSIX (pthreads) denominada el *ámbito de contención*.

Hay dos enfoques respecto a la *contención* que deben tener los hilos, esto es: En el momento que un proceso separa su ejecución en dos hilos, ¿debe cada uno de estos recibir la misma atención que recibiría un proceso completo?

Ámbito de contención de proceso (*Process Contention Scope, PCS*; en POSIX, `PTHREAD_SCOPE_PROCESS`)

Una respuesta es que todos los hilos deben ser atendidos sin exceder el tiempo que sería asignado a un sólo proceso. Un proceso que consta de varios hilos siguiendo el modelo *muchos a uno*, o uno que multiplexa varios *hilos no unidos* bajo un modelo *muchos a muchos*, se ejecuta bajo este ámbito.

Ámbito de contención de sistema (System Contention Scope, SCS; en POSIX, PTHREAD_SCOPE_SYSTEM)

Este ámbito es cuando, en contraposición, cada hilo es visto por el planificador como un proceso independiente; este es el ámbito en el que se ejecutarían los hilos bajo el modelo *uno a uno*, o cada uno de los *hilos unidos* bajo un modelo *muchos a muchos*, dado que los hilos son tratados, para propósitos de planificación, como procesos normales.

La definición de pthreads apunta a que, si bien el programador puede solicitar que sus hilos sean tratados bajo cualquiera de estos procesos, una implementación específica puede implementar ambos o sólo uno de los ámbitos. Un proceso que solicite que sus hilos sean programados bajo un ámbito no implementado serán ejecutados bajo el otro, notificando del error (pero permitiendo continuar con la operación).

Las implementaciones de pthreads tanto en Windows como en Linux sólo contemplan SCS.

Respecto a los otros aspectos mencionados en este capítulo, la especificación pthreads incluye funciones por medio de las cuales el programador puede solicitar al núcleo la prioridad de cada uno de los hilos por separado (`pthread_setschedprio`) e incluso solicitar el empleo de determinado algoritmo de planificación (`sched_setscheduler`).

5.4 Planificación de multiprocesadores

Hasta este punto, el enfoque de este capítulo se ha concentrado en la planificación asumiendo un sólo procesador. Del mismo modo que lo que se ha visto hasta este momento, no hay una sola estrategia que pueda ser vista como superior a las demás en todos los casos.

Para trabajar en multiprocesadores, puede mantenerse una sola lista de procesos e ir despachándolos a cada uno de los procesadores como unidades de ejecución equivalentes e idénticas, o pueden mantenerse listas separadas de procesos. A continuación se presentan algunos argumentos respecto a estos enfoques.

5.4.1 Afinidad a procesador

En un entorno multiprocesador, después de que un proceso se ejecutó por cierto tiempo, tendrá parte importante de sus datos copiados en el caché del procesador en el que fue ejecutado. Si el despachador decidiera lanzarlo en un procesador que no compartiera dicho caché, estos datos tendrían que ser *invalidados* para mantener la coherencia, y muy probablemente (por *localidad de referencia*) serían vuelto a cargar al caché del nuevo procesador.

Los procesadores actuales normalmente tienen disponibles varios niveles de caché; si un proceso es migrado entre dos núcleos del mismo procesador, probablemente sólo haga falta invalidar los datos en el caché más interno (L1), dado que el caché en chip (L2) es compartido entre los varios núcleos del mismo chip; si un proceso es migrado a un CPU físicamente separado, será necesario invalidar también el caché en chip (L2), y mantener únicamente el del controlador de memoria (L3).

Pero dado que la situación antes descrita varía de computadora a computadora, no se puede enunciar una regla general — Más allá de que el sistema operativo debe conocer cómo están estructurados los diversos procesadores que tiene a su disposición, y buscar realizar las migraciones *más baratas*, aquellas que tengan lugar entre los procesadores más cercanos.

Resulta obvio por esto que un proceso que fue ejecutado en determinado procesador vuelve a ser ejecutado en el mismo, esto es, el proceso *tiene afinidad* por cierto procesador. Un proceso que *preferentemente* será ejecutado en determinado procesador se dice que *tiene afinidad suave* por ese procesador, pero determinados patrones de carga (por ejemplo,

una mucho mayor cantidad de procesos afines a cierto procesador que a otro, saturando su cola de procesos listos, mientras que el segundo procesador tiene tiempo disponible) pueden llevar a que el despachador decida activarlo en otro procesador.

Por otro lado, algunos sistemas operativos ofrecen la posibilidad de declarar *afinidad dura*, con lo cual se *garantiza* a un proceso que siempre será ejecutado en un procesador, o en un conjunto de procesadores.

Un entorno NUMA, por ejemplo, funcionará mucho mejor si el sistema que lo emplea maneja tanto un esquema de afinidad dura como algoritmos de asignación de memoria que le aseguren que un proceso siempre se ejecutará en el procesador que tenga mejor acceso a sus datos.

5.4.2 Balanceo de cargas

En un sistema multiprocesador, la situación ideal es que todos los procesadores estén despachando trabajos al 100% de su capacidad. Sin embargo, ante una definición tan rígida, la realidad es que siempre habrá uno o más procesadores con menos del 100% de carga, o uno o más procesadores con procesos encolados y a la espera, o incluso ambas situaciones.

La divergencia entre la carga de cada uno de los procesadores debe ser lo más pequeña posible. Para lograr esto, se pueden emplear esquemas de *balanceo de cargas*: Algoritmos que analicen el estado de las colas de procesos y, de ser el caso, transfieran procesos entre las colas para homogeneizarlas. Claro está, el balanceo de cargas puede actuar precisamente en sentido contrario de la afinidad al procesador, y efectivamente puede reubicar a los procesos con afinidad suave.

Hay dos estrategias primarias de balanceo: Por un lado, la migración activa o migración *por empuje* (*push migration*) consiste en una tarea que ejecuta como parte del núcleo y periódicamente revisa el estado de los procesadores, y en caso de encontrar un desbalance mayor a cierto umbral, *empuja* a uno o más procesos de la cola del procesador más ocupado a la del procesador más libre. Linux ejecuta este algoritmo cada 200 milisegundos.

Por otro lado, está la migración pasiva o migración *por jalón* (*pull migration*). Cuando algún procesador queda sin tareas pendientes, ejecuta al proceso especial *desocupado* (*idle*). Ahora, el proceso *desocupado* no significa que el procesador detenga su actividad — Ese tiempo puede utilizarse para ejecutar tareas del núcleo. Una de esas tareas puede consistir en averiguar si hay procesos en espera en algún otro de los procesadores, y de ser así, *jalarlo* a la cola de este procesador.

Ambos mecanismos pueden emplearse —y normalmente lo hacen— en el mismo sistema. Los principales sistemas operativos modernos emplean casi siempre ambos mecanismos.

Como sea, debe mantenerse en mente que todo balanceo de cargas que se haga entre los procesadores conllevará una penalización en términos de afinidad al CPU.

5.4.3 Colas de procesos: ¿Una o varias?

En los puntos relativos al multiprocesamiento hasta ahora abordados se parte del supuesto que hay una cola de procesos listos por cada procesador. Si, en cambio, hubiera una cola global de procesos listos de la cual el siguiente proceso a ejecutarse fuera asignándose al siguiente procesador, fuera éste cualquiera de los disponibles, podría ahorrarse incluso elegir entre una estrategia de migración *por empuje* o *por jalón* — Mientras hubiera procesos pendientes, éstos serían asignados al siguiente procesador que tuviera tiempo disponible.

El enfoque de una sola cola, sin embargo, no se usa en ningún sistema en uso amplio. Esto es en buena medida porque un mecanismo así haría mucho más difícil mantener la afinidad al procesador y restaría flexibilidad al sistema completo.

5.4.4 Procesadores con soporte a hilos hardware (*hyperthreading*)

El término de *hilos* como abstracción general de algo que corre con mayor frecuencia y dentro de un mismo proceso puede llevar a una confusión, dado que en esta sección se tocan dos temas relacionados. Para esta subsección en particular, se hace referencia a los *hilos en hardware* que forman parte de ciertos procesadores, ofreciendo al sistema una *casi* concurrencia adicional.

Conforme han subido las frecuencias de reloj en el cómputo más allá de lo que permite llevar al sistema entero como una sola unidad, una respuesta recurrente ha sido incrementar el paralelismo. Y esto no sólo se hace proveyendo componentes completos adicionales, sino que separando las *unidades funcionales* de un procesador.

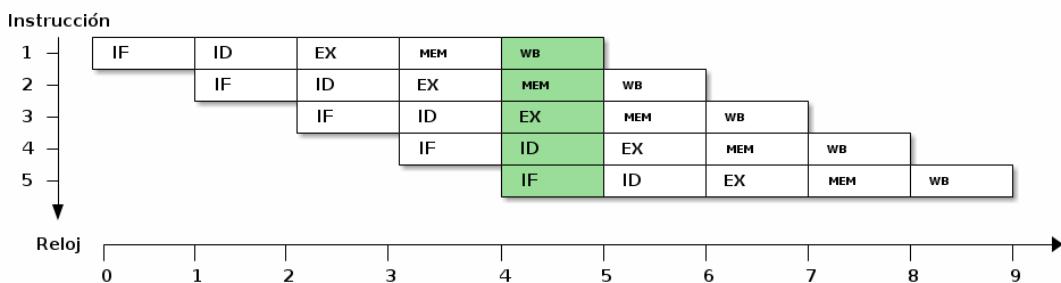


Figura 5.20: Descomposición de una instrucción en sus cinco pasos clásicos para organizarse en un *pipeline*

El flujo de una sola instrucción a través de un procesador simple como el MIPS puede ser dividido en cinco secciones principales, creando una estructura conocida como *pipeline* (tubería). Idealmente, en todo momento habrá una instrucción diferente ejecutando en cada una de las secciones del procesador, como lo ilustra la figura 5.20. Las secciones en los procesadores MIPS clásicos son:

- Recuperación de la instrucción (*Instruction Fetch, IF*)
- Decodificación de la instrucción (*Instruction Decode, ID*)
- Ejecución (*Execution, EX*)
- Acceso a datos (MEM)
- Almacenamiento (*Writeback, WB*)

La complejidad de los procesadores actuales ha crecido ya por encima de lo aquí delineado (el Pentium 4 tiene más de 20 etapas), sin embargo se presenta esta separación como base para la explicación. Un procesador puede iniciar el procesamiento de una instrucción cuando la siguiente apenas avanzó la quinta parte de su recorrido — De este modo, puede lograr un paralelismo interno, manteniendo idealmente siempre ocupadas a sus partes funcionales.

Sin embargo, se ha observado que hay patrones recurrentes que intercalan operaciones que requieren servicio de diferentes componentes del procesador, o que requieren de la inserción de *burbujas* porque una unidad es más lenta que las otras — Lo cual lleva a que incluso empleando *pipelines*, un procesador puede pasar hasta el 50% del tiempo esperando a que haya datos disponibles solicitados a la memoria.

Para remediar esto, varias de las principales familias de procesadores presentan a un mismo *núcleo* de procesador como si estuviera compuesto de dos o más *hilos hardware* (conocidos en el mercado como *hyper-threads*). Esto puede llevar a una mayor utilización del procesador, siguiendo patrones como el ilustrado en la figura 5.21

Hay que recordar que, a pesar de que se *presenten* como hilos independientes, el rendimiento de cada uno depende de la secuencia particular de instrucciones del otro — No puede esperarse que el incremento en el procesamiento sea de 2x; la figura presenta

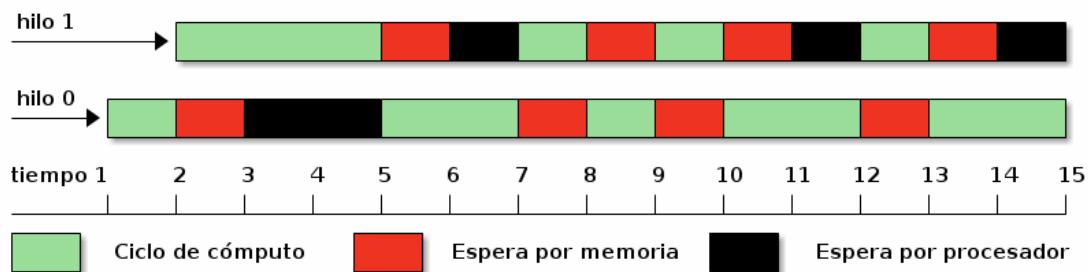


Figura 5.21: Alternando ciclos de cómputo y espera por memoria, un procesador que implementa hilos hardware (*hyperthreaded*) se presenta como dos procesadores

varios puntos en que un hilo está en *espera por procesador*, dado que el otro está empleando las unidades funcionales que éste requiere.

La planificación de los hilos hardware sale del ámbito del presente material, y este tema se presenta únicamente para aclarar un concepto que probablemente confunda al alumno por su similitud; los hilos en hardware implican cuestiones de complejidad tal como el ordenamiento específico de las instrucciones, predicción de ramas de ejecución, e incluso asuntos relativos a la seguridad, dado que se han presentado *goteos* que permiten a un proceso ejecutando en un hilo *espiar* el estado del procesador correspondiente a otro de los hilos. Para abundar al respecto, el ámbito adecuado podría ser un texto orientado a la construcción de compiladores (ordenamiento de instrucciones, aprovechamiento del paralelismo), o uno de arquitectura de sistemas (estudio del pipeline, aspectos del hardware).

Esta estrategia guarda gran similitud, y no puede evitar hacerse el paralelo, con la *compartición de procesador* empleada por la CDC6600, presentada en la sección 5.2.9.

5.5 Tiempo real

Todos los esquemas de manejo de tiempo hasta este momento se han enfocado a repartir el tiempo disponible entre todos los procesos que requieren atención. Es necesario también abordar a los procesos que *requieren garantías de tiempo*: procesos que para poder ejecutarse deben garantizar el haber tenido determinado tiempo de proceso antes de un tiempo límite. Los procesos con estas características se conocen como *de tiempo real*.

Hay ejemplos de procesos que requieren este tipo de planificación a todo nivel; los ejemplos más comunes son los controladores de dispositivos y los recodificadores o reproductores de medios (audio, video). La lógica general es la misma:

Para agendarse como un proceso con requisitos de tiempo real, éste debe declarar sus requisitos de tiempo (formalmente, *efectuar su reserva de recursos*) al iniciar su ejecución o en el transcurso de la misma. Claro está, siendo que los procesos de tiempo real obtienen una *prioridad* mucho mayor a otros, normalmente se requerirá al iniciar el proceso que éste declare que durante parte de su ejecución trabajará con restricciones de tiempo real.

5.5.1 Tiempo real duro y suave

Supóngase que un dispositivo genera periódicamente determinada cantidad de información y la va colocando en un área determinada de memoria compartida (en un *buffer*). Al inicializarse, su controlador declarará al sistema operativo cuánto tiempo de ejecución le tomará recoger y procesar dicha información, liberando el *buffer* para el siguiente ciclo de escritura del dispositivo, y la frecuencia con que dicha operación tiene que ocurrir.

En un sistema capaz de operar con garantías de tiempo real, si el sistema operativo puede garantizar que en ese intervalo le otorgará al proceso en cuestión suficiente tiempo para procesar esta información, el proceso se ejecuta; en caso contrario, recibe un error antes de que esto ocurra por medio del cual podrá alertar al usuario.

Los sistemas en que el tiempo máximo es garantizable son conocidos como de *tiempo real duro*.

La necesidad de atención en tiempo real puede manejarse *periódica* (por ejemplo, *requiero del procesador por 30ms cada segundo*), o *aperiódica*, por ocurrencia única (*necesito que este proceso, que tomará 600ms, termine de ejecutarse en menos de 2s*).

Realizar una reserva de recursos requiere que el planificador sepa con certeza cuánto tiempo toma realizar las tareas de sistema que ocurrirán en el periodo en cuestión. Cuando entran en juego algunos componentes de los sistemas de propósito general que tienen una latencia con variaciones impredecibles (como el almacenamiento en disco o la memoria virtual) se vuelve imposible mantener las garantías de tiempo ofrecidas. Por esta razón, en un sistema operativo de propósito general empleando hardware estándar *no es posible implementar tiempo real duro*.

Para solventar necesidades como las expresadas en sistemas de uso general, el *tiempo real suave* sigue requiriendo que los procesos críticos reciban un trato prioritario por encima de los procesos comunes; agendar a un proceso con esta prioridad puede llevar a la inanición de procesos de menor prioridad y un comportamiento que bajo ciertas métricas resultaría *injusto*. Un esquema de tiempo real suave puede implementarse a través de un esquema similar al de la *retroalimentación multinivel*, con las siguientes particularidades:

- La cola de tiempo real recibe prioridad sobre todas las demás colas
- La prioridad de un proceso de tiempo real *no se degrada* conforme se ejecuta repetidamente
- La prioridad de los demás procesos *nunca llegan a subir* al nivel de tiempo real por un proceso automático (aunque sí puede hacerse por una llamada explícita)
- La latencia de despacho debe ser mínima

Casi todos los sistemas operativos en uso amplio hoy en día ofrecen facilidades básicas de tiempo real suave.

5.5.2 Sistema operativo interrumpible (*prevenible*)

Para que la implementación de tiempo real suave sea apta para estos requisitos es necesario modificar el comportamiento del sistema operativo. Cuando un proceso de usuario hace una llamada al sistema, o cuando una interrupción corta el flujo de ejecución, hace falta que el sistema procese completa la rutina que da servicio a dicha solicitud antes de que continúe operando. Se dice entonces que el sistema operativo *no es prevenible* o *no es interrumpible*.

Para lograr que el núcleo pueda ser interrumpido para dar el control de vuelta a procesos de usuario, un enfoque fue el poner *puntos de interrupción* en los puntos de las funciones del sistema donde fuera seguro, tras asegurarse que las estructuras estaban en un estado estable. Esto, sin embargo, no modifica en mucho la situación porque estos puntos son relativamente pocos, y es muy difícil reestructurar su lógica para permitir puntos de prevención adicionales.

Otro enfoque es hacer al núcleo entero completamente interrumpible, asegurándose de que, a lo largo de todo su código, todas las modificaciones a estructuras internas estén protegidas por mecanismos de sincronización, como los estudiados en la sección 4.3. Este método ralentiza varios procesos del núcleo, pero es mucho más flexible, y ha sido adoptado por los diversos sistemas operativos. Tiene la ventaja adicional de que permite que haya *hilos* del núcleo corriendo de forma concurrente en todos los procesadores del sistema.

5.5.3 Inversión de prioridades

Un efecto colateral de que las estructuras del núcleo estén protegidas por mecanismos de sincronización es que puede presentarse la *inversión de prioridades*. Esto es:

- Un proceso *A* de baja prioridad hace una llamada al sistema, y es interrumpido a la mitad de dicha llamada
- Un proceso *B* de prioridad *tiempo real* hace una segunda llamada al sistema, que requiere de la misma estructura que la que tiene bloqueada *A*

Al presentarse esta situación, *B* se quedará esperando hasta que *A* pueda ser nuevamente agendado — Esto es, un proceso de alta prioridad no podrá avanzar hasta que uno de baja prioridad libere su recurso.

La respuesta introducida por Solaris 2 a esta situación a este fenómeno es la *herencia de prioridades*: Todos los procesos que estén accediendo (y, por tanto, bloqueando) recursos requeridos por un proceso de mayor prioridad, serán tratados como si fueran de la prioridad de dicho recurso *hasta que terminen de utilizar el recurso en cuestión*, tras lo cual volverán a su prioridad nativa.

5.6 Otros recursos

- *Simulation of CPU Process scheduling*

<http://stimulationofcp.sourceforge.net/>

P. A. Krishnan (1999-2009); programa en Java desarrollado para un curso de Sistemas Operativos, presentando la simulación de distintos algoritmos de planificación.

- *Thread Scheduling (ctd): quanta, switching and scheduling algorithms*

http://www.javamex.com/tutorials/threads/thread_scheduling_2.shtml

Neil Coffey (2013); Javamex tutorial and performance information

- *Microprocessor Design / Pipelined Processors*

http://en.wikibooks.org/wiki/Microprocessor_Design/Pipelined_Processors

WikiBooks.org (2007-2013)

- *Thread scheduling and synchronization*

<http://www.cs.indiana.edu/classes/b534-plal/ClassNotes/sched-synch-details4.pdf>

Beth Plale (2003); Indiana University

- Páginas de manual de Linux:

- Hilos POSIX (pthreads)

<http://man7.org/linux/man-pages/man7/pthreads.7.html>

- Modificación del ámbito de contención (pthread_attr_setscope)

http://man7.org/linux/man-pages/man3/pthread_attr_setscope.3.html

- Interfaz para solicitar cambios a los parámetros y políticas de planificación (sched_setscheduler)

http://man7.org/linux/man-pages/man2/sched_setscheduler.2.html

- *Windows Internals, 6th edition*

<http://technet.microsoft.com/en-us/sysinternals/bb963901.aspx>

Mark Russinovich, David A. Solomon y Alex Ionescu (2012); por Microsoft Press. El capítulo 5 aborda a profundidad los temas de hilos y procesos bajo Windows, y está disponible como ejemplo del libro para su descarga en la página referida.

- *Optimizing preemption*

<https://lwn.net/Articles/563185/>

Jonathan Corbet (2013), *Linux Weekly News*

6 — Administración de memoria

6.1 Funciones y operaciones del administrador de memoria

El único espacio de almacenamiento que el procesador puede utilizar directamente, más allá de los registros (que si bien le son internos y sumamente rápidos, son de capacidad demasiado limitada) es la memoria física. Todas las arquitecturas de procesador tienen instrucciones para interactuar con la memoria, pero ninguna lo tiene para hacerlo con medios *persistentes* de almacenamiento, como las unidades de disco. Cabe mencionar que cuando se encuentre en un texto referencia al *almacenamiento primario* siempre se referirá a la memoria, mientras que el *almacenamiento secundario* se refiere a los discos u otros medios de almacenamiento persistente.

Todos los programas a ejecutar deben cargarse a la memoria del sistema antes de ser utilizados. En este capítulo se mostrará cómo el sistema operativo administra la memoria para permitir que varios procesos la compartan — Esta tarea debe preverse desde el proceso de compilación de los programas (en particular, la fase de *ligado*). Hoy en día, además, casi todos los sistemas operativos emplean implementaciones que requieren de hardware especializado — La *Unidad de Manejo de Memoria* (MMU). Se describirá cómo se manejaban los sistemas multitarea antes de la universalización de las MMU, y qué rol juegan hoy en día.

En esta primer sección se presentarán algunos conceptos base que se emplearán en las secciones subsecuentes.

6.1.1 Espacio de direccionamiento

La memoria está estructurada como un arreglo direccionable de bytes. Esto es, al solicitar el contenido de una dirección específica de memoria, el hardware entregará un byte (8 bits), y no menos. Si se requiere hacer una operación sobre *bits* específicos, se deberá solicitar y almacenar bytes enteros. En algunas arquitecturas, el *tamaño de palabra* es mayor — Por ejemplo, los procesadores Alpha incurrían en *fallas de alineación* si se solicitaba una dirección de memoria no alineada a 64 bits, y toda llamada a direcciones *mal alineadas* tenía que ser *atrapada* por el sistema operativo, re-alineada, y entregada.

Un procesador que soporta un *espacio de direccionamiento* de 16 bits puede referirse *directamente* a hasta 2^{16} bytes, esto es, a hasta 65,536 bytes (64KB). Estos procesadores fueron comunes en las décadas de 1970 y 1980 — Los más conocidos incluyen al Intel 8080 y 8085, Zilog Z80, MOS 6502 y 6510, y Motorola 6800. Hay que recalcar que estos procesadores son reconocidos como procesadores de 8 *bits*, pero con *espacio de direccionamiento* de 16 bits. El procesador empleado en las primeras PC, el Intel 8086, manejaba un direccionamiento de 20 bits (hasta 1024KB), pero al ser una arquitectura *real* de 16 bits requería del empleo de *segmentación* para alcanzar toda su memoria.

Con la llegada de la era de las *computadoras personales*, diversos fabricantes introdujeron a mediados de los años 1980 los procesadores con arquitectura de 32 bits. Por ejemplo, la arquitectura IA-32 de Intel tiene su inicio oficial con el procesador 80386 (o simplemente 386). Este procesador podía referenciar desde el punto de vista teórico hasta 2^{32} bytes (4GB)

de RAM. No obstante, debido a las limitaciones tecnológicas (y tal vez estratégicas) para producir memorias con esta capacidad, tomó más de veinte años para que las memorias ampliamente disponibles alcanzaran dicha capacidad.

Hoy en día, los procesadores dominantes son de 32 o 64 bits. En el caso de los procesadores de 32 bits, sus registros pueden referenciar hasta 4,294,967,296 bytes (4GB) de RAM, que está ya dentro de los parámetros de lo esperable hoy en día. Una arquitectura de 32 bits sin extensiones adicionales no puede emplear una memoria de mayor capacidad. No obstante, a través de un mecanismo llamado *PAE* (Extensión de Direcciones Físicas, *Physical Address Extension*) permite extender esto a rangos de hasta 2^{52} bytes a cambio de un nivel más de indirección.

Un procesador de 64 bits podría direccionar hasta 18,446,744,073,709,551,616 bytes (16 Exabytes). Los procesadores comercialmente hoy en día no ofrecen esta capacidad de direccionamiento principalmente por un criterio económico: Al resultar tan poco probable que exista un sistema con estas capacidades, los chips actuales están limitados a entre 2^{40} y 2^{48} bits — 1 y 256 terabytes. Esta restricción debe seguir teniendo sentido económico por muchos años aún.

6.1.2 Hardware: de la unidad de manejo de memoria (MMU)

Con la introducción de sistemas multitarea, es decir, dos o más programas ejecutándose, se vio la necesidad de tener más de un programa cargado en memoria. Esto lleva que el sistema operativo junto con información del programa a ejecutar debe resolver cómo ubicar los programas en la memoria física disponible.

Luego ha sido necesario emplear más memoria de la que está directamente disponible, con el propósito de ofrecer a los procesos más espacio de lo que puede direccionar /la arquitectura (hardware) empleada. Por otro lado, la abstracción de un espacio virtualmente ilimitado para realizar sus operaciones incluso cuando la memoria *real* es mucho menor a la solicitada, y por último, la ilusión de tener un bloque contiguo e ininterrumpido de memoria, cuando en realidad puede haber alta *fragmentación*.

Se explicará cómo la MMU cubre estas necesidades, y qué mecanismos emplea para lograrlo — Y qué cuidados se deben conservar, incluso como programadores de aplicaciones en lenguajes de alto nivel, para aprovechar de la mejor manera estas funciones (y evitar, por el contrario, que los programas se vuelvan lentos por no manejar la memoria correctamente).

La MMU es también la encargada de verificar que un proceso no tenga acceso a leer o modificar los datos de otro — Si el sistema operativo tuviera que verificar cada una de las instrucciones ejecutadas por un programa para evitar errores en el acceso a la memoria, la penalización en velocidad sería demasiado severa¹.

Una primer aproximación a la protección de acceso se implementa usando un *registro base* y un *registro límite*: Si la arquitectura ofrece dos registros del procesador que sólo pueden ser modificados por el sistema operativo (Esto es, el hardware define la modificación de dichos registros como una operación privilegiada que requiere estar ejecutando en *modo supervisor*), la MMU puede comparar sin penalidad *cada acceso a memoria* para verificar que esté en el rango permitido.

Por ejemplo, si a un proceso le fue asignado un espacio de memoria de 64K (65535 bytes) a partir de la dirección 504214 (492K), el *registro base* contendría 504214, y el *registro límite* 65535. Si hubiera una instrucción por parte de dicho proceso que solicitara una dirección menor a 504214 o mayor a 569749 (556K), la MMU lanzaría una excepción o *trampa* interrumpiendo el procesamiento, e indicando al sistema operativo que ocurrió

¹Y de hecho está demostrado que no puede garantizarse que una verificación estática sea suficientemente exhaustiva

una *violación de segmento (segmentation fault)*². El sistema operativo entonces procedería a terminar la ejecución del proceso, reclamando todos los recursos que tuviera asignados y notificando a su usuario.

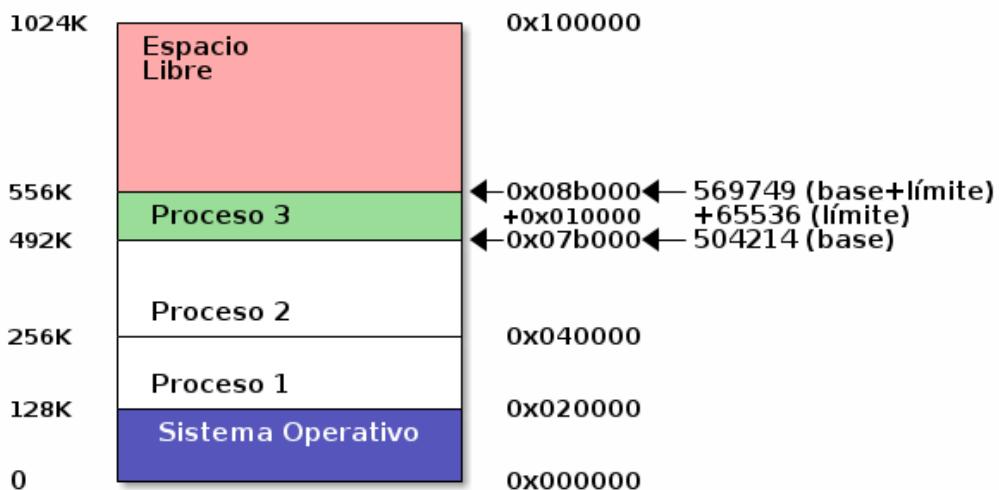


Figura 6.1: Espacio de direcciones válidas para el proceso 3 definido por un registro base y un registro límite

6.1.3 La memoria caché

Hay otro elemento en la actualidad se asume como un hecho: La memoria *caché*. Si bien su manejo es (casi) transparente para el sistema operativo, es muy importante mantenerlo en mente.

Conforme el procesador avanza en la ejecución de las instrucciones (aumentando el valor almacenado en el registro de conteo de instrucción), se producen accesos a memoria. Por un lado, tiene que buscar en memoria la siguiente instrucción a ejecutar. Por otro lado, estas instrucciones pueden requerirle uno o más operadores adicionales que deban ser leídos de la memoria. Por último, la instrucción puede requerir guardar su resultado en cierta dirección de memoria.

Hace años esto no era un problema — La velocidad del procesador estaba básicamente sincronizada con la del manejador de memoria, y el flujo podía mantenerse básicamente estable. Pero conforme los procesadores se fueron haciendo más rápidos, y conforme se ha popularizado el procesamiento en paralelo, la tecnología de la memoria no ha progresado a la misma velocidad. La memoria de alta velocidad es demasiado cara, e incluso las distancias de unos pocos centímetros se convierten en obstáculos insalvables por la velocidad máxima de los electrones viajando por pistas conductoras.

Cuando el procesador solicita el contenido de una dirección de memoria y esta no está aún disponible, *tiene que detener su ejecución (stall)* hasta que los datos estén disponibles. El CPU no puede, a diferencia del sistema operativo, “congelar” todo y guardar el estado para atender a otro proceso: Para el procesador, la lista de instrucciones a ejecutar es estrictamente secuencial, y todo tiempo que requiere esperar a una transferencia de datos es tiempo perdido.

La respuesta para reducir esa espera es la *memoria caché*. Esta es una memoria de alta velocidad, situada *entre* la memoria principal y el procesador propiamente, que guarda

²¿Por qué de segmento? Ver la sección 6.3

copias de las *páginas* que van siendo accesadas, partiendo del principio de la *localidad de referencia*:

Localidad temporal Es probable que un recurso que fue empleado recientemente vuelva a ser empleado en un futuro cercano.

Localidad espacial La probabilidad de que un recurso *aún no requerido* sea accesado es mucho mayor si fue requerido algún recurso cercano.

Localidad secuencial Un recurso, y muy particularmente la memoria, tiende a ser requerido de forma secuencial.

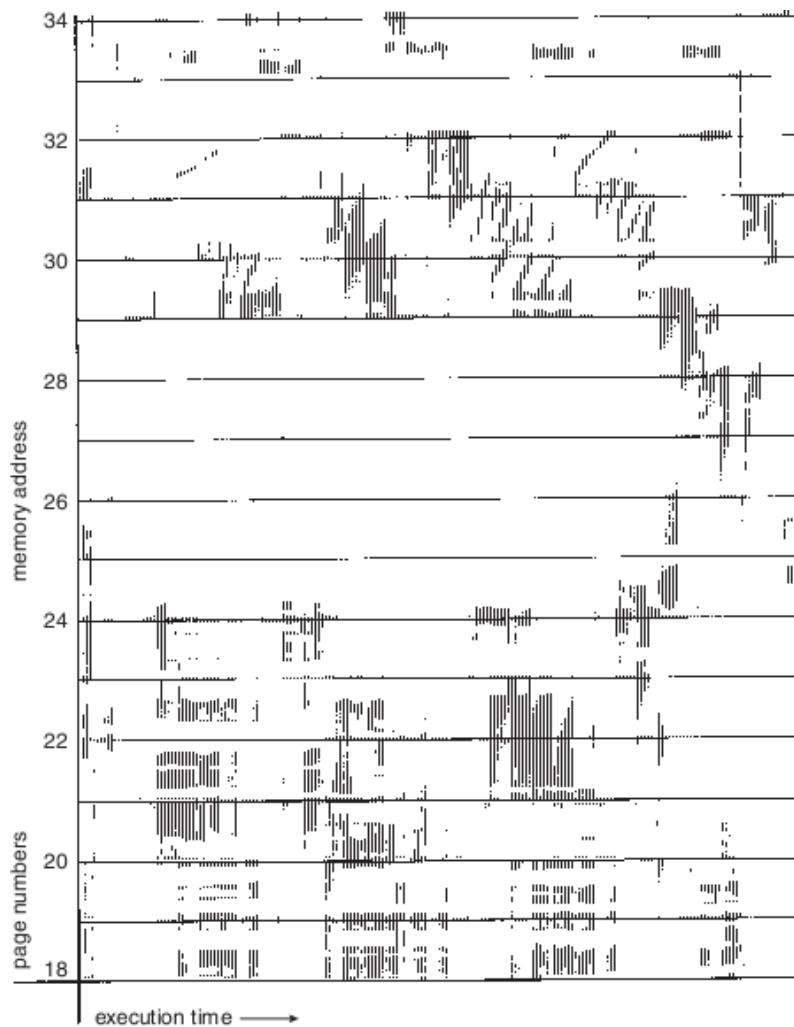


Figura 6.2: Patrones de acceso a memoria, demostrando la localidad espacial / temporal (Silberschatz, p.350)

Aplicando el concepto de localidad de referencia, cuando el procesador solicita al hardware determinada dirección de memoria, el hardware no sólo transfiere a la memoria caché el byte o palabra solicitado, sino que transfiere un bloque o *página* completo.

Cabe mencionar que hoy en día (particularmente desde que se detuvo la *guerra de los Megahertz*), parte importante del diferencial en precios de los procesadores líderes en el mercado es la cantidad de memoria caché de primero y segundo nivel con que cuentan.

6.1.4 El espacio en memoria de un proceso

Cuando un sistema operativo inicia un proceso, no se limita a volcar el archivo ejecutable a memoria, sino que tiene que proporcionar la estructura para que éste vaya guardando la información de estado relativa a su ejecución.

Sección (o segmento) de texto Es el nombre que recibe la imagen en memoria de las instrucciones a ser ejecutadas. Usualmente, la sección de texto ocupa las direcciones más bajas del espacio en memoria.

Sección de datos Espacio fijo preasignado para las variables globales y datos inicializados (como las cadena de caracteres por ejemplo). Este espacio es fijado en tiempo de compilación, y no puede cambiar (aunque los datos que cargados allí sí cambian en el tiempo de vida del proceso)

Espacio de libres Espacio de memoria que se emplea para la asignación dinámica de memoria durante la ejecución del proceso. Este espacio se ubica por encima de la sección de datos, y crece hacia arriba. Este espacio es conocido en inglés como el *Heap*. Cuando el programa es escrito en lenguajes que requieren *manejo dinámico manual de la memoria* (como C), esta área es la que se maneja a través de las llamadas de la familia de `malloc` y `free`. En lenguajes con gestión automática, esta área es monitoreada por los *recolectores de basura*.

Pila de llamadas Consiste en un espacio de memoria que se usa para almacenar la secuencia de funciones que han sido llamadas dentro del proceso, con sus parámetros, direcciones de retorno, variables locales, etc. La pila ocupa la parte más alta del espacio en memoria, y crece hacia abajo.

En inglés, la pila de llamadas es denominada *Stack*.

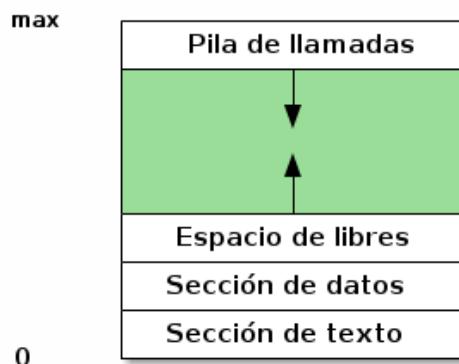


Figura 6.3: Regiones de la memoria para un proceso

6.1.5 Resolución de direcciones

Un programa compilado no emplea nombres simbólicos para las variables o funciones que llama³; el compilador, al convertir el programa a lenguaje máquina, las substituye por la dirección en memoria donde se encuentra la variable o la función⁴.

Ahora bien, en los sistemas actuales, los procesos requieren coexistir con otros, para lo cual las direcciones indicadas en el *texto* del programa pueden requerir ser traducidas al lugar *relativo al sitio de inicio del proceso en memoria* — Esto es, las direcciones son *resueltas* o traducidas. Existen diferentes estrategias de resolución, que se pueden clasificar a grandes rasgos⁵ en:

³Cuando se hace *ligado dinámico* a bibliotecas externas sí se mantiene la referencia por nombre, pero para los propósitos de esta sección, se habla de las referencias internas únicamente

⁴de hecho, una vez que el programa se pasa a lenguaje de máquina, no existe diferencia real entre la dirección que ocupa una variable o código ejecutable. La diferencia se establece por el uso que se dé a la referencia de memoria. En la sección 6.6.1 se abordará un ejemplo de cómo esto puede tener importantes consecuencias.

⁵Esta explicación simplifica muchos detalles; para el lector interesado en profundizar en este tema, se

En tiempo de compilación El texto del programa tiene la dirección *absoluta* de las variables y funciones. Esto era muy común en las computadoras previas al multiprocesamiento. En la arquitectura compatible con PC, el formato ejecutable .COM es un volcado de memoria directo de un archivo objeto con las direcciones indicadas de forma absoluta. Esto puede verse hoy principalmente en sistemas embebidos o de función específica.

En tiempo de carga Al cargarse a memoria el programa y antes de iniciar su ejecución, el *cargador* (componente del sistema operativo) actualiza las referencias a memoria dentro del texto para que apunten al lugar correcto — Claro está, esto depende de que el compilador indique dónde están todas las referencias a variables y funciones.

En tiempo de ejecución El programa nunca hace referencia a una ubicación absoluta de memoria, sino que lo hace siempre relativo a una *base* y un *desplazamiento* (offset). Esto permite que el proceso sea incluso reubicado en la memoria mientras está siendo ejecutado sin tener que sufrir cambios, pero requiere de hardware específico (como una MMU).

Esto es, los nombres simbólicos (por ejemplo, una variable llamada *contador*) para ser traducidos ya sea a ubicaciones en la memoria, pueden resolverse en tiempo de compilación (y quedar plasmada en el programa en disco con una ubicación explícita y definitiva: 510200), en tiempo de carga (sería guardada en el programa en disco como *inicio* + 5986 bytes, y el proceso de carga incluiría substituirla por la dirección resuelta a la suma del registro base, 504214, y el *desplazamiento*, 5986, esto es, 510200).

Por último, para emplear la resolución en tiempo de ejecución, se mantiene en las instrucciones a ser ejecutadas por el proceso la etiqueta relativa al módulo actual, *inicio* + 5986 bytes, y es resuelta cada vez que sea requerido.

6.2 Asignación de memoria contigua

En los sistemas de ejecución en lotes, así como en las primeras computadoras personales, sólo un programa se ejecutaba a la vez. Por lo que, más allá de la carga del programa y la satisfacción de alguna eventual llamada al sistema solicitando recursos, el sistema operativo no tenía que ocuparse de la asignación de memoria.

Al nacer los primeros sistemas operativos multitarea, se hizo necesario resolver cómo asignar el espacio en memoria a diferentes procesos.

6.2.1 Partición de la memoria

La primer respuesta, claro está, es la más sencilla: Asignar a cada programa a ser ejecutado un bloque *contiguo* de memoria de un tamaño fijo. En tanto los programas permitieran la resolución de direcciones en tiempo de carga o de ejecución, podrían emplear este esquema.

El sistema operativo emplearía una región específica de la memoria del sistema (típicamente la *región baja* — Desde la dirección en memoria 0x00000000 hacia arriba), y una vez terminado el espacio necesario para el núcleo y sus estructuras, el sistema asigna espacio a cada uno de los procesos. Si la arquitectura en cuestión permite limitar los segmentos disponibles a cada uno de los procesos (por ejemplo, con los registros *base* y *límite* discutidos anteriormente), esto sería suficiente para alojar en memoria a varios procesos y evitar que interfieran entre sí.

Desde la perspectiva del sistema operativo, cada uno de los espacios asignados a un proceso es una *partición*. Cuando el sistema operativo inicia, toda la memoria disponible es vista como un sólo bloque, y conforme se van ejecutando procesos, este bloque va siendo

recomienda el libro *Linkers and Loaders* (Ligadores y cargadores) de John R. Levine (1999). El libro está disponible en línea desde el sitio Web del autor, <http://www.iecc.com/linker/>

subdividido para satisfacer sus requisitos. Al cargar un programa el sistema operativo calcula cuánta memoria va a requerir a lo largo de su vida prevista. Esto incluye el espacio requerido para la asignación dinámica de memoria con la familia de funciones `malloc` y `free`.

Fragmentación

La *fragmentación* comienza a aparecer cuando más procesos terminan su ejecución, y el sistema operativo libera la memoria asignada a cada uno de ellos. A medida que los procesos finalizan, aparecen regiones de memoria disponible, *interrumpidas* por regiones de memoria usada por los procesos que aún se encuentran activos.

Si la computadora no tiene hardware específico que permita que los procesos resuelvan sus direcciones en tiempo de ejecución, el sistema operativo no puede reasignar los bloques existentes, y aunque pudiera hacerlo, mover un proceso entero en memoria puede resultar una operación costosa en tiempo de procesamiento.

Al crear un nuevo proceso, el sistema operativo tiene tres estrategias según las cuales podría asignarle uno de los bloques disponibles:

Primer ajuste El sistema toma el primer bloque con el tamaño suficiente para alojar el nuevo proceso. Este es el mecanismo más simple de implementar y el de más rápida ejecución. No obstante, esta estrategia puede causar el desperdicio de memoria, si el bloque no es exactamente del tamaño requerido.

Mejor ajuste El sistema busca entre todos los bloques disponibles cuál es el que mejor se ajusta al tamaño requerido por el nuevo proceso. Esto implica la revisión completa de la lista de bloques, pero permite que los bloques remanentes, una vez que se ubicó al nuevo proceso, sean tan pequeños como sea posible (esto es, que haya de hecho un *mejor ajuste*).

Peor ajuste El sistema busca cuál es el bloque más grande disponible, y se lo asigna al nuevo proceso. Empleando una estructura de datos como un *montículo*, esta operación puede ser incluso más rápida que la de primer espacio. Con este mecanismo se busca que los bloques que queden después de otorgarlos a un proceso sean tan grandes como sea posible, de cierto modo balanceando su tamaño.

La *fragmentación externa* se produce cuando hay muchos bloques libres entre bloques asignados a procesos; la *fragmentación interna* se refiere a la cantidad de memoria dentro de un bloque que nunca se usará — Por ejemplo, si el sistema operativo maneja *bloques* de 512 bytes y un proceso requiere sólo 768 bytes para su ejecución, el sistema le entregará dos bloques (1024 bytes), con lo cual desperdicia 256 bytes. En el peor de los casos, con un bloque de n bytes, un proceso podría solicitar $kn + 1$ bytes de memoria, desperdimando por fragmentación interna $n - 1$ bytes.

Según análisis estadísticos (Silberschatz, p.289), por cada N bloques asignados, se perderán del orden de $0.5N$ bloques por fragmentación interna y externa.

Compactación

Un problema importante que va surgiendo como resultado de esta fragmentación es que el espacio total libre de memoria puede ser mucho mayor que lo que requiere un nuevo proceso, pero al estar *fragmentada* en muchos bloques, éste no encontrará una partición contigua donde ser cargado.

Si los procesos emplean resolución de direcciones en tiempo de ejecución, cuando el sistema operativo comience a detectar un alto índice de fragmentación, puede lanzar una operación de *compresión* o *compactación*. Esta operación consiste en mover los contenidos en memoria de los bloques asignados para que ocupen espacios contiguos, permitiendo unificar varios bloques libres contiguos en uno solo.

La compactación tiene un costo alto — Involucra mover prácticamente la totalidad de la memoria (probablemente más de una vez por bloque).

Intercambio (*swap*) con el almacenamiento secundario

Siguiendo de cierto modo la lógica requerida por la compactación se encuentran los sistemas que utilizan *intercambio (swap)* entre la memoria primaria y secundaria. En estos sistemas, el sistema operativo puede *comprometer* más espacio de memoria del que tiene físicamente disponible. Cuando la memoria se acaba, el sistema suspende a un proceso (usualmente un proceso “bloqueado”) y almacena una copia de su imagen en memoria en almacenamiento secundario para luego poder restaurarlo.

Hay algunas restricciones a observar previo a suspender un proceso. Por ejemplo, se debe considerar si el proceso tiene pendiente alguna operación de entrada/salida, en la cual el resultado se deberá copiar en su espacio de memoria. Si el proceso resultara suspendido (retirándolo de la memoria principal), el sistema operativo no tendría dónde continuar almacenando estos datos conforme llegaran. Una estrategia ante esta situación podría ser que todas las operaciones se realicen únicamente a *buffers* (regiones de memoria de almacenamiento temporal) en el espacio del sistema operativo, y éste las transfiera el contenido del buffer al espacio de memoria del proceso suspendido una vez que la operación finalice.

Esta técnica se popularizó en los sistemas de escritorio hacia fines de los 1980 y principios de los 1990, en que las computadoras personales tenían típicamente entre 1 y 8MB de memoria.

Se debe considerar que las unidades de disco son del orden de decenas de miles de veces más lentas que la memoria, por lo que este proceso resulta muy caro — Por ejemplo, si la imagen en memoria de un proceso es de 100MB, bastante conservador hoy en día, y la tasa de transferencia sostenida al disco de 50MB/s, intercambiar un proceso al disco toma dos segundos. Cargarlo de vuelta a memoria toma otros dos segundos — Y a esto debe sumarse el tiempo de posicionamiento de la cabeza de lectura/escritura, especialmente si el espacio a emplear no es secuencial y contiguo. Resulta obvio por qué esta técnica ha caído en desuso conforme aumenta el tamaño de los procesos.

6.3 Segmentación

Al desarrollar un programa en un lenguaje de alto nivel, el programador usualmente no se preocupa por la ubicación en la memoria física de los diferentes elementos que lo componen. Esto se debe a que en estos lenguajes las variables y funciones son referenciadas por sus *nombres*, no por su ubicación⁶. No obstante, cuando se compila el programa para una arquitectura que soporte segmentación, el compilador ubicará a cada una de las secciones presentadas en la sección 6.1.4 en un segmento diferente.

Esto permite activar los mecanismos que evitan la escritura accidental de las secciones de memoria del proceso que no se deberían modificar (aquellas que contienen código o de sólo lectura), y permitir la escritura de aquellas que sí (en las cuales se encuentran las variables globales, la pila o *stack* y el espacio de asignación dinámica o *heap*).

Así, los elementos que conforman un programa se organizan en *secciones*: una sección contiene el espacio para las variables globales, otra sección contiene el código compilado, otra sección contiene la *tabla de símbolos*, etc.

Luego, cuando el sistema operativo crea un proceso a partir del programa, debe organizar el contenido del archivo ejecutable en memoria. Para ello carga en memoria algunas secciones del archivo ejecutable (como mínimo la sección para las variables globales y la sección de código) y puede configurar otras secciones como la pila o la sección de libres. Para garantizar la protección de cada una de estas secciones en la

⁶Al programar en lenguaje C por ejemplo, un programador puede trabajar a este mismo nivel de abstracción, puede referirse directamente a las ubicaciones en memoria de sus datos empleando *aritmética de apuntadores*.

memoria del proceso, el sistema puede definir que cada *sección* del programa se encuentra en un *segmento* diferente, con diferentes tipos de acceso.

La segmentación es un concepto que se aplica directamente a la arquitectura del procesador. Permite separar las regiones de la memoria *lineal* en *segmentos*, cada uno de los cuales puede tener diferentes permisos de acceso, como se explicará en la siguiente sección. La segmentación también ayuda a incrementar la *modularidad* de un programa: Es muy común que las bibliotecas *ligadas dinámicamente* estén representadas en segmentos independientes.

Un código compilado para procesadores que implementen segmentación siempre generará referencias a la memoria en un espacio *segmentado*. Este tipo de referencias se denominan *direcciones lógicas* y están formadas por un *selector* de segmento y un *desplazamiento* dentro del segmento. Para interpretar esta dirección, la MMU debe tomar el selector, y usando alguna estructura de datos, obtiene la dirección base, el tamaño del segmento y sus atributos de protección. Luego, aplicando el mecanismo explicado en secciones anteriores, toma la dirección base del segmento y le suma el desplazamiento para obtener una *dirección lineal física*.

La traducción de una dirección lógica a una dirección lineal puede fallar por diferentes razones: Si el segmento no se encuentra en memoria, ocurrirá una *excepción* del tipo *segmento no presente*. Por otro lado, si el desplazamiento especificado es mayor al tamaño definido para el segmento, ocurrirá una excepción del tipo *violación de segmento*.

6.3.1 Permisos

Una de las principales ventajas del uso de segmentación consiste en permitir que cada uno de los segmentos tenga un distinto *juego de permisos* para el proceso en cuestión: El sistema operativo puede indicar, por ejemplo, que el *segmento de texto* (el código del programa) sea de lectura y ejecución, mientras que las secciones de datos, libres y pila (donde se almacena y trabaja la información misma del programa) serán de lectura y escritura, pero la ejecución estará prohibida⁷.

De este modo, se puede evitar que un error en la programación resulte en que datos proporcionados por el usuario o por el entorno modifiquen el código que está siendo ejecutado.⁸ Es más, dado que el acceso de *ejecución* está limitado a sólo los segmentos cargados del disco por el sistema operativo, un atacante no podrá introducir código ejecutable tan fácilmente — Tendría que cargarlo como un segmento adicional con los permisos correspondientes.

La segmentación también permite distinguir *niveles* de acceso a la memoria: Para que un proceso pueda efectuar *llamadas al sistema*, debe tener acceso a determinadas estructuras compartidas del núcleo. Claro está, al ser memoria privilegiada, su acceso requiere que el procesador esté ejecutando en *modo supervisor*.

6.3.2 Intercambio parcial

Un uso muy común de la segmentación, particularmente en los sistemas de los 1980, era el de permitir que sólo *certas regiones* de un programa sean intercambiadas al disco: Si un programa está compuesto por porciones de código que nunca se ejecutarán aproximadamente al mismo tiempo en sucesión, puede separar su texto (e incluso los datos

⁷Si bien este es el manejo clásico, no es una regla inquebrantable: El código *automodificable* conlleva importantes riesgos de seguridad, pero bajo ciertos supuestos, el sistema debe permitir su ejecución. Además, muchos lenguajes de programación permiten la *metaprogramación*, que requiere la ejecución de código construido en tiempo de ejecución.

⁸Sin embargo, incluso bajo este esquema, dado que la *pila de llamadas (stack)* debe mantenerse como escribible, es común encontrar ataques que permiten modificar la *dirección de retorno* de una subrutina, como será descrito en la sección 6.6.1.

correspondientes) en diferentes segmentos.

A lo largo de la ejecución del programa, algunos de sus segmentos pueden no emplearse por largos períodos de tiempo. Éstos pueden ser enviadas al *espacio de intercambio (swap)* ya sea a solicitud del proceso o por iniciativa del sistema operativo.

Rendimiento

En la sección 6.2.1, donde se presenta el concepto de intercambio, se explicó que intercambiar un proceso completo resultaba demasiado caro. Cuando se tiene de un espacio de memoria segmentado, y muy particularmente cuando se usan bibliotecas de carga dinámica, la sobrecarga es mucho menor:

En primer término, se puede hablar de la cantidad de información a intercambiar: En un sistema que sólo maneja regiones contiguas de memoria, intercambiar un proceso significa mover toda su información al disco. En un sistema con segmentación, puede enviarse a disco cada uno de los segmentos por separado, según el sistema operativo lo juzgue necesario. Podría *sacar* de memoria a alguno de los segmentos, eligiendo no necesariamente al que más *estorbe* (esto es, el más grande), sino el que más probablemente no esté siendo utilizado: Emplear el principio de localidad de referencia para intercambiar al segmento *menos recientemente utilizado (LRU, Least Recently Used)*.

Además de esto, si se tiene un segmento de texto (sea el código programa base o alguna de las bibliotecas) y su acceso es de sólo lectura, una vez que éste fue copiado una vez al disco, ya no hace falta volver a hacerlo: Se tiene la certeza de que no será modificado por el proceso en ejecución, por lo que basta marcarlo como *no presente* en las tablas de segmentos en memoria para que cualquier acceso ocasione que el sistema operativo lo traiga de disco.

Por otro lado, si la biblioteca en cuestión reside en disco (antes de ser cargada) como una imagen directa de su representación en memoria, al sistema operativo le bastará identificar el archivo en cuestión al cargar el proceso; no hace falta siquiera cargarlo en la memoria principal y guardarla al área de intercambio, puede quedar referido directamente al espacio en disco en que reside el archivo.

Claro está, el acceso a disco sigue siendo una fuerte penalización cada vez que un segmento tiene que ser cargado del disco (sea del sistema de archivos o del espacio de intercambio), pero este mecanismo reduce dicha penalización, haciendo más atractiva la flexibilidad del intercambio por segmentos.

6.3.3 Ejemplificando

A modo de ejemplo, y conjuntando los conceptos presentados en esta sección, si un proceso tuviera la siguiente tabla de segmentos:

Segmento	Inicio	Tamaño	Permisos	Presente
0	15210	150	RWX	sí
1	1401	100	R	sí
2	965	96	RX	sí
3	-	184	W	no
4	10000	320	RWX	sí

En la columna de permisos, *R* se refiere a lectura, *W* a escritura y *X* a ejecución.

Un segmento que ha sido enviado al espacio de intercambio (en este caso, el 3), deja de estar presente en memoria y, por tanto, no tiene ya dirección de inicio registrada.

El resultado de hacer referencia a las siguientes direcciones y tipos de acceso:

Dirección lógica	Tipo de acceso	Dirección física
0-100	R	15310
2-82	X	1533
2-82	W	Atrapada: Violación de seguridad
2-132	R	Atrapada: Desplazamiento fuera de rango
3-15	W	Atrapada: Segmento faltante
3-130	R	Atrapada: Segmento faltante; violación de seguridad
4-130	X	10130
5-15	X	Atrapada: Segmento invalido

Cuando se atrapa una situación de excepción, el sistema operativo debe intervenir. Por ejemplo, la solicitud de un segmento inválido, de un desplazamiento mayor al tamaño del segmento, o de un tipo de acceso que no esté autorizado, típicamente llevan a la terminación del proceso, en tanto que una de segmento faltante (indicando un segmento que está en el espacio de intercambio) llevaría a la suspensión del proceso, lectura del segmento de disco a memoria, y una vez que éste estuviera listo, se permitiría continuación de la ejecución.

En caso de haber más de una excepción, como se observa en la solicitud de lectura de la dirección 3-94, el sistema debe reaccionar primero a la *más severa*: Si como resultado de esa solicitud iniciará el proceso de carga del segmento, sólo para abortar la ejecución del proceso al detectarse la violación de tipo de acceso, sería un desperdicio injustificado de recursos.

6.4 Paginación

La fragmentación externa y, por tanto, la necesidad de compactación pueden evitarse por completo empleando la *paginación*. Ésta consiste en que cada proceso está dividido en varios bloques de tamaño fijo (más pequeños que los segmentos) llamados *páginas*, dejando de requerir que la asignación sea de un área *contigua* de memoria. Claro está, esto requiere de mayor espacialización por parte del hardware, y mayor información relacionada a cada uno de los procesos: No basta sólo con indicar dónde inicia y dónde termina el área de memoria de cada proceso, sino que se debe establecer un *mapeo* entre la ubicación real (*física*) y la presentada a cada uno de los procesos (*lógica*). La memoria se presentará a cada proceso como si fuera de su uso exclusivo.

La memoria física se divide en una serie de *marcos (frames)*, todos ellos del mismo tamaño, y el espacio cada proceso se divide en una serie de *páginas (pages)*, del mismo tamaño que los marcos. La MMU se encarga del mapeo entre páginas y marcos a través de *tablas de páginas*.

Cuando se trabaja bajo una arquitectura que maneja paginación, las direcciones que maneja el CPU ya no son presentadas de forma absoluta. Los bits de cada dirección se separan en un *identificador de página* y un *desplazamiento*, de forma similar a lo presentado al hablar de resolución de instrucciones en tiempo de ejecución. La principal diferencia con lo entonces abordado es que cada proceso tendrá ya no un único espacio en memoria, sino una multitud de páginas.

El tamaño de los marcos (y, por tanto, las páginas) debe ser una *potencia de 2*, de modo que la MMU pueda discernir fácilmente la porción de una dirección de memoria que se refiere a la *página del desplazamiento*. El rango varía, según el hardware, entre los 512 bytes (2^9) y 16MB (2^{24}); al ser una potencia de 2, la MMU puede separar la dirección en memoria entre los primeros m bits (referentes a la página) y los últimos n bits (referentes al desplazamiento).

Para poder realizar este mapeo, la MMU requiere de una estructura de datos denominada *tabla de páginas (page table)*, que *resuelve* la relación entre páginas y marcos, convirtiendo una *dirección lógica* (en el espacio del proceso) en la *dirección física* (la ubicación en que *realmente* se encuentra en la memoria del sistema).

Se puede tomar como ejemplo para explicar este mecanismo el esquema presentado en la figura 6.9 (Silberschatz, p.292). Éste presenta un esquema minúsculo de paginación: Un *espacio de direccionamiento* de 32 bytes (5 bits), organizado en 8 páginas de 4 bytes cada una (esto es, la página es representada con los 3 bits *más significativos* de la dirección, y el desplazamiento con los 2 bits *menos significativos*).

El proceso que se presenta tiene una visión de la memoria como la columna del lado izquierdo: Para el proceso existen 4 páginas, y tiene sus datos distribuidos en orden desde la dirección 00000 (0) hasta la 01111 (15), aunque en realidad en el sistema éstas se encuentren desordenadas y ubicadas en posiciones no contiguas.

Cuando el proceso quiere referirse a la letra *f*, lo hace indicando la dirección 00101 (5). De esta dirección, los tres bits más significativos (001, 1 — Y para la computadora, lo *natural* es comenzar a contar por el 0) se refieren a la página número 1, y los dos bits menos significativos (01, 1) indican al *desplazamiento* dentro de ésta.

La MMU verifica en la tabla de páginas, y encuentra que la página 1 corresponde al marco número 6 (110), por lo que traduce la dirección lógica 00101 (5) a la física 11001 (26).

Se puede tomar la paginación como una suerte de resolución o traducción de direcciones en tiempo de ejecución, pero con una *base* distinta para cada una de las páginas.

6.4.1 Tamaño de la página

Ahora, si bien la fragmentación externa se resuelve al emplear paginación, el problema de la fragmentación interna persiste: Al dividir la memoria en bloques de longitud preestablecida de 2^n bytes, un proceso en promedio desperdiciará $\frac{2^n}{2}$ (y, en el peor de los casos, hasta $2^n - 1$). Multiplicando esto por el número de procesos que están en ejecución en todo momento en el sistema, para evitar que una proporción sensible de la memoria se pierda en fragmentación interna, se podría tomar como estrategia emplear un tamaño de página tan pequeño como fuera posible.

Sin embargo, la sobrecarga administrativa (el tamaño de la tabla de paginación) en que se incurre por gestionar demasiadas páginas pequeñas se vuelve una limitante en sentido opuesto:

- Las transferencias entre unidades de disco y memoria son mucho más eficientes si pueden mantenerse como recorridos continuos. El controlador de disco puede responder a solicitudes de acceso directo a memoria (DMA) siempre que tanto los fragmentos en disco como en memoria sean continuos; fragmentar la memoria demasiado jugaría en contra de la eficiencia de estas solicitudes.
- El bloque de control de proceso (PCB) incluye la información de memoria. Entre más páginas tenga un proceso (aunque estas fueran muy pequeñas), más grande es su PCB, y más información requerirá intercambiar en un cambio de contexto.

Estas consideraciones opuestas apuntan a que se debe mantener el tamaño de página más grande, y se regulan con las primeras expuestas en esta sección.

Hoy en día, el tamaño habitual de las páginas es de 4KB u 8KB (2^{12} o 2^{13} bytes). Hay algunos sistemas operativos que soportan múltiples tamaños de página — Por ejemplo, Solaris puede emplear páginas de 8KB y 4MB (2^{13} o 2^{22} bytes), dependiendo del tipo de información que se declare que almacenarán.

6.4.2 Almacenamiento de la tabla de páginas

Algunos de los primeros equipos en manejar memoria paginada empleaban un conjunto especial de registros para representar la tabla de páginas. Esto era posible dado que

eran sistemas de 16 bits, con páginas de 8KB (2^{13}). Esto significa que en el sistema había únicamente 8 páginas posibles (2^3), por lo que resultaba sensato dedicar un registro a cada una.

En los equipos actuales, mantener la tabla de páginas en registros resultaría claramente imposible: Teniendo un procesador de 32 bits, e incluso si se definiera un tamaño de página *muy* grande (por ejemplo, 4MB), existirían 1024 páginas posibles⁹; con un tamaño de páginas mucho más común (4KB, 2^{12} bytes), la tabla de páginas llega a ocupar 5MB.¹⁰ Los registros son muy rápidos, sin embargo, son correspondientemente muy caros. El manejo de páginas más pequeñas (que es lo normal), y muy especialmente el uso de espacios de direccionamiento de 64 bits, harían prohibitivo este enfoque. Además, nuevamente, cada proceso tiene una tabla de páginas distinta — Se haría necesario hacer una transferencia de información muy grande en cada cambio de contexto.

Otra estrategia para enfrentar esta situación es almacenar la propia tabla de páginas en memoria, y apuntar al inicio de la tabla con un juego de registros especiales: el *registro de base de la tabla de páginas* (PTBR, *Page Table Base Register*) y el *registro de longitud de la tabla de páginas* (PTLR, *Page Table Length Register*).¹¹ De esta manera, en el cambio de contexto sólo hay que cambiar estos dos registros, y además se cuenta con un espacio muy amplio para guardar las tablas de páginas que se necesiten. El problema con este mecanismo es la velocidad: Se estaría penalizando a *cada acceso a memoria* con un acceso de memoria adicional — Si para resolver una dirección lógica a su correspondiente dirección física hace falta consultar la tabla de páginas en memoria, el tiempo efectivo de acceso a memoria se duplica.

El buffer de traducción adelantada (TLB)

La salida obvia a este dilema es el uso de un caché. Sin embargo, más que un caché genérico, la MMU utiliza un caché especializado en el tipo de información que maneja: El *buffer de traducción adelantada* o *anticipada* (Translation Lookaside Buffer, TLB). El TLB es una tabla asociativa (un *hash*) en memoria de alta velocidad, una suerte de registros residentes dentro de la MMU, donde las *llaves* son las páginas y los *valores* son los marcos correspondientes. De este modo, las búsquedas se efectúan en tiempo constante.

El TLB típicamente tiene entre 64 y 1024 entradas. Cuando el procesador efectúa un acceso a memoria, si la página solicitada está en el TLB, la MMU tiene la dirección física de inmediato.¹² En caso de no encontrarse la página en el TLB, la MMU lanza un *fallo de página* (*page fault*), con lo cual consulta de la memoria principal cuál es el marco correspondiente. Esta nueva entrada es agregada al TLB; por las propiedades de *localidad de referencia* que se presentaron anteriormente, la probabilidad de que las regiones más empleadas de la memoria durante un área específica de ejecución del programa sean cubiertas por relativamente pocas entradas del TLB son muy altas.

Como sea, dado que el TLB es limitado, es necesario explicitar un mecanismo que indique dónde guardar las nuevas entradas una vez que el TLB está lleno y se produce un fallo de página. Uno de los esquemas más comunes es emplear la entrada *menos recientemente utilizada* (LRU, *Least Recently Used*), nuevamente apelando a la localidad de referencia. Esto tiene como consecuencia necesaria que debe haber un mecanismo que contabilice los accesos dentro del TLB (lo cual agrega tanto latencia como costo). Otro

⁹4MB es 2^{22} bytes; $\frac{2^{32}}{2^{22}} = 2^{10} = 1024$

¹⁰ $\frac{2^{32}}{2^{12}} = 2^{20} = 1048576$, cada entrada con un mínimo de 20 bits para la página y 20 bits para el marco. ¡La tabla de páginas misma ocuparía 1280 páginas!

¹¹¿Por qué es necesario el segundo? Porque es prácticamente imposible que un proceso emplee su espacio de direccionamiento completo; al indicar el límite máximo de su tabla de páginas por medio del PTLR se evita desperdiciar grandes cantidades de memoria indicando todo el espacio no utilizado.

¹²Indica Silberschatz (p.295) que el tiempo efectivo de acceso puede ser 10% superior al que tomaría sin emplear paginación.

mecanismo (con obvias desventajas) es el reemplazar una página al azar. Se explicará con mayor detalle más adelante algunos de los mecanismos más empleados para este fin, comparando sus puntos a favor y en contra.

Subdividiendo la tabla de páginas

Incluso empleando un TLB, el espacio empleado por las páginas sigue siendo demasiado grande. Si se considera un escenario más frecuente que el propuesto anteriormente: Empleando un procesador con espacio de direccionamiento de 32 bits, y un tamaño de página estándar (4KB, 2^{12}), se tendría 1,048,576 (2^{20}) páginas. Si cada entrada de la página ocupa 40 bits¹³ (esto es, 5 bytes), cada proceso requeriría de 5MB (5 bytes por cada una de las páginas) sólamente para representar su mapeo de memoria. Esto, especialmente en procesos pequeños, resultaría más gravoso para el sistema que los beneficios obtenidos de la paginación.

Aprovechando que la mayor parte del espacio de direccionamiento de un proceso está típicamente vacío (la pila de llamadas y el heap), se puede subdividir el identificador de página en dos (o más) niveles, por ejemplo, separando una dirección de 32 bits en una *tabla externa* de 10 bits, una *tabla interna* de 10 bits, y el *desplazamiento* de 12 bits.

Este esquema funciona adecuadamente para computadoras con direccionamiento de hasta 32 bits. Sin embargo, se debe considerar que cada nivel de páginas conlleva un acceso adicional a memoria en caso de fallo de página — Emplear paginación jerárquica con un nivel externo y uno interno implica que un fallo de página *triplica* (y no duplica, como sería con un esquema de paginación directa) el tiempo de acceso a memoria. Para obtener una tabla de páginas manejable bajo los parámetros aquí descritos en un sistema de 64 bits, se puede *septuplicar* el tiempo de acceso (cinco accesos *en cascada* para fragmentos de 10 bits, y un tamaño de página de 14 bits, mas el acceso a la página destino).

Otra alternativa es el emplear *funciones digestoras (hash functions)*¹⁴ para mapear cada una de las páginas a un *espacio muestral* mucho más pequeño. Cada página es mapeada a una lista de correspondencias simples¹⁵.

Un esquema basado en funciones digestoras ofrece características muy deseables: El tamaño de la tabla de páginas puede variar según crece el uso de memoria de un proceso (aunque esto requiera recalcular la tabla con diferentes parámetros) y el número de accesos a memoria en espacios tan grandes como el de un procesador de 64 bits se mantiene mucho más tratable. Sin embargo, por la alta frecuencia de accesos a esta tabla, debe elegirse un algoritmo digestor muy ágil para evitar que el tiempo que tome calcular la posición en la tabla resulte significativo frente a las alternativas.

6.4.3 Memoria compartida

Hay muchos escenarios en que diferentes procesos pueden beneficiarse de compartir áreas de su memoria. Uno de ellos es como mecanismo de comunicación entre procesos (*IPC, Inter Process Communication*), en que dos o más procesos pueden intercambiar estructuras de datos complejas sin incurrir en el costo de copiado que implicaría copiarlas a través del sistema operativo.

Otro caso, mucho más frecuente, es el de *compartir código*. Si un mismo programa es ejecutado varias veces, y dicho programa no emplea mecanismos de *código auto-modificable*,

¹³20 bits identificando a la página y otros 20 bits identificando al marco; omitiendo aquí la necesidad de alinear los accesos a memoria a *bytes* individuales, que lo aumentarían a 24

¹⁴Una función digestora puede definirse como $H : U \rightarrow M$, una función que *mapea* o *proyecta* al conjunto U en un conjunto M mucho menor; una característica muy deseable de toda función hash es que la *distribución resultante* en M resulte homogénea y tan poco dependiente de la secuencialidad de la entrada como sea posible.

¹⁵A una lista y no a un valor único dado que una función digestora es necesariamente proclive a presentar *colisiones*; el sistema debe poder resolver dichas colisiones sin pérdida de información.

no tiene sentido que las páginas en que se representa cada una de dichas instancias ocupe un marco independiente — El sistema operativo puede asignar a páginas de diversos procesos *el mismo conjunto de marcos*, con lo cual puede aumentar la capacidad percibida de memoria.

Y si bien es muy común compartir los *segmentos de texto* de los diversos programas que están en un momento dado en ejecución en la computadora, este mecanismo es todavía más útil cuando se usan *bibliotecas del sistema*: Hay bibliotecas que son empleadas por una gran cantidad de programas¹⁶.

Claro está, para ofrecer este modelo, el sistema operativo debe garantizar que las páginas correspondientes a las *secciones de texto* (el código del programa) sean de sólo lectura.

Un programa que está programado y compilado de forma que permita que todo su código sea de sólo lectura es *reentrante*, dado que posibilita que diversos procesos entren a su espacio en memoria sin tener que sincronizarse con otros procesos que lo estén empleando.

Copiar al escribir (*Copy on Write, CoW*)

En los sistemas Unix, el mecanismo más frecuentemente utilizado para crear un nuevo proceso es el empleo de la llamada al sistema `fork()`. Cuando es invocado por un proceso, el sistema operativo crea a un nuevo proceso *idéntico* al que lo llamó, diferenciándose únicamente en *el valor entregado* por la llamada a `fork()`. Si ocurre algún error, el sistema entrega un número negativo (indicando la causa del error). En caso de ser exitoso, el proceso nuevo (o proceso *hijo*) recibe el valor `0`, mientras que el proceso preexistente (o proceso *padre*) recibe el PID (número identificador de proceso) del hijo. Es frecuente encontrar el siguiente código:

```
/* (...) */
int pid;
/* (...) */
pid = fork();
if (pid == 0) {
    /* Soy el proceso hijo */
    /* (...) */
} else if (pid < 0) {
    /* Ocurrió un error, no se creó un proceso hijo */
} else {
    /* Soy el proceso padre */
    /* La variable 'pid' tiene el PID del proceso hijo */
    /* (...) */
}
```

Este método es incluso utilizado normalmente para crear nuevos procesos, transfiriendo el *ambiente* (variables, por ejemplo, que incluyen cuál es la *entrada* y *salida* estándar de un proceso, esto es, a qué terminal están conectados, indispensable en un sistema multiusuario). Frecuentemente, la siguiente instrucción que ejecuta un proceso hijo es `execve()`, que carga a un nuevo programa sobre del actual y transfiere la ejecución a su primer instrucción.

¹⁶Algunos ejemplos sobresalientes podrían ser la `libc` o `glibc`, que proporciona las funciones estándar del lenguaje C y es, por tanto, requerida por casi todos los programas del sistema; los diferentes entornos gráficos (en los Unixes modernos, los principales son Qt y Gtk); bibliotecas para el manejo de cifrado (`openssl`), compresión (`zlib`), imágenes (`libpng`, `libjpeg`), etc.

Cuesta trabajo comprender el por qué de esta lógica si no es por el empleo de la memoria compartida: El costo de `fork()` en un sistema Unix es muy bajo, se limita a crear las estructuras necesarias en la memoria del núcleo. Tanto el proceso padre como el proceso hijo comparten *todas* sus páginas de memoria — Sin embargo, siendo dos procesos independientes, no deben poder modificarse más que por los canales explícitos de comunicación entre procesos.

Esto ocurre así gracias a un mecanismo llamado *Copiar al escribir* (*Copy on Write* o *CoW*). Las páginas de memoria de ambos procesos son las mismas *mientras sean sólo leídas*. Sin embargo, si uno de los procesos modifica a cualquier dato en una de estas páginas, ésta se copia a un nuevo marco, y deja de ser una página compartida. El resto de las páginas seguirá siendo compartido. Esto se puede lograr marcando *todas* las páginas compartidas como *sólo lectura*, con lo cual cuando uno de los dos procesos intente modificar la información de alguna página se generará un fallo. El sistema operativo, al notar que esto ocurre sobre un espacio *CoW*, en vez de responder al fallo terminando al proceso, copiará sólo la página en la cual se encuentra la dirección de memoria que causó el fallo, y esta vez marcará la página como *lectura y escritura*.

Incluso cuando se ejecutan nuevos programas a través de `execve()`, es posible que una buena parte de la memoria se mantenga compartida, por ejemplo, al referirse a copias de bibliotecas de sistema.

6.5 Memoria virtual

Varios de los aspectos mencionados en la sección 6.4 (*Paginación*) van conformando a lo que se conoce como *memoria virtual*: En un sistema que emplea paginación, un proceso no conoce su dirección en memoria relativa a otros procesos, sino que trabajan con una *idealización* de la memoria, en la cual ocupan el espacio completo de direccionamiento, desde el cero hasta el límite lógico de la arquitectura, independientemente del tamaño físico de la memoria disponible.

Y si bien en el modelo mencionado de paginación los diferentes procesos pueden *compartir* regiones de memoria y *direccionar* más memoria de la físicamente disponible, no se ha presentado aún la estrategia que se emplearía cuando el total de páginas solicitadas por todos los procesos activos en el sistema superara el total de espacio físico. Es ahí donde entra en juego la *memoria virtual*: Para ofrecer a los procesos mayor espacio en memoria del que existe físicamente, el sistema emplea espacio en *almacenamiento secundario* (típicamente, disco duro), a través de un esquema de *intercambio* (*swap*) guardando y trayendo páginas enteras.

Es importante apuntar que la memoria virtual es gestionada *de forma automática y transparente* por el sistema operativo. No se hablaría de memoria virtual, por ejemplo, si un proceso pide explícitamente intercambiar determinadas páginas.

Puesto de otra manera: Del mismo modo que la segmentación (sección 6.3) permitió hacer mucho más cómodo y útil al intercambio (6.2.1) a través del intercambio parcial (6.3.2), permitiendo que continuara la ejecución del proceso incluso con ciertos segmentos *intercambiados* (*swapeados*) a disco, la memoria virtual lo hace aún más conveniente al aumentar la *granularidad* del intercambio: Ahora ya no se enviarán a disco secciones lógicas completas del proceso (segmentos), sino que se podrá reemplazar página por página, aumentando significativamente el rendimiento resultante. Al emplear la memoria virtual, de cierto modo la memoria física se vuelve sólo una *proyección parcial* de la memoria lógica, potencialmente mucho mayor a ésta.

Técnicamente, cuando se habla de memoria virtual, no se está haciendo referencia a un *intercambiador* (*swapper*), sino al *paginador*.

6.5.1 Paginación sobre demanda

La memoria virtual entra en juego desde la carga misma del proceso. Se debe considerar que existe una gran cantidad de *código durmiente* o *inalcanzable*: Código que sólo se emplea eventualmente, como el que responde ante una situación de excepción o el que se emplea sólo ante circunstancias particulares (por ejemplo, la exportación de un documento a determinados formatos, o la verificación de que no haya tareas pendientes antes de cerrar un programa). Y si bien a una computadora le sería imposible ejecutar código que no esté cargado en memoria,¹⁷ el código sí puede comenzar ejecutarse sin estar *completamente* en memoria: Basta con haber cargado la página donde están las instrucciones que permiten continuar con su ejecución actual.

La *paginación sobre demanda* significa que, para comenzar a ejecutar un proceso, el sistema operativo carga sólamente la *porción necesaria* para comenzar la ejecución (posiblemente una página o ninguna), y que a lo largo de la ejecución, el paginador es *flojo*:¹⁸ Sólo carga a memoria las páginas cuando van a ser utilizadas. Al emplear un paginador *flojo*, las páginas que no sean requeridas nunca serán siquiera cargadas a memoria.

La estructura empleada por la MMU para implementar un paginador *flojo* es muy parecida a la descrita al hablar del buffer de traducción adelantada (sección 6.4.2): La *tabla de páginas* incluirá un *bit de validez*, indicando para cada página del proceso si está presente o no en memoria. Si el proceso intenta emplear una página que esté marcada como no válida, esto causa un fallo de página, que lleva a que el sistema operativo lo suspenda y traiga a memoria la página solicitada para luego continuar con su ejecución:

1. Verifica en el PCB si esta solicitud corresponde a una página que ya ha sido asignada a este proceso.
2. En caso de que la referencia sea inválida, se termina el proceso.
3. Procede a traer la página del disco a la memoria. El primer paso es buscar un marco disponible (por ejemplo, a través de una tabla de asignación de marcos)
4. Solicita al disco la lectura de la página en cuestión hacia el marco especificado
5. Una vez que finaliza la lectura de disco, modifica tanto al PCB como al TLB para indicar que la tabla está en memoria.
6. Termina la suspensión del proceso, continuando con la instrucción que desencadenó al fallo. El proceso puede continuar sin notar que la página había sido intercambiada.

Llevando este proceso al extremo, se puede pensar en un sistema de *paginación puramente sobre demanda* (*pure demand paging*): En un sistema así, *ninguna* página llegará al espacio de un proceso si no es a través de un fallo de página. Un proceso, al iniciarse, comienza su ejecución *sin ninguna página en memoria*, y con el apuntador de siguiente instrucción del procesador apuntando a una dirección que no está en memoria (la dirección de la rutina de *inicio*). El sistema operativo responde cargando esta primer página, y conforme avanza el flujo del programa, el proceso irá ocupando el espacio real que empleará.

6.5.2 Rendimiento

La paginación sobre demanda puede impactar fuertemente el rendimiento de un proceso. Se ha explicado ya que un acceso a disco es varios miles de veces más lento que el acceso a memoria. Es posible calcular el tiempo de acceso efectivo a memoria (t_e) a partir de la probabilidad que en un proceso se presente un fallo de página ($0 \leq p \leq 1$), conociendo el tiempo de acceso a memoria (t_a) y el tiempo que toma atender a un fallo de página (t_f):

¹⁷Una computadora basada en la arquitectura von Neumann, como prácticamente todas las existen hoy en día, no puede *ver* directamente más que a la memoria principal

¹⁸En cómputo, muchos procesos pueden determinarse como *ansiosos* (*eager*), cuando buscan realizar todo el trabajo que puedan desde el inicio, o *flojos* (*lazy*), si buscan hacer el trabajo mínimo en un principio y diferir para más tarde tanto como sea posible

$$t_e = (1 - p)t_a + pt_f$$

Ahora bien, dado que t_a ronda hoy en día entre los 10 y 200ns, mientras que t_f está más bien cerca de los 8ms (la latencia típica de un disco duro es de 3ms, el tiempo de posicionamiento de cabeza de 5ms, y el tiempo de transferencia es de 0.05ms), para propósitos prácticos se puede ignorar a t_a . Con los valores presentados, seleccionando el mayor de los t_a presentados, si sólo un acceso a memoria de cada 1000 ocasiona un fallo de página (esto es, $p = \frac{1}{1000}$):

$$t_e = (1 - \frac{1}{1000}) \times 200\text{ns} + \frac{1}{1000} \times 8,000,000\text{ns}$$

$$t_e = 199,8\text{ns} + 8000\text{ns} = 8199,8\text{ns}$$

Esto es, en promedio, se tiene un tiempo efectivo de acceso a memoria 40 veces mayor a que si no se empleara este mecanismo. Con estos mismos números, para mantener la degradación de rendimiento por acceso a memoria por debajo del 10%, se debería reducir la probabilidad de fallos de página a $\frac{1}{399,990}$.

Cabe mencionar que este impacto al rendimiento no necesariamente significa que una proporción relativamente alta de fallos de página para un proceso impacte negativamente a todo el sistema — El mecanismo de paginación sobre demanda permite, al no requerir que se tengan en memoria todas las páginas de un proceso, que haya *más procesos activos* en el mismo espacio en memoria, aumentando el grado de multiprogramación del equipo. De este modo, si un proceso se ve obligado a esperar por 8ms a que se resuelva un fallo de página, durante ese tiempo pueden seguirse ejecutando los demás procesos.

Acomodo de las páginas en disco

El cálculo recién presentado, además, asume que el acomodo de las páginas en disco es óptimo. Sin embargo, si para llegar a una página hay que resolver la dirección que ocupa en un sistema de archivos (posiblemente navegar una estructura de directorio), y si el espacio asignado a la memoria virtual es compartido con los archivos en disco, el rendimiento sufrirá adicionalmente.

Una de las principales deficiencias estructurales en este sentido de los sistemas de la familia Windows es que el espacio de almacenamiento se asigna en el espacio libre del sistema de archivos. Esto lleva a que, conforme crece la fragmentación del disco, la memoria virtual quede esparcida por todo el disco duro. La generalidad de sistemas tipo Unix, en contraposición, reservan una partición de disco *exclusivamente* para paginación.

6.5.3 Reemplazo de páginas

Si se aprovechan las características de la memoria virtual para aumentar el grado de multiprogramación, como se explicó en la sección anterior, se presenta un problema: Al *sobre-comprometer* memoria, en determinado momento, los procesos que están en ejecución pueden caer en un patrón que requiera cargarse a memoria física páginas por un mayor uso de memoria que el que hay físicamente disponible.

Y si se tiene en cuenta que uno de los objetivos del sistema operativo es otorgar a los usuarios la *ilusión* de una computadora dedicada a sus procesos, no sería aceptable terminar la ejecución de un proceso ya aceptado y cuyos requisitos han sido aprobados porque no existe suficiente memoria. Se hace necesario encontrar una forma justa y adecuada de llevar a cabo un *reemplazo de páginas* que permita continuar satisfaciendo sus necesidades.

El reemplazo de páginas es una parte fundamental de la paginación, ya que es la pieza que posibilita una verdadera separación entre memoria lógica y física. El mecanismo

básico a ejecutar es simple: Si todos los marcos están ocupados, el sistema deberá encontrar una página que pueda liberar (una *página víctima*) y grabarla al espacio de intercambio en el disco. Luego, se puede emplear el espacio recién liberado para cargar la página requerida, y continuar con la ejecución del proceso.

Esto implica a una *doble* transferencia al disco (una para grabar la página víctima y una para traer la página de reemplazo), y por tanto, a una doble demora.

Se puede, con un mínimo de *burocracia* adicional (aunque requiere de apoyo de la MMU): implementar un mecanismo que disminuya la probabilidad de tener que realizar esta doble transferencia: Agregar un *bit de modificación* o *bit de página sucia* (*dirty bit*) a la tabla de páginas. Este bit se marca como apagado siempre que se carga una página a memoria, y es automáticamente encendido por hardware cuando se realiza un acceso de escritura a dicha página.

Cuando el sistema operativo elige una página víctima, si su *bit de página sucia* está encendido, es necesario grabarla al disco, pero si está apagado, se garantiza que la información en disco es idéntica a su copia en memoria, y permite ahorrar la mitad del tiempo de transferencia.

Ahora bien, ¿cómo decidir qué páginas remplazar marcándolas como *victimas* cuando hace falta? Para esto se debe implementar un *algoritmo de reemplazo de páginas*. La característica que se busca en este algoritmo es que, para una patrón de accesos dado, permita obtener el menor número de fallos de página.

De la misma forma como se realizó la descripción de los algoritmos de planificación de procesos, para analizar los algoritmos de reemplazo se usará una *cadena de referencia*. Esto es, sobre una lista de referencias a memoria. Estas cadenas modelan el comportamiento de un conjunto de procesos en el sistema, y, obviamente, diferentes comportamientos llevarán a diferentes resultados.

Hacer un volcado y trazado de ejecución en un sistema real puede dar una enorme cantidad de información, del orden de un millón de accesos por segundo. Para reducir esta información en un número más tratable, se puede simplificar basado en que no interesa cada referencia a una *dirección* de memoria, sino cada referencia a una *página* diferente.

Además, varios accesos a direcciones de memoria en la misma página no causan efecto en el estado. Se puede tomar como un sólo acceso a todos aquellos que ocurren de forma consecutiva (esto es, sin llamar a ninguna otra página, no es necesario que sean en instrucciones consecutivas) a una misma página.

Para analizar a un algoritmo de reemplazo, si se busca la cantidad de fallos de página producidos, además de la cadena de referencia, es necesario conocer la cantidad de páginas y marcos del sistema que se está modelando. Por ejemplo, si para la siguiente cadena:

1, 4, 3, 4, 1, 2, 4, 2, 1, 3, 1, 4

Al recorrer esta cadena en un sistema con cuatro o más marcos, sólo se presentarían cuatro fallos (el fallo inicial que hace que se cargue por primera vez cada una de las páginas). Si, en el otro extremo, se cuenta con sólo un marco, se presentarían 12 fallos, dado que a cada solicitud se debería reemplazar el único marco disponible. El rendimiento evaluado sería en los casos de que se cuenta con dos o tres marcos.

Un fenómeno interesante que se presenta con algunos algoritmos es la *anomalía de Belady*, publicada en 1969: Si bien la lógica indica que a mayor número de marcos disponibles se tendrá una menor cantidad de fallos de página, como lo ilustra la figura 6.17, con algunas de cadenas de referencia y bajo ciertos algoritmos puede haber una *regresión* o degradación, en la cual la cantidad de fallos aumenta aún con una mayor cantidad de marcos, como se puede ver en la figura 6.18.

Es importante recalcar que si bien la anomalía de Belady se presenta como un problema importante ante la evaluación de los algoritmos, en el texto de Luis La Red (p.559-569) se

puede observar que en simulaciones con características más cercanas a las de los patrones reales de los programas, su efecto observado es prácticamente nulo.

Para los algoritmos que se presentan a continuación, se asumirá una memoria con tres marcos, y con la siguiente cadena de referencia:

7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1

Primero en entrar, primero en salir (FIFO)

El algoritmo de más simple y de obvia implementación es, nuevamente, el FIFO: Al cargar una página en memoria, se toma nota de en qué momento fue cargada, y cuando llegue el momento de reemplazar una página vieja, se elige la que haya sido cargada hace más tiempo.

Partiendo de un estado inicial en que las tres páginas están vacías, necesariamente las tres primeras referencias a distintas páginas de memoria (7, 0, 1) causarán fallos de página. La siguiente (2) causará uno, pero la quinta referencia (0) puede ser satisfecha sin requerir una nueva transferencia.

La principal ventaja de este algoritmo es, como ya se ha mencionado, la simplicidad, tanto para programarlo como para comprenderlo. Su implementación puede ser tan simple como una lista ligada circular, cada elemento que va recibiendo se agrega en el último elemento de la lista, y se “empuja” el apuntador para convertirlo en la cabeza. Su desventaja, claro está, es que no toma en cuenta a la historia de las últimas solicitudes, por lo que puede causar un bajo rendimiento. Todas las páginas tienen la misma probabilidad de ser reemplazadas, sin importar su frecuencia de uso.

Con las condiciones aquí presentadas, un esquema FIFO causará 15 fallos de página en un total de 20 accesos requeridos.

El algoritmo FIFO es vulnerable a la anomalía de Belady. La figura 6.18 ilustra este fenómeno al pasar de 3 a 4 marcos.

La prevalencia de cadenas que desencadenan la anomalía de Belady fue uno de los factores principales que llevaron al diseño de nuevos algoritmos de reemplazo de páginas.

Reemplazo de páginas óptimo (OPT, MIN)

Un segundo algoritmo, de interés puramente teórico, fue propuesto, y es típicamente conocido como OPT o MIN. Bajo este algoritmo, el enunciado será elegir como página víctima a aquella página que *no vaya a ser utilizada* por un tiempo máximo (o nunca más).

Si bien este algoritmo está demostrado como óptimo o mínimo, se mantiene como curiosidad teórica porque requiere conocimiento a priori de las necesidades a futuro del sistema — Y si esto es impracticable ya en los algoritmos de despachadores, lo será mucho más con un recurso de reemplazo tan dinámico como la memoria.

Su principal utilidad reside en que ofrece una cota mínima: Calculando el número de fallos que se presentan al seguir OPT, es posible ver qué tan cercano resulta otro algoritmo respecto al caso óptimo. Para esta cadena de referencia, y con tres páginas, se tiene un total de nueve fallos.

Menos recientemente utilizado (LRU)

Este esquema se ha revisado en diversos mecanismos relacionados con la administración de memoria. Busca acercarse a OPT *prediciendo* cuándo será la próxima vez en que se emplee cada una de las páginas que tiene en memoria basado en la *historia reciente* de su ejecución.

Cuando necesita elegir una página víctima, LRU elige la página que no ha sido empleada hace más tiempo.

Para la cadena de referencia, LRU genera 12 fallos, en el punto medio entre OPT y FIFO.

Una observación interesante puede ser que para una cadena S y su *cadena espejo* (invertida) R^S , el resultado de evaluar S por LRU es igual al de evaluar R^S por OPT, y viceversa.

La principal debilidad de LRU es que para su implementación requiere apoyo en hardware¹⁹ sensiblemente más complejo que FIFO. Una implementación podría ser agregar un contador a cada uno de los marcos, actualizarlo siempre al hacer una referenciar a dicha página, y elegir como víctima a la página con un menor conteo. Este mecanismo tiene la desventaja de que, en presencia de una gran cantidad de páginas, tiene que recorrerlas a todas para buscar a la más envejecida.

Otro mecanismo es emplear una lista doblemente ligada con dos métodos de acceso: Lista y stack. Cada vez que se haga referencia a una página, se mueve a la cabeza del stack, y cada vez que se busque a una página víctima, se selecciona a aquella que esté en el extremo *inferior* del stack. Este mecanismo hace un poco más cara la actualización (pueden requerirse hasta seis modificaciones), pero encuentra a la página víctima en tiempo constante.

Se ha demostrado que LRU y OPT están libres de la anomalía de Belady, dado que, para n marcos, las páginas que estarían en memoria son un subconjunto estricto de las que estarían con $n + 1$ marcos.

Más frecuentemente utilizada (MFU) / Menos frecuentemente utilizada (LFU)

Estos dos algoritmos se basan en mantener un contador, como lo hace LRU, pero en vez de medir el tiempo, miden la *cantidad* de referencias que se han hecho a cada página.

MFU parte de la lógica que, si una página fue empleada muchas veces, probablemente vuelva a ser empleada muchas veces más; LFU parte de que una página que ha sido empleada pocas veces es probablemente una página recién cargada, y va a ser empleada en el futuro cercano.

Estos dos algoritmos son tan caros de implementar como LRU, y su rendimiento respecto a OPT no es tan cercana, por lo cual casi no son empleados.

Aproximaciones a LRU

Dada la complejidad que presenta la implementación de LRU en hardware, algunos sistemas implementan una *aproximación* a éste.

Bit de referencia Esta es una aproximación bastante común. Consiste en que todas las entradas de la tabla de páginas tengan un bit adicional, al que llamaremos *de referencia* o *de acceso*. Al iniciar la ejecución, todos los bits de referencia están apagados (0). Cada vez que se referencia a un marco, su bit de referencia se enciende (esto en general lo realiza el hardware).

El sistema operativo invoca periódicamente a que se apaguen nuevamente todos los bits de referencia. En caso de presentarse un fallo de página, se elige por FIFO sobre el subconjunto de marcos que no hayan sido referenciados en el periodo actual (esto es, entre todos aquellos para los cuales el bit de referencia sea 0).

Columna de referencia Una mejoría casi trivial sobre la anterior consiste en agregar *varios* bits de referencia, conformándose como una *columna*: En vez de descartar su valor cada vez que transcurre el periodo determinado, el valor de la columna de referencia es recorrido a la derecha, descartando el bit más bajo. Por ejemplo, con una implementación de 4 bits, un marco que no ha sido empleado en los últimos 4 periodos tendría el valor 0000, mientras que un marco que sí ha sido referenciado los últimos cuatro periodos tendría 1111. Un marco que fue empleado hace cuatro y tres periodos, pero desde entonces no, tendría el 0011.

¹⁹Dada la frecuencia con que se efectúan referencias a memoria, emplear un mecanismo puramente en software para actualizar las entradas de los marcos resultaría inaceptablemente lento

Cuando el sistema tenga que elegir a una nueva página víctima, lo hará de entre el conjunto que tenga un número más bajo.

La parte de mantenimiento de este algoritmo es muy simple; recorrer una serie de bits es una operación muy sencilla. Seleccionar el número más pequeño requiere una pequeña búsqueda, pero sigue resultando mucho más sencillo que LRU.

Segunda oportunidad (o reloj) El algoritmo de la segunda oportunidad trabaja también basado en un bit de referencia y un recorrido tipo FIFO. La diferencia en este caso es que, al igual que hay eventos que *encienden* a este bit (efectuar una referencia al marco), también hay eventos que lo *apagan*:

Se mantiene un apuntador a la *próxima víctima*, y cuando el sistema requiera efectuar un reemplazo, éste verificará si el marco al que apunta tiene el bit de referencia encendido o apagado. En caso de estar apagado, el marco es seleccionado como víctima, pero en caso de estar encendido (indicando que fue utilizado recientemente), se le da una *segunda oportunidad*: El bit de referencia se apaga, el apuntador de víctima potencial avanza una posición, y vuelve a intentarlo.

A este algoritmo se le llama también *de reloj* porque puede implementarse como una lista ligada circular, y el apuntador puede ser visto como una manecilla. La manecilla avanza sobre la lista de marcos buscando uno con el bit de referencia apagado, y apagando a todos a su paso.

En el peor caso, el algoritmo de *segunda oportunidad* degenera en FIFO.

Segunda oportunidad mejorada El bit de referencia puede ampliarse con un *bit de modificación*, dándonos las siguientes combinaciones, en orden de preferencia:

- (0, 0) No ha sido utilizado ni modificado recientemente. Candidato ideal para su reemplazo.
- (0, 1) No ha sido utilizada recientemente, pero está modificada. No es tan buena opción, porque es necesario escribir la página a disco antes de reemplazarla, pero puede ser elegida.
- (1, 0) El marco está *limpio*, pero fue empleado recientemente, por lo que probablemente se vuelva a requerir pronto.
- (1, 1) Empleada recientemente y *sucia* — Sería necesario escribir la página a disco antes de reemplazar, y probablemente vuelva a ser requerida pronto. Hay que evitar reemplazarla.

La lógica para encontrar una página víctima es similar a la *segunda oportunidad*, pero busca reducir el costo de E/S. Esto puede requerir, sin embargo, dar hasta cuatro vueltas (por cada una de las listas) para elegir la página víctima.

Algoritmos con manejo de buffers

Un mecanismo que se emplea con cada vez mayor frecuencia es que el sistema no espere a enfrentarse a la necesidad de reemplazar un marco, sino que proactivamente busque tener siempre espacio vacío en memoria. Para hacerlo, conforme la carga lo permite, el sistema operativo busca las páginas *sucias* más proclives a ser paginadas a disco y va actualizando el disco (y marcándolas nuevamente como *limpias*). De este modo, cuando tenga que traer una página nueva del disco, siempre habrá espacio donde ubicarla sin tener que esperar a que se transfiera una para liberarla.

6.5.4 Asignación de marcos

Abordando el problema prácticamente por el lado opuesto al del reemplazo de páginas, ¿cómo se asignan los marcos existentes a los procesos del sistema? Esto es, ¿qué esquemas se pueden definir para que la asignación inicial (y, de ser posible, en el transcurso de la ejecución) sea adecuada?

Por ejemplo, usando esquema sencillo: Un sistema con 1024KB de memoria, compuesta

de 256 páginas de 4096 bytes cada una, y basado en paginación puramente sobre demanda.

Si el sistema operativo ocupa 248KB, el primer paso será reservar las 62 páginas que éste requiere, y destinar las 194 páginas restantes para los procesos a ejecutar.

Conforme se van lanzando y comienzan a ejecutar los procesos, cada vez que uno de ellos genere un fallo de página, se le irá asignando uno de los marcos disponibles hasta causar que la memoria entera esté ocupada. Claro está, cuando un proceso termine su ejecución, todos los marcos que tenía asignados volverán a la lista de marcos libres.

Una vez que la memoria esté completamente ocupada (esto es, que haya 194 páginas ocupadas por procesos), el siguiente fallo de página invocará a un algoritmo de reemplazo de página, que elegirá una de las 194.²⁰

Este esquema, si bien es simple, al requerir una gran cantidad de fallos de página explícitos puede penalizar el rendimiento del sistema — El esquema puede resultar demasiado flojo, no le vendría mal ser un poco más *ansioso* y asignar, de inicio, un número determinado como mínimo utilizable de marcos.

Mínimo de marcos

Si un proceso tiene asignados muy pocos marcos, su rendimiento indudablemente se verá afectado. Hasta ahora se ha supuesto que cada instrucción puede causar un sólo fallo de página, pero la realidad es más compleja. Cada instrucción del procesador puede, dependiendo de la arquitectura, desencadenar varias solicitudes y potencialmente varios fallos de página.

Todas las arquitecturas proporcionan instrucciones de referencia directa a memoria (instrucciones que permiten especificar una dirección de memoria para leer o escribir) — Esto significa que todas requerirán que, para que un proceso funcione adecuadamente, tenga por lo menos dos marcos asignados: En caso de que se le permitiera solo uno, si la instrucción ubicada en 0x00A2C8 solicita la carga de 0x043F00, esta causaría dos fallos: El primero, cargar al marco la página 0x043, y el segundo, cargar nuevamente la página 0x00A, necesario para leer la siguiente instrucción a ejecutar del programa (0x00A2CC, asumiendo palabras de 32 bits).

Algunas arquitecturas, además, permiten *referencias indirectas a memoria*, esto es, la dirección de carga puede solicitar la dirección *que está referenciada* en 0x043F00. El procesador tendría que recuperar esta dirección, y podría encontrarse con que hace referencia a una dirección en otra página (por ejemplo, 0x010F80). Cada nivel de indirección que se permite aumenta en uno el número de páginas que se deben reservar como mínimo por proceso.

Algunas arquitecturas, particularmente las más antiguas,²¹ permiten que tanto los operandos de algunas instrucciones aritméticas como su resultado sean direcciones de memoria (y no operan estrictamente sobre los registros, como las arquitecturas RISC). En dichas arquitecturas, el mínimo debe también tener este factor en cuenta: Si en una sola instrucción es posible sumar dos direcciones de memoria y guardar el resultado en una adicional, el mínimo a reservar es de cuatro marcos: Uno para el flujo del programa, uno para el primer operando, uno para el segundo operando, y uno para el resultado.

Esquemas de asignación

Ahora, una vez establecido el número mínimo de marcos por proceso, ¿cómo determinar el nivel deseable?

²⁰En realidad, dentro de la memoria del sistema operativo, al igual que la de cualquier otro proceso, hay regiones que deben mantenerse residentes y regiones que pueden paginarse. Se puede, simplificando, omitir por ahora esa complicación y asumir que el sistema operativo completo se mantendrá siempre en memoria

²¹Aquellas diseñadas antes de que la velocidad del procesador se distanciara tanto del tiempo de acceso a memoria

Partiendo de que el rendimiento de un proceso será mejor entre menos fallos de paginación cause, se podría intentar otorgar a cada proceso el total de marcos que solicita — Pero esto tendría como resultado disminuir el grado de multiprogramación, y por tanto, reducir el uso efectivo total del procesador.

Otra alternativa es la *asignación igualitaria*: Se divide el total de espacio en memoria física entre todos los procesos en ejecución, en partes iguales. Esto es, volviendo a la computadora hipotética que se presentó al inicio de esta sección, si existen 4 procesos que requieren ser ejecutados, de los 194 marcos disponibles, el sistema asignará 48 marcos (192KB) a dos de los procesos y 49 (196KB) a los otros dos (es imposible asignar fracciones de marcos). De este modo, el espacio será compartido por igual.

La asignación igualitaria resulta ser un esquema deficiente para casi todas las distribuciones de procesos: Bajo este esquema, si P_1 es un gestor de bases de datos que puede estar empleando 2048KB (512 páginas) de memoria virtual (a pesar de que el sistema tiene sólo 1MB de memoria física) y P_2 es un lector de texto que está empleando un usuario, requiriendo apenas 112KB (28 páginas), con lo cual incluso dejaría algunos de sus marcos sin utilizar.

Un segundo esquema, que resuelve mejor esta situación, es la *asignación proporcional*: Dar a cada proceso una porción del espacio de memoria física proporcional a su uso de memoria virtual.

De tal suerte, si además de los procesos anteriores se tiene a P_3 empleando 560KB (140 páginas) y a P_4 con 320KB (80 páginas) de memoria virtual, el uso total de memoria virtual sería de $V_T = 512 + 28 + 140 + 80 = 760$ páginas, esto es, el sistema tendría comprometido a través de la memoria virtual un sobreuso cercano a 4:1 sobre la memoria física²².

Cada proceso recibirá entonces $F_P = \frac{V_P}{V_T} \times m$, donde F_P indica el espacio de memoria física que el proceso recibirá, V_P la cantidad de memoria virtual que está empleando, y m la cantidad total de marcos de memoria disponibles. De este modo, P_1 recibirá 130 marcos, P_2 7, P_3 35 y P_4 20, proporcionalmente a su uso de memoria virtual.

Cabe apuntar que este mecanismo debe observar ciertos parámetros mínimos: Por un lado, si el mínimo de marcos definido para esta arquitectura es de 4, por más que entrara en ejecución un proceso de 32KB (8 páginas) o aumentara al doble el grado de multiprocesamiento, ningún proceso debe tener asignado menos del mínimo definido.

La asignación proporcional también debe cuidar no sobre-asignar recursos a un proceso *obeso*: P_1 es ya mucho más grande que todos los procesos del sistema. En caso de que esta creciera mucho más, por ejemplo, si multiplicara por 4 su uso de memoria virtual, esto llevaría a que se *castigara* desproporcionadamente a todos los demás procesos del sistema.

Por otro lado, este esquema ignora por completo las prioridades que hoy en día manejan todos los sistemas operativos; si se quisiera considerar, podría incluirse como factor la prioridad, multiplicando junto con V_P .

El esquema de asignación proporcional sufre, sin embargo, cuando cambia el nivel de multiprogramación — Esto es, cuando se inicia un nuevo proceso o finaliza un proceso en ejecución, deben recalcularse los espacios en memoria física asignados a cada uno de los procesos restantes. Si finaliza un proceso, el problema es menor, pues sólo se asignan los marcos y puede esperarse a que se vayan poblando por paginación sobre demanda, pero si inicia uno nuevo, es necesario reducir de golpe la asignación de todos los demás procesos hasta abrir suficiente espacio para que quepa.

Por último, el esquema de la asignación proporcional también tiende a desperdiciar recursos: Si bien hay procesos que mantienen un patrón estable de actividad a lo largo de su ejecución, muchos otros pueden tener períodos de mucho menor requisitos. Por ejemplo, un proceso servidor de documentos pasa la mayor parte de su tiempo simplemente

²²Ya que de los 1024KB, o 256 páginas, que tiene el sistema descrito, descontando los 248KB, o 62 páginas, que ocupa el sistema operativo, quedan 194 páginas disponibles para los procesos

esperando solicitudes, y podría reducirse a un uso mínimo de memoria física, sin embargo, al solicitársele un documento, se le deberían poder asignar más marcos (para trabajar en una *ráfaga*) hasta que termine con su tarea. En la sección 6.5.5 se retomará este tema.

Ámbitos del algoritmo de reemplazo de páginas

Para atender a los problemas no resueltos que se describieron en la sección anterior, se puede discutir el ámbito en que operará el algoritmo de reemplazo de páginas.

Reemplazo local Mantener tan estable como sea posible el cálculo hecho por el esquema de asignación empleado. Esto significa que cuando se presente un fallo de página, las páginas que serán consideradas para su intercambio serán únicamente aquellas pertenecientes *al mismo proceso* que el que causó el fallo.

Un proceso tiene asignado su espacio de memoria física, y se mantendrá estable mientras el sistema operativo no tome alguna decisión por cambiarlo.

Reemplazo global Los algoritmos de asignación determinan el espacio asignado a los procesos al ser inicializados, e influyen a los algoritmos de reemplazo (por ejemplo, dando mayor peso para ser elegidas como páginas víctima a aquellas que pertenezcan a un proceso que excede de su asignación en memoria física).

Los algoritmos de reemplazo de páginas operan sobre el espacio completo de memoria, y la asignación física de cada proceso puede variar según el estado del sistema momento a momento.

Reemplazo global con prioridad Es un esquema mixto, en el que un proceso puede *sobrepasar* su límite siempre que le robe espacio en memoria física exclusivamente a procesos de prioridad inferior a él. Esto es consistente con el comportamiento de los algoritmos planificadores, que siempre dan preferencia a un proceso de mayor prioridad por sobre de uno de prioridad más baja.

El reemplazo local es más rígido y no permite aprovechar para mejorar el rendimiento los períodos de inactividad de algunos de los procesos. En contraposición, los esquemas basados en reemplazo global pueden llevar a rendimiento inconsistente: Dado que la asignación de memoria física sale del control de cada proceso, puede que la misma sección de código presente tiempos de ejecución muy distintos si porciones importantes de su memoria fueron paginadas a disco.

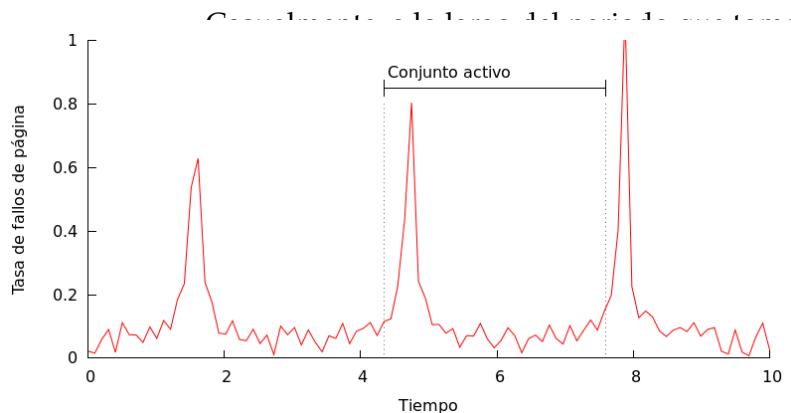
6.5.5 Hiperpaginación

Es un fenómeno que se puede presentar por varias razones: cuando (bajo un esquema de reemplazo local) un proceso tiene asignadas pocas páginas para llevar a cabo su trabajo, y genera fallos de página con tal frecuencia que le imposibilita realizar trabajo real. Bajo un esquema de reemplazo global, cuando hay demasiados procesos en ejecución en el sistema y los constantes fallos y reemplazos hacen imposible a todos los procesos involucrados avanzar, también se presenta hiperpaginación²³.

Hay varios escenarios que pueden desencadenar la hiperpaginación, y su impacto es tan claro e identificable que prácticamente cualquier usuario de cómputo lo sabrá reconocer. A continuación se presentará un escenario ejemplo en que las malas decisiones del sistema operativo pueden conducirlo a este estado.

Suponga un sistema que está con una carga media normal, con un esquema de reemplazo global de marcos. Se lanza un nuevo proceso, que como parte de su inicialización requiere poblar diversas estructuras a lo largo de su espacio de memoria virtual. Para hacerlo, lanza una serie de fallos de página, a las que el sistema operativo responde reemplazando a varios marcos pertenecientes a otros procesos.

²³Una traducción literal del término *thrashing*, empleado en inglés para designar a este fenómeno, resulta más gráfica: *Paliza*



parte del disco duro. El sistema completo avanza cada vez más lento.

6.23: Los picos y valles en la cantidad de fallos de página de un proceso definen a su *conjunto activo*. Los síntomas de la hiperpaginación son muy claros, y no son difíciles de detectar. ¿Qué estrategia puede emplear el sistema operativo una vez que se da cuenta que se presentó esta situación?

Una salida sería reducir el nivel de multiprogramación — Si la paginación se presentó debido a que los requisitos de memoria de los procesos actualmente en ejecución no pueden ser satisfechos con la memoria física disponible, el sistema puede seleccionar a uno (o más) de los procesos y suspenderlos por completo hasta que el sistema vuelva a un estado normal. Podría seleccionarse, por ejemplo, al proceso con menor prioridad, al que esté causando más cantidad de fallos, o al que esté ocupando más memoria.

Modelando el *conjunto activo*

Un pico en la cantidad de fallos de página no necesariamente significa que se va a presentar una situación de hiperpaginación — Muchas veces indica que el proceso cambió su atención de un conjunto de páginas a otro, o dicho de otro modo, que cambió el *conjunto activo* del proceso — Y resulta natural que, al cambiar el conjunto activo, el proceso accese de golpe una serie de páginas que no había referenciado en cierto tiempo.

El *conjunto activo* es, pues, la aproximación más clara a la *localidad de referencia* de un proceso dado: El conjunto de páginas sobre los que está iterando en un momento dado.

Idealmente, para evitar los problemas relacionados con la hiperpaginación, el sistema debe asignar a cada proceso suficientes páginas como para que mantenga en memoria física su conjunto activo — Y si no es posible hacerlo, el proceso es un buen candidato para ser suspendido. Sin embargo, detectar con suficiente claridad como para efectuar este diagnóstico cuál es el conjunto activo es una tarea muy compleja, que típicamente implica rastrear y verificar del orden de los últimos miles a decenas de miles de accesos a memoria.

6.6 Consideraciones de seguridad

Para una cobertura a mayor profundidad del material presentado en esta sección, se sugiere estudiar los siguientes textos:

- [Smashing The Stack For Fun And Profit](#) (Aleph One, revista Phrack, 1996)
- [The Tao of Buffer Overflows](#) (Enrique Sánchez, inédito, pero disponible en Web)

6.6.1 Desbordamientos de buffer (*buffer overflows*)

Una de las funciones principales de los sistemas operativos en la que se ha insistido a lo largo del libro es la de implementar protección entre los procesos pertenecientes a diferentes usuarios, o ejecutándose con distinto nivel de privilegios. Y si bien el enfoque general que se ha propuesto es el de analizar por separado subsistema por subsistema, al

esta inicialización (que puede parecer s de veces más lento que la memoria) memoria que acaban de ser enviados a

lización del procesador decrece, puede antenimiento. Se lanzan estos procesos, isponible para cada uno de los procesos

s de paginación, algunas veces contratar NOOP (esto es, no tiene trabajo que asa en espera de un nuevo marco por parte del disco duro. El sistema completo avanza cada vez más lento.

hablar de administración de memoria es necesario mencionar también las implicaciones de seguridad que del presente tema se pueden desprender.

En las computadoras de arquitectura von Neumann, todo dato a ser procesado (sean instrucciones o datos) debe pasar por la memoria, por el *almacenamiento primario*. Sólo desde ahí puede el procesador leer la información directamente.

A lo largo del presente capítulo se ha mencionado que la MMU incluye ya desde el hardware el concepto de *permisos*, separando claramente las regiones de memoria donde se ubica el código del programa (y son, por tanto, ejecutables y de sólo lectura) de aquellas donde se encuentran los datos (de lectura y escritura). Esto, sin embargo, no los pone a salvo de los *desbordamientos de buffer* (*buffer overflows*), errores de programación (típicamente, la falta de verificación de límites) que pueden convertirse en vulnerabilidades²⁴.

La pila de llamadas (stack)

Recordando lo mencionado en la sección 6.1.4, en que se presentó el espacio en memoria de un proceso, es conveniente profundizar un poco más acerca de cómo está estructurada la *pila de llamadas (stack)*.

El *stack* es el mecanismo que brinda un sentido local a la representación del código estructurado. Está dividido en *marcos de activación* (sin relación con el concepto de marcos empleado al hablar de memoria virtual); durante el periodo en que es el marco *activo* (esto es, cuando no se ha transferido el control a ninguna otra función), está delimitado por dos valores, almacenados en registros:

Apuntador a la pila (*Stack pointer, SP*) Apunta al *final actual* (dirección inferior) de la pila.

En arquitecturas x86, emplea el registro `ESP`; cuando se pide al procesador que actúe sobre el *stack* (con las operaciones `pushl` o `popl`), lo hace sobre este registro

Apuntador del marco (*Frame pointer, FP*, o *Base local, LB*) Apunta al *inicio* del marco actual, o lo que es lo mismo, al final del marco anterior. En arquitecturas x86, emplea el registro `EBP`.

A cada función a que va entrando la ejecución del proceso, se va creando un *marco de activación* en el stack, que incluye:

- Los argumentos recibidos por la función
- La dirección de retorno al código que la invocó
- Las variables locales creadas en la función

Con esto en mente, es posible analizar la traducción de una llamada a función en C a su equivalente en ensamblador, y en segundo término ver el marco del stack resultante:

```
void func(int a, int b, int c) {
    char buffer1[5];
    char buffer2[10];
}

void main() {
    func(1, 2, 3);
}
```

Y lo que el código resultante en ensamblador efectúa es:

1. El procesador *empuja* (`pushl`) los tres argumentos al stack (`ESP`). La notación empleada (`$1, $2, $3`) indica que el número indicado se expresa de forma literal. Cada uno de estos tres valores restará 4 bytes (el tamaño de un valor entero en x86-32) a `ESP`.

²⁴Citando a Theo de Raadt, autor principal del sistema operativo OpenBSD, todo error es una vulnerabilidad esperando a ser descubierta

2. En ensamblador, los nombres asignados a las variables y funciones no significan nada. La llamada `call` no es lo que se entendería como una llamada a función en un lenguaje de alto nivel — Lo que hace el procesador es *empujar* al stack la dirección de la siguiente instrucción, y cargar a éste la dirección en el fuente donde está la etiqueta de la función (esto es, transferir la ejecución hacia allá).
3. Lo primero que hace la función al ser invocada es asegurarse de saber a dónde volver: *empuja* al stack el viejo apuntador al marco (`EBP`), y lo reemplaza (`movl`) por el actual. A esta ubicación se le llama `SFP` (*Saved Frame Pointer, apuntador al marco grabado*)
4. Por último, con `subl`, resta el espacio necesario para alojar las variables locales, `buffer1` y `buffer2`. Notarán que, si bien éstas son de 5 y 10 bytes, está recorriendo 20 bytes — Esto porque, en la arquitectura x86-32, los accesos a memoria deben estar *alineados a 32 bits*.

```
; main
    pushl $3
    pushl $2
    pushl $1
    call func

func:
    pushl %ebp
    movl %esp, %ebp
    subl $20, %esp
```

C y las funciones de manejo de cadenas

El lenguaje de programación C fue creado con el propósito de ser tan simple como sea posible, manteniéndose tan cerca del hardware como se pudiera, para que pudiera ser empleado como un lenguaje de programación para un sistema operativo portable. Y si bien en 1970 era visto como un lenguaje bastante de alto nivel, hoy en día es lo más bajo nivel en que programa la mayor parte de los desarrolladores del mundo.

C no tiene soporte nativo para *cadenas* de caracteres. El soporte es provisto a través de *familias* de funciones en la biblioteca estándar del lenguaje, que están siempre disponibles en cualquier implementación estándar de C. Las familias principales son `strcat`, `strcpy`, `printf` y `gets`. Estas funciones trabajan con cadenas que siguen la siguiente estructura:

- Son arreglos de 1 o más caracteres (`char`, 8 bits)
- Deben terminar con el byte de terminación NUL (`\0`)

El problema con estas funciones es que sólo algunas de las funciones derivadas implementan verificaciones de límites, y algunas son incluso capaces de crear cadenas ilegales (que no concluyan con el terminador `\0`).

El problema aparece cuando el programador no tiene el cuidado necesario al trabajar con datos de los cuales no tiene *certeza*. Esto se demuestra con el siguiente código vulnerable:

```
#include <stdio.h>
int main(int argc, char **argv) {
    char buffer[256];
    if(argc > 1) strcpy(buffer, argv[1]);
    printf("Escribiste %s\n", buffer);
    return 0;
}
```

El problema con este código reside en el `strcpy(buffer, argv[1])` — Dado que el código es recibido del usuario, no se tiene la *certeza* de que el argumento que recibe el programa por línea de comandos (empleando `argv[1]`) quepa en el arreglo `buffer[256]`. Esto es, si se ejecuta el programa ejemplo con una cadena de 120 caracteres:

```
$ ./ejemplo1 'perl -e 'print "A" x 120''
Escribiste: AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAA
$
```

La ejecución resulta exitosa. Sin embargo, si se ejecuta el programa con un parámetro demasiado largo para el arreglo:

```
$ ./ejemplo1 'perl -e 'print "A" x 500'
Escribiste: AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAAAAAAAAAAA
Segmentation fault
$
```

De una falla a un ataque

En el ejemplo recién presentado, parecería que el sistema *atrapó* al error exitosamente y detuvo la ejecución — Pero no lo hizo: El `Segmentation fault` no fue generado al sobreescribir el buffer ni al intentar procesarlo, sino después de terminar de hacerlo: Al llegar la ejecución del código al `return 0`.

Para volver de una función a quien la invocó, incluso si dicha función es `main()`, lo que hace `return` es restaurar el viejo SFP y hacer que el apuntador a siguiente dirección *salte* a la dirección que tiene en RET. Sin embargo, como se observa en el esquema, RET fue sobreescrito por la dirección `0x41414141` (AAAA). Dado que esa dirección no forma parte del espacio del proceso actual, se lanza una excepción por violación de segmento, y el proceso es terminado.

Ahora, lo expuesto anteriormente implica que el código *está demostrado vulnerable*, pero se ha *explorado* aún. El siguiente paso es, conociendo el acomodo exacto de la memoria, sobreescribir únicamente lo necesario para alterar el flujo del programa — Esto es, sobreescribir RET con una dirección válida. Para esto, es necesario conocer la longitud desde el inicio del buffer hasta donde terminan RET y SFP, en este caso particular, 264 bytes (256 del buffer mas 4 de RET mas 4 de SFP).

Citando al texto de Enrique Sánchez,

¿Por qué ocurre un desbordamiento de stack? Imagina un vaso y una botella de cerveza. ¿Qué ocurre si sirves la botella completa en el vaso? Se va a derramar. Imagina que tu variable es el vaso, y la entrada del usuario es la cerveza. Puede ocurrir que el usuario sirva tanto líquido como el que cabe en el vaso, pero puede también seguir sirviendo hasta que se derrame. La cerveza se derramaría en todas direcciones, pero la memoria no crece de esa manera, es sólo un arreglo bidimensional, y sólo crece en una dirección.

Ahora, ¿qué más pasa cuando desbordas a un contenedor? El líquido sobrante va a mojar la botana, los papeles, la mesa, etc. En el caso de los papeles, destruirá cualquier cosa que hubieras apuntado (como el teléfono que acabas de anotar de esa linda chica). Cuando tu variable se desborde, ¿qué va a sobrescribir? Al EBP, al EIP, y lo que les siga, dependiendo de la función, y si es la última función, las variables de ambiente. Puede que el programa aborte y tu shell resulte inutilizado a causa de las variables sobreescritas.

Hay dos técnicas principales: *Saltar* a un punto determinado del programa, y *saltar* hacia dentro del stack.

Un ejemplo de la primera técnica se muestra a continuación. Si el atacante está intentando burlar la siguiente validación simple de nombre de usuario y contraseña,

```
if (valid_user(usr, pass)) {
    /* (...) */
} else {
    printf("Error!\n");
    exit 1;
}
```

Y detecta que `valid_user()` es susceptible a un desbordamiento, le bastaría con incrementar en 4 la dirección de retorno. La conversión de este `if` a ensamblador es, primero, saltar hacia la etiqueta `valid_user`, y (empleando al valor que ésta regrese en `%EBX`) ir a la siguiente instrucción, o saltar a la etiqueta `FAIL`. Esto puede hacerse con la instrucción `BNE $0, %EBX, FAIL` (*Branch if Not Equal, Saltar si no es igual*, que recibe como argumentos dos valores a ser comparados, en este caso el registro `%EBX` y el número 0, y la etiqueta destino, `FAIL`). Cambiar la dirección destino significa burlar la verificación.

Por otro lado, el atacante podría usar la segunda técnica para lograr que el sistema haga algo más complejo — Por ejemplo, que ejecute código arbitrario que él proporcione. Para esto, el ataque más frecuente es saltar *hacia dentro del stack*.

Para hacerlo, si en vez de proporcionar simplemente una cadena suficientemente grande para sobrepasar el buffer se *inyecta* una cadena con código ejecutable válido, y sobreescribiera la dirección de retorno con la dirección de su código *dentro del buffer*, tendría 256 bytes de espacio para especificar código arbitrario. Este código típicamente se llama *shellcode*, pues se emplea para obtener un *shell* (un intérprete de comandos) que ejecuta con los privilegios del proceso explotado.

Mecanismos de mitigación

Claro está, el mundo no se queda quieto. Una vez que estos mecanismos de ataque se dieron a conocer, comenzó un fuerte trabajo para crear mecanismos de mitigación de daños.

La principal y más importante medida es crear una cultura de programadores conscientes y prácticas seguras. Esto cruza necesariamente el no emplear funciones que no hagan verificación de límites. La desventaja de esto es que hace falta cambiar al *factor humano*, lo cual resulta prácticamente imposible de lograr con suficiente profundidad²⁵. Muchos desarrolladores esgrimen argumentos en contra de estas prácticas, como la pérdida de rendimiento que estas funciones requieren, y muchos otros sencillamente nunca se dieron por enterados de la necesidad de programar correctamente.

Por esto, se han ido creando diversos mecanismos automatizados de protección ante los desbordamientos de buffer. Ninguno de estos mecanismos es *perfecto*, pero sí ayudan a reducir los riesgos ante los atacantes menos persistentes o habilidosos.

Secciones de datos no ejecutables

En secciones anteriores se describió la protección que puede imponer la MMU por regiones, evitando la modificación de código ejecutable.

En la arquitectura x86, dominante en el mercado de computadoras personales desde hace muchos años, esta característica existía en varios procesadores basados en el modelo de segmentación de memoria, pero desapareció al cambiarse el modelo predominante por uno de memoria plana paginada, y fue hasta alrededor del 2001 en que fue introducida de vuelta, bajo los nombres *bit NX* (*Never eXecute*, Nunca ejecutar) o *bit XD* (*eXecute Disable*, Deshabilitar ejecución), como una característica particular de las extensiones PAE.

Empleando este mecanismo, la MMU puede evitar la ejecución de código en el área de stack, lo cual anula la posibilidad de *saltar al stack*. Esta protección desafortunadamente no es muy efectiva: Una vez que tiene acceso a un buffer vulnerable, el atacante puede *saltar a libc*, esto es, por ejemplo, proporcionar como parámetro el nombre de un programa a ejecutar, e indicar como retorno la dirección de la función `system` o `execve` de la `libc`.

Las secciones de datos no ejecutables son, pues, un obstáculo ante un atacante, aunque no representan una dificultad mucho mayor.

Aleatorización del espacio de direcciones

Otra técnica es que, en tiempo de carga y a cada ejecución, el proceso reciba diferentes direcciones base para sus diferentes áreas. Esto hace más difícil para el atacante poder indicar a qué dirección destino se debe saltar.

Un atacante puede emplear varias técnicas para ayudarse a *adivinar* detalles acerca del acomodo en memoria de un proceso, y, con un buffer suficientemente grande, es común ver *cadenas de NOP*, esto es, una extensión grande de operaciones nulas, seguidas del *shellcode*, para aumentar las probabilidades de que el control se transfiera a un punto útil.

Empleo de canarios

Se llama *canario* a un valor aleatorio de protección²⁵, insertado entre los buffers y la dirección de retorno, que es verificado antes de regresar de una función. Si se presentó un desbordamiento de buffer, el valor del *canario* será reemplazado por basura, y el sistema podrá detener la ejecución del proceso comprometido antes de que brinde privilegios elevados al atacante.

Un atacante tiene dos mecanismos ante un sistema que requiere del canario: Uno es el atacar no directamente a la función en cuestión, sino que al *manejador de señales* que es notificado de la anomalía, y otro es, ya que se tiene acceso a la memoria del proceso, *averiguar el valor del canario*. Esto requiere ataques bastante más sofisticados que los vistos en esta sección, pero definitivamente ya no fuera del alcance de los atacantes.

6.6.2 Ligado estático y dinámico de bibliotecas

Las *bibliotecas de código* (o simplemente *bibliotecas*) implementan el código de una serie de funcionalidades generales, que pueden ser usadas en diferentes programas y contextos. Un ejemplo clásico sería la biblioteca estándar de C, la cual ofrece funciones básicas de entrada / salida, manejo de cadenas, entre otras.

A medida que el software crece en complejidad, los programadores recurren a la *reutilización de código* para aprovechar la implementación de la funcionalidad que ofrecen las distintas bibliotecas. De esta forma, evitan “reinventar la rueda”, y se concentran en la funcionalidad específica del software que están construyendo.

²⁵El ejemplo más claro de este problema es la función `gets`, la cual sigue siendo enseñada y usada en los cursos básicos de programación en C.

²⁶Este uso proviene de la costumbre antigua de los mineros de tener un canario en una jaula en las minas. Como el canario es muy sensible ante la falta de oxígeno, si el canario moría servía como indicador a los mineros de que debían abandonar la mina de inmediato, antes de correr la misma suerte.

El concepto de *ligado* se refiere al proceso mediante el cual, se toma el *codigo objeto* de un programa junto con el código de las bibliotecas que éste usa para crear un archivo ejecutable. De forma general existen dos tipos de ligado, que se explican a continuación.

El *ligado estático* consiste en tomar el código de una biblioteca e integrarlo al código del programa para generar el archivo ejecutable. Esto implica que *cada programa* tiene su propia copia del código de la biblioteca, lo cual puede causar un desperdicio de memoria y disco si existen muchos programas que usan la misma versión.

Por su parte, en el *ligado dinámico* el código de las bibliotecas no se copia dentro de la imagen ejecutable del programa, pero requiere establecer algún mecanismo para informar que el programa necesita un código externo. Esto se puede implementar de diferentes formas. Por ejemplo, se puede incluir un fragmento de código dentro del programa que usa la biblioteca denominado *stub*, el cual en tiempo de ejecución solicita que se cargue la biblioteca requerida. Otra estrategia que se puede utilizar consiste en incluir algunas indicaciones que le permiten al sistema operativo, en el momento de crear el proceso, ubicar las bibliotecas que este requerirá para su ejecución. En cualquier caso, el ligado dinámico busca que las bibliotecas sólo sean cargadas cuando sean requeridas.

La figura 6.4, presentada al principio de este capítulo, ilustra el momento en que ocurre cada uno de estos ligados: El ligado estático es realizado por el *editor de ligado*, uniendo en un sólo *módulo cargable* al programa compilado (*módulo objeto*) con las bibliotecas (*otros objetos*); el ligado dinámico es realizado parcialmente en tiempo de carga (para las *bibliotecas del sistema*) y parcialmente en tiempo de ejecución (para las *bibliotecas de carga dinámica*).²⁷

Las bibliotecas y la seguridad

El ligado dinámico puede traer consigo una serie de problemas, entre los cuales se destacan el manejo de versiones de las bibliotecas y potenciales vulnerabilidades. El primer problema es conocido, en ambientes Windows, como el *infierro de las DLL*. Este *infierro* se puede causar de muchas formas. Por ejemplo, si al instalar un nuevo programa, se instala también una versión incompatible de una biblioteca que es usada por otros programas. Esto causa que los demás programas no se puedan ejecutar — Y lo que es más, hace que la depuración del fallo sea muy difícil. Por otro lado, si no se tienen los controles suficientes, al desinstalar un programa se puede borrar una biblioteca compartida, lo cual puede llevar a que otros programas dejen de funcionar.

El *infierro de las DLL* puede ser prevenido mediante estrategias como el *versionamiento* de las biblioteca de ligado dinámico (esto es, hacer que cada componente de las bibliotecas lleve la versión que implementa o *nivel de compatibilidad* que implementa),²⁸ y mediante el uso de scripts de instalación o *gestores de dependencias* que verifican si existe en el sistema una versión compatible de la biblioteca. Teniendo esta información, la biblioteca en cuestión se instalará únicamente en caso necesario.

El ligado dinámico puede presentar problemas o vulnerabilidades debido a que el programa usa un código proporcionado por terceros, y *confía* en que la biblioteca funciona tal como se espera sin incluir código malicioso. Por tal razón, desde el punto de vista teórico bastaría que un atacante instale su propia versión de una biblioteca para que pueda tener el control de los programas que la usan e incluso del mismo sistema operativo.²⁹ En el caso de bibliotecas ligadas *estáticamente*, dado que estas forman ya parte del programa,

²⁷Refiérase al libro *Linkers and Loaders* (ligadores y cargadores) de John R. Levine (1999) para mayores detalles respecto a este proceso.

²⁸Este nivel de compatibilidad incluye no sólo a la {interfaz de aplicación al programador} (API, definida en las secciones 3.7 y 3.7.1), sino también la {interfaz de aplicación binaria} (ABI), esto es, no sólo la información de el nombre de las funciones que expone y los tipos de argumentos que reciben, sino también la ubicación en memoria de su definición en un archivo ya compilado.

²⁹Esto se puede lograr, por ejemplo, alterando la configuración del entorno en la cual el sistema busca las bibliotecas.

un atacante tendría que modificar al archivo objeto mismo del programa para alterar las bibliotecas.

Así las cosas, más allá de la economía de espacio en memoria, ¿cómo se explica que sea tanto más popular el ligado dinámico en los sistemas operativos modernos?

Parte muy importante de la respuesta es la *mantenibilidad*: Si es encontrado un fallo en una biblioteca de carga dinámica, basta con que los desarrolladores lo corrijan una vez (cuidando, claro, de mantener la compatibilidad binaria) y reemplazar a dicha biblioteca en disco *una sola vez*. Todos los programas que liguen dinámicamente con esta biblioteca tendrán disponible de inmediato la versión actualizada. En muchos sistemas operativos, el *gestor de paquetes* puede detectar cuáles de los procesos en ejecución emplean a determinada biblioteca dinámica, y reiniciarlos de forma transparente al administrador.

En contraste, de estar el fallo en una biblioteca de ligado estático, el código afectado estaría incluido como parte de *cada uno de los programas* ligados con ella. Como consecuencia, para corregir este defecto, cada uno de los programas afectados tendría que ser recompilado (o, por lo menos, *relijado*) antes de poderse beneficiar de las correcciones.

Y si bien este proceso resulta manual y tedioso para un administrador de sistemas con acceso a las fuentes de los programas que tiene instalados, resulta mucho más oneroso aún para quienes emplean software *propietario* (En la sección A.1.3 se aborda con mayor detenimiento lo que significa el software propietario en contraposición al software libre).

Cabe mencionar que el comportamiento del sistema ante la actualización de una biblioteca descrita ilustra una de las diferencias semánticas entre sistemas Windows y sistemas Unix que serán abordadas en el capítulo 7: Mientras un sistema Unix permite la eliminación de un archivo *que está siendo utilizado*, Windows no lo permite. Esto explica por qué las actualizaciones de bibliotecas en sistemas Windows se aplican *durante el proceso de apagado*: Mientras haya procesos que tienen abierta una biblioteca, ésta no puede ser reemplazada. Caso contrario en sistemas Unix, en que el archivo puede ser reemplazado, pero mientras no sean reiniciados los procesos en cuestión, éstos seguirán ejecutando la versión de la biblioteca con el error.

6.7 Otros recursos

- *The Grumpy Editor goes 64-bit*

<https://lwn.net/Articles/79036/>

Jonathan Corbet (2004); Linux Weekly News. Experiencia del editor de Linux Weekly News al migrar a una arquitectura de 64 bits en 2004. Lo más interesante del artículo son los comentarios, ilustran buena parte de los pros y contras de una migración de 32 a 64 bits.

- *Using Valgrind to debug Xen Toolstacks*

<http://www.hellion.org.uk/blog/posts/using-valgrind-on-xen-toolstacks/>

Ian J. Campbell (2013). Presenta un ejemplo de uso de la herramienta *Valgrind*, para encontrar problemas en la asignación, uso y liberación de memoria en un programa en C.

- *Process memory usage*

<http://troysunix.blogspot.mx/2011/07/process-memory-usage.html>

ejemplos de pmap en diferentes Unixes

- *Página de manual de pmap en NetBSD*

<http://www.daemon-systems.org/man/pmap.1.html>

Más allá de simplemente mostrar la operación de una herramienta del sistema en Unix, esta página de manual ilustra claramente la estructura de la organización de la memoria.

- *Linkers and Loaders*

<http://www.iecc.com/linker/>

John R. Levine (1999). Libro de libre descarga y redistribución, dedicado a la tarea de los editores de ligado y el proceso de carga de los programas.

- *An anomaly in space-time characteristics of certain programs running in a paging machine*
<http://dl.acm.org/citation.cfm?doid=363011.363155>
Belady, Nelson, Shedler (1969); Communications of the ACM
- *Understanding the Linux Virtual Memory Manager*
http://ptgmedia.pearsoncmg.com/images/0131453483/downloads/gorman_book.pdf

Mel Gorman (2004). Libro de libre descarga y redistribución, parte de la colección *Bruce Perens' Open Source Series*. Aborda a profundidad los mecanismos y algoritmos relativos a la memoria empleados por el sistema operativo Linux. Entra en detalles técnicos a profundidad, presentándolos poco a poco, por lo que no resulta demasiado complejo de leer. El primer tercio del libro describe los mecanismos, y los dos tercios restantes siguen el código comentado que los implementa.

- *The art of picking Intel registers*
<http://www.swansontec.com/sregisters.html>
William Swanson (2003).
- *The Tao of Buffer Overflows*
http://sistop.gwolf.org/biblio/The_Tao_of_Buffer_Overflows_-_Enrique_Sanchez.pdf
Enrique Sánchez (trabajo en proceso, facilitado por el autor)
- *Smashing the Stack for fun and profit*
<http://www.phrack.org/issues.html?issue=49&id=14>
Aleph One (1996): Uno de los primeros artículos publicados acerca de los *buffers desbordados*
- *Attacking the Windows 7/8 Address Space Randomization*
<http://kingcope.wordpress.com/2013/01/24/>
Kingcopes (2013): Explica cómo puede burlarse la protección ALSR (aleatorización de direcciones) en Windows 7 y 8, logrando una dirección predecible de memoria hacia la cual saltar.
- *An overview of Non-Uniform Memory Access*
<https://dl.acm.org/citation.cfm?doid=2500468.2500477>
Cristoph Lameter (2013); *Communications of the ACM*
- *Anatomy of a killer bug: How just 5 characters can murder iPhone, Mac apps*
http://www.theregister.co.uk/2013/09/04/unicode_of_death_crash/
Chris Williams (2013); *The Register*. A fines de agosto del 2013, se descubrió una *cadena Unicode de la muerte*, que causa que cualquier programa que intente desplegarla en pantalla en la línea de productos Apple se caiga. Este artículo relata el proceso de averiguar, a partir del *reporte de falla* generado por el sistema y analizando el contenido de la memoria que éste reporta, cómo se puede encontrar un *desbordamiento de entero*.

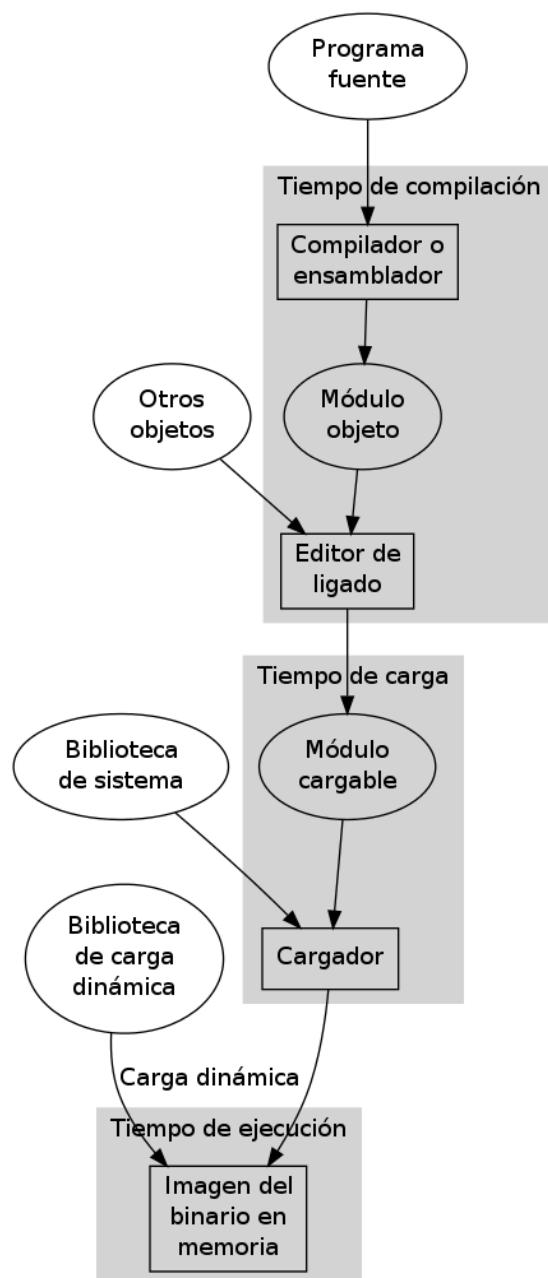


Figura 6.4: Proceso de compilación y carga de un programa, indicando el tipo de resolución de direcciones (Silberschatz, p.281)

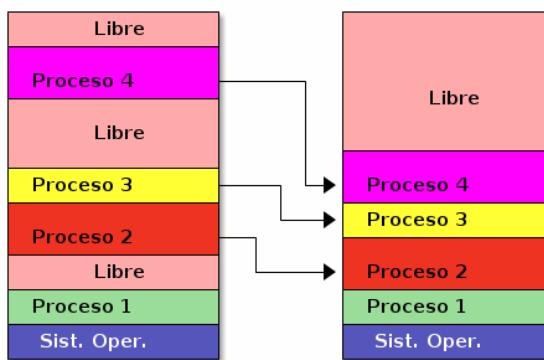


Figura 6.5: Compactación de la memoria de procesos en ejecución

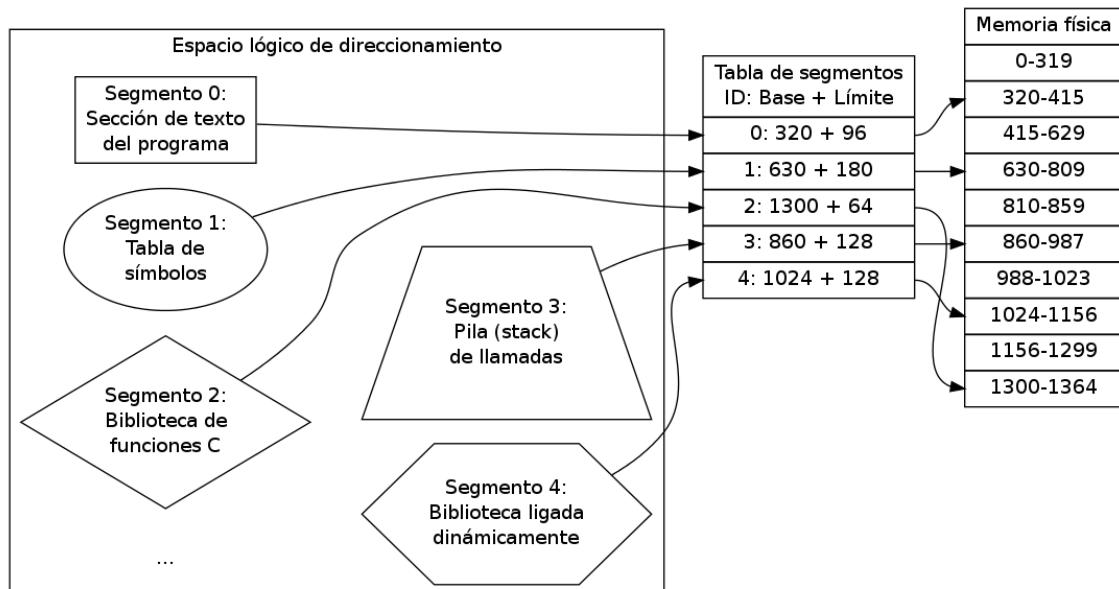


Figura 6.6: Ejemplo de segmentación

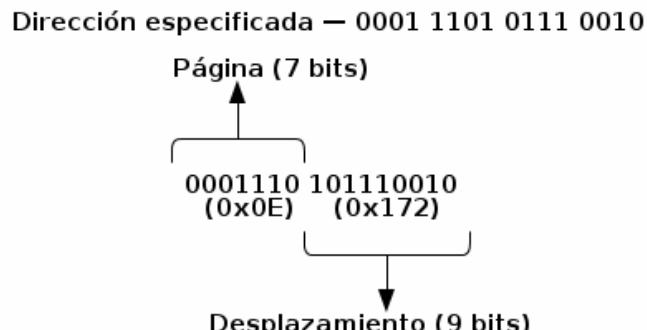


Figura 6.7: Página y desplazamiento, en un esquema de direccionamiento de 16 bits y páginas de 512 bytes

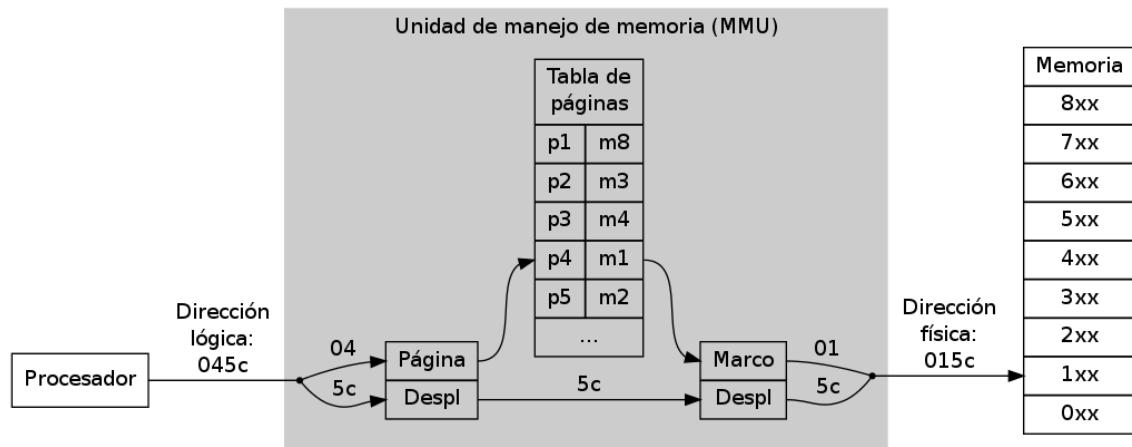


Figura 6.8: Esquema del proceso de pagenación, ilustrando el rol de la MMU

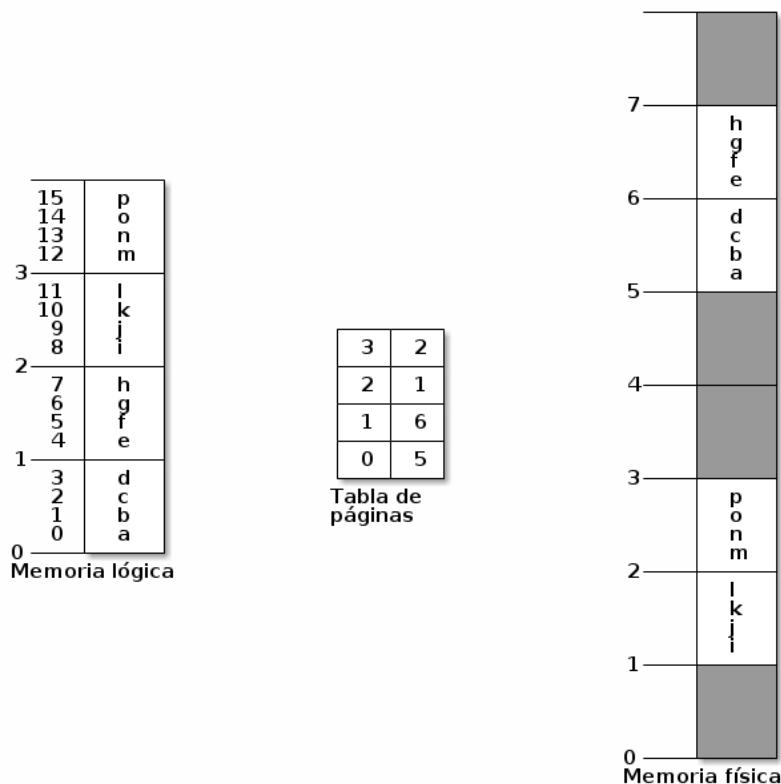


Figura 6.9: Ejemplo (minúsculo) de pagenación, con un espacio de direccionamiento de 32 bytes y páginas de 4 bytes

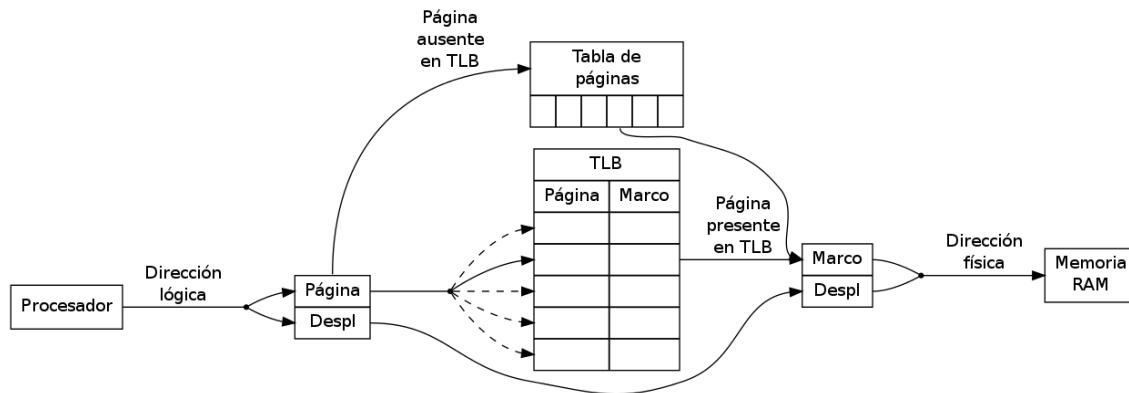


Figura 6.10: Esquema de paginación empleando un *buffer de traducción adelantada* (TLB)

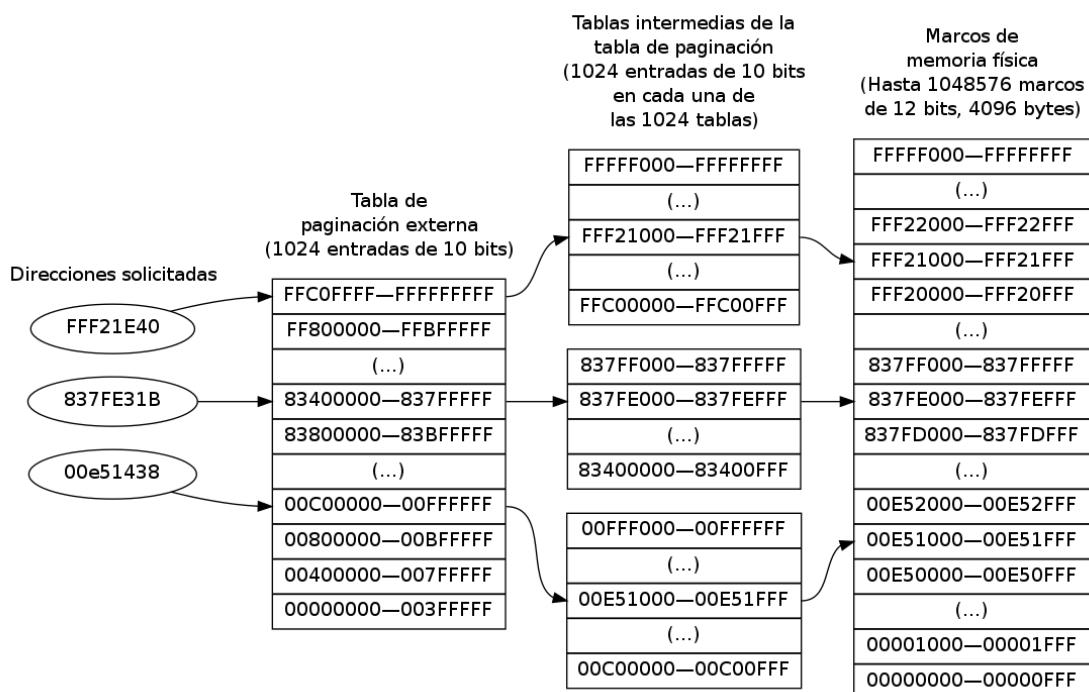


Figura 6.11: Paginación en dos niveles: Una tabla externa de 10 bits, tablas intermedias de 10 bits, y marcos de 12 bits (esquema común para procesadores de 32 bits)



Figura 6.12: Uso de memoria compartida: Tres procesos comparten la memoria ocupada por el texto del programa (azul), difieren sólo en los datos.

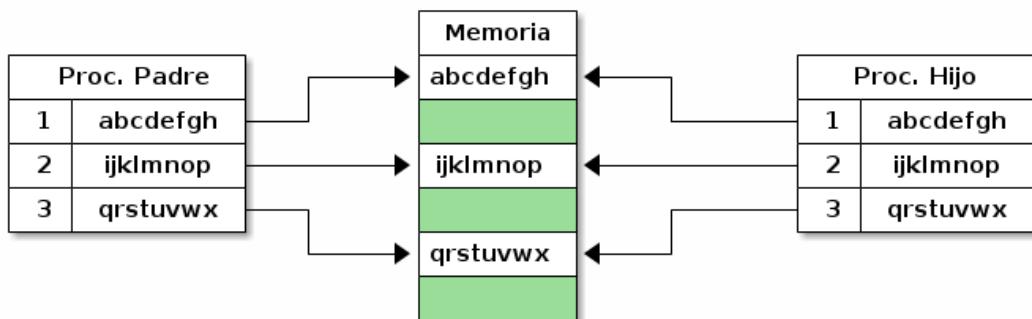


Figura 6.13: Memoria de dos procesos inmediatamente después de la creación del proceso hijo por `fork()`

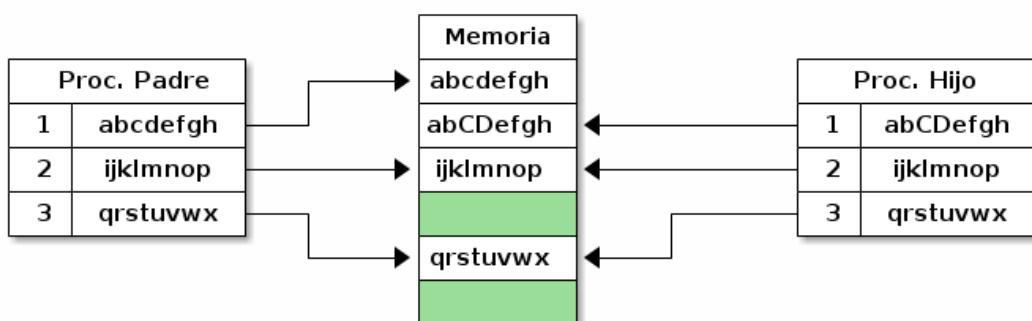


Figura 6.14: Cuando el proceso hijo modifica información en la primer página de su memoria, se crea como una página nueva.

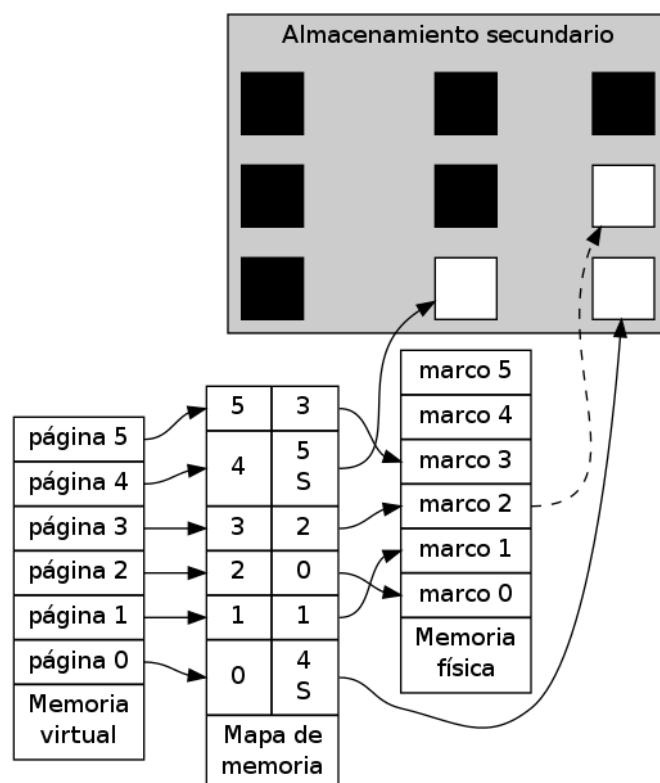


Figura 6.15: Esquema general de la memoria, incorporando espacio en almacenamiento secundario, representando la memoria virtual

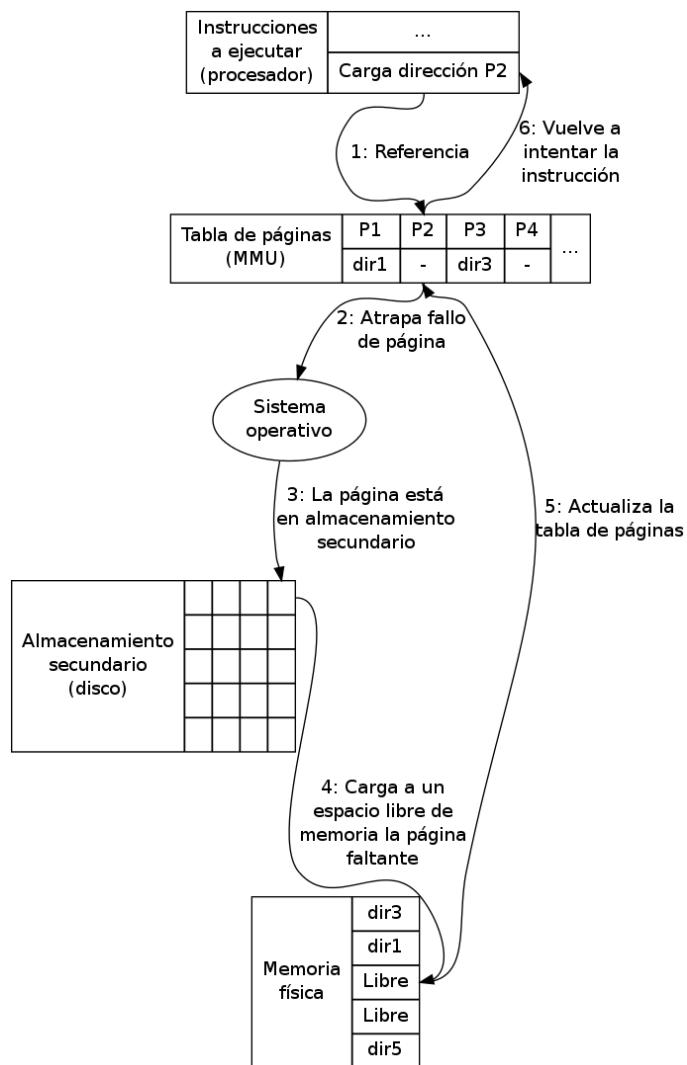


Figura 6.16: Pasos que atraviesa la respuesta a un fallo de página

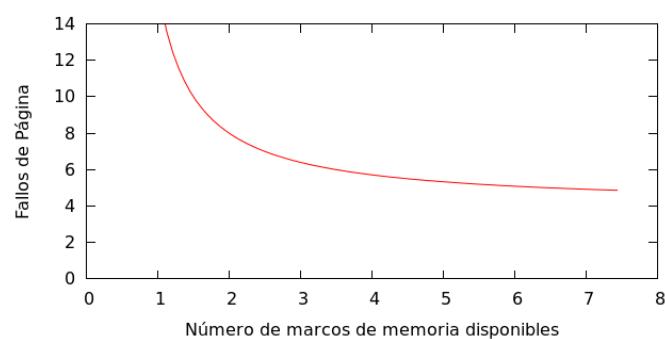


Figura 6.17: Relación ideal entre el número de marcos y fallos de página



Figura 6.18: Comportamiento del algoritmo FIFO que exhibe la anomalía de Belady al pasar de 3 a 4 marcos. La cadena de referencia que genera este comportamiento es 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5 (Belady, 1969)

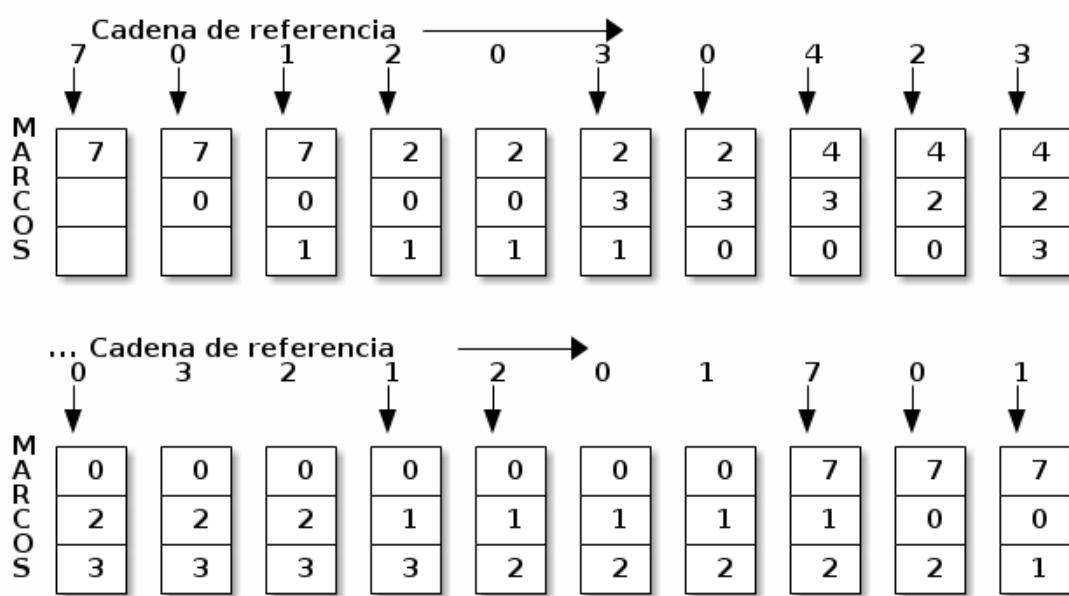


Figura 6.19: Algoritmo FIFO de reemplazo de páginas

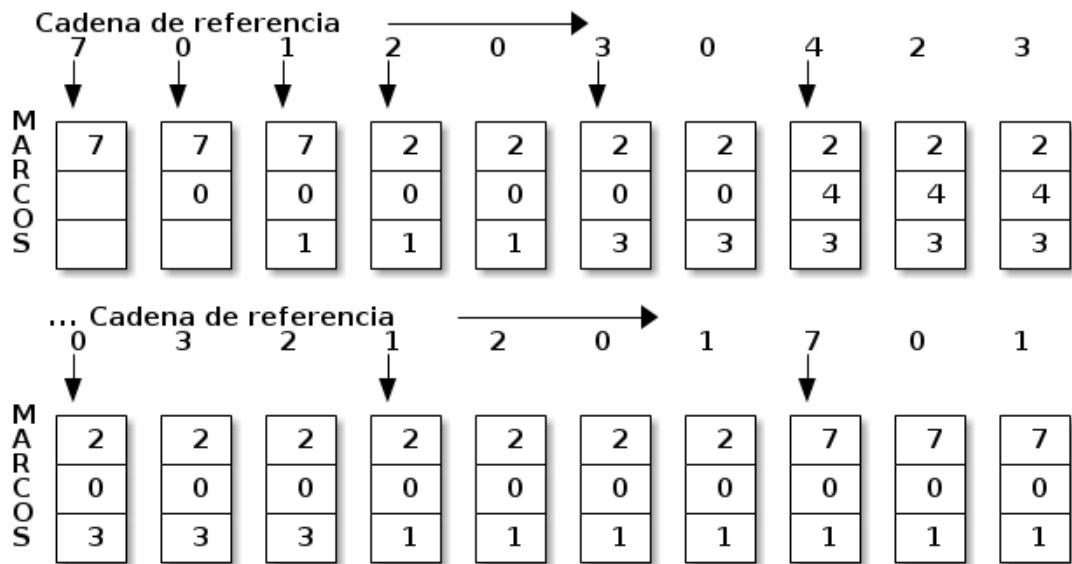


Figura 6.20: Algoritmo óptimo de reemplazo de páginas (OPT)

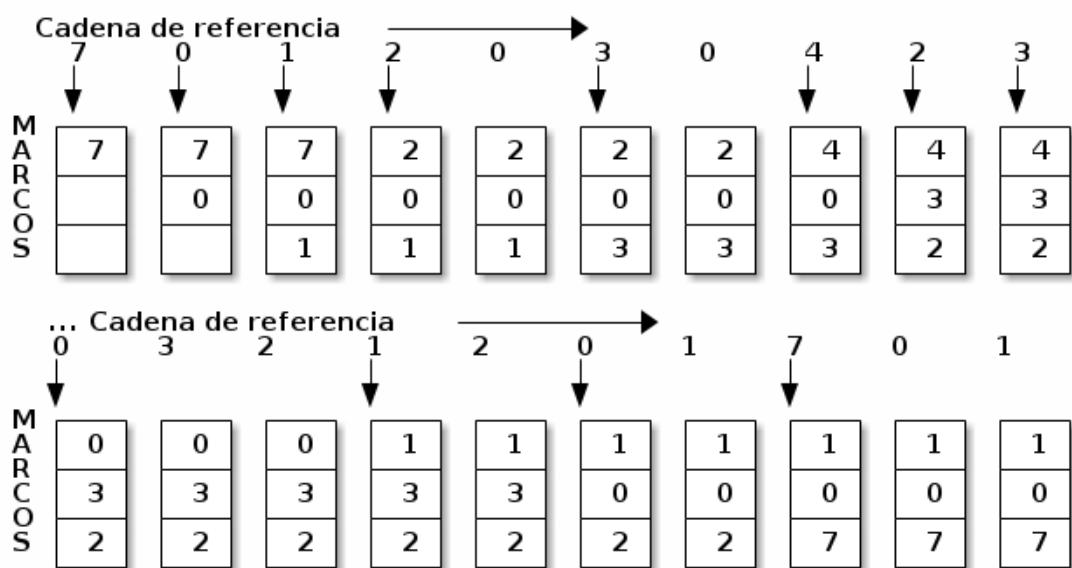


Figura 6.21: Algoritmo reemplazo de páginas menos recientemente utilizadas (LRU)

7 — Organización de archivos

7.1 Introducción

De los roles que cumple el sistema operativo, probablemente el que más consciente tengan en general sus usuarios es el de la gestión del espacio de almacenamiento, esto es, la organización de la información en un *sistema de archivos*. Al día de hoy, todos los usuarios de equipo de cómputo dan por sentado y comprenden a grandes rasgos la organización del espacio de almacenamiento en un *directorio jerárquico*, con unidades de almacenamiento llamadas *archivos*, de diferentes tipos según su función. En el presente capítulo se revisará la semántica que compone a este modelo, para en el capítulo 8 continuar con los detalles de la gestión del espacio físico donde éstos están alojados.

La abstracción que hoy se conoce como *sistemas de archivos* es una de las que más tiempo ha vivido y se ha mantenido a lo largo de la historia de la computación, sobreviviendo a lo largo de prácticamente todas las generaciones de sistemas operativos. Sin embargo, para poder analizar cómo es que el sistema operativo representa la información en el dispositivo físico, el presente capítulo inicia discutiendo cómo es que esta información es comprendida por los niveles más altos — Por los programas en espacio de usuario.

La información *cruda* tiene que pasar una serie de transformaciones. Yendo de niveles superiores a niveles más bajos, un programa estructura sus datos en *archivos*, siguiendo el *formato* que resulte mas pertinente al tipo de información a representar. Un conjunto de archivos hoy en día es típicamente representado en una estructura de *directorios*,¹.

Cada dispositivo empleado para almacenar archivos tiene un directorio. Cuando un sistema opera con más de un dispositivo físico, existen principalmente dos mecanismos para integrar a dichos dispositivos en un *sistema de archivos virtual*,² brindando al usuario una interfaz uniforme. Por último, los archivos son una estructura meramente lógica; deben ser convertidos para ser representados en un *dispositivo de bloques* como los diversos tipos de unidades –aunque esta nomenclatura es a veces incorrecta– como *discos*. Este último paso será abordado en el capítulo 8.

Del diagrama presentado en la figura 7.1, toca al objeto de nuestro estudio –el sistema operativo– recibir del espacio de usuario las llamadas al sistema que presentan la interfaz de archivos y directorios, integrar el sistema de archivos virtual, y traducir la información resultante a un sistema de archivos.

Cabe mencionar que varias de las capas aquí presentadas podrían perfectamente ser subdivididas, analizadas por separado, e incluso tratarse de forma completamente modular — De hecho, este es precisamente el modo en que se implementan de forma transparente características hoy en día tan comunes como sistemas de archivos en red, o compresión y cifrado de la información. Una referencia más detallada acerca de ventajas, desventajas, técnicas y mecanismos de la división y comunicación entre capas puede ubicarse en el artículo de [Heidemann y Popek \(1994\)](#).

¹Existen otros mecanismos para su organización, pero estos no están tan ampliamente difundidos

²Esto será abordado en la sección 7.3.3

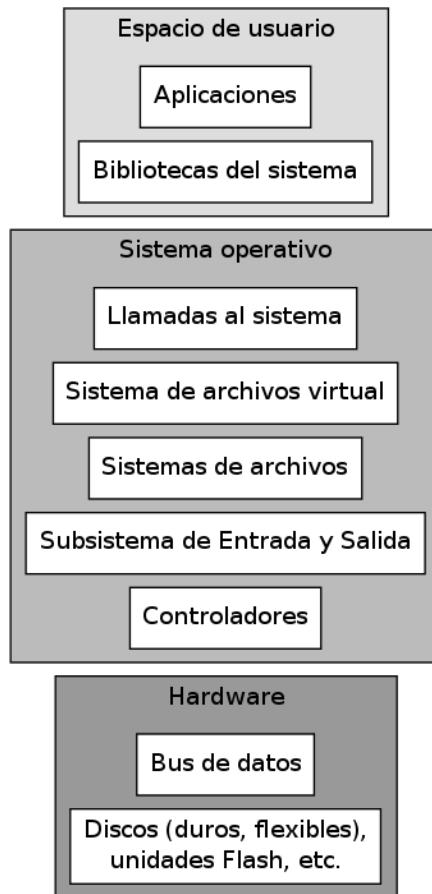


Figura 7.1: Capas de abstracción para implementar los sistemas de archivos

7.2 Concepto de archivo

En primer término, un archivo es un *tipo de datos abstracto* — Esto es, podría verse como una estructura que exclusivamente permite la manipulación por medio de una interfaz *orientada a objetos*: Los procesos en el sistema sólo pueden tener acceso a los archivos por medio de la interfaz ofrecida por el sistema operativo.³ La siguiente sección describe las principales operaciones provistas por esta interfaz.

Para el usuario, los archivos son la *unidad lógica mínima* al hablar de almacenamiento: Todo el almacenamiento *persistente* (que sobrevive en el tiempo, sea a reinicios del sistema, a pérdida de corriente o a otras circunstancias en el transcurso normal de ejecución) en el sistema al que tiene acceso, se efectúa dentro de archivos; el espacio libre en los diferentes dispositivos no tiene mayor existencia fuera de saber que está *potencialmente* disponible.

Dentro de cada *volúmen* (cada medio de almacenamiento), los archivos disponibles conforman a un *directorio*, y son típicamente identificados por un *nombre* o una *ruta*. Más adelante se presentarán de las diferentes construcciones semánticas que pueden conformar a los directorios.

7.2.1 Operaciones con archivos

Cada sistema operativo definirá la interfaz de archivos acorde con su semántica, pero en líneas generales, las operaciones que siempre estarán disponibles con un archivo son:

³Como se verá en la sección 8.1.3, esto no es *necesariamente* así, sin embargo, el uso de los dispositivos *en crudo* es muy bajo. Este capítulo está enfocado exclusivamente al uso estructurado en sistemas de archivos.

Borrar Elimina al archivo del directorio y, de ser procedente, libera el espacio del dispositivo

Abrir Solicita al sistema operativo verificar si el archivo existe o puede ser creado (dependiendo del modo requerido) y se cuenta con el acceso para el *modo de acceso* al archivo indicado y si el medio lo soporta (por ejemplo, a pesar de contar con todos los permisos necesarios, el sistema operativo no debe permitir abrir para escritura un archivo en un CD-ROM u otro medio de sólo lectura). En C, esto se hace con la función `fopen()`.

Al abrir un archivo, el sistema operativo asigna un *descriptor de archivo* que identifica la relación entre el proceso y el archivo en cuestión; estos serán definidos propiamente en la sección 7.2.2.

Todas las operaciones descritas a continuación operan sobre el descriptor de archivo, no con su nombre o ruta.

Cerrar Indica al sistema que el *proceso en cuestión* terminó de trabajar con el archivo; el sistema entonces debe escribir los buffers a disco y eliminar la entrada que representa a esta combinación archivo-proceso de las tablas activas, invalidando al *descriptor de archivo*. En C, para cerrar un descriptor de archivo se usa `fclose()`.

Dado que todas las operaciones se realizan a través del descriptor de archivo, si un proceso cierra un archivo y requiere seguir utilizándolo, tendrá que abrirlo de nuevo para obtener un nuevo descriptor.

Leer Si se solicita al sistema la lectura leer de un archivo hacia determinado buffer, éste copia el siguiente *pedazo* de información a éste. Este *pedazo* podría ser una línea o un bloque de longitud definida, dependiendo del modo en que se solicite la lectura. El sistema mantiene un apuntador a la última posición leída, para poder *continuar* con la lectura de forma secuencial.

La función que implementa la lectura en C es `fread()`. Cabe mencionar que `fread()` entrega el *número de caracteres* especificado; para trabajar con líneas de texto hace falta trabajar con bibliotecas que implementen esta funcionalidad, como `readline`.

Escribir Teniendo un archivo abierto, guarda información en él. Puede ser que escriba desde su primer posición (*truncando* al archivo, esto es, borrando toda la información que pudiera ya tener), o *agregando* al archivo, esto es, iniciando con el apuntador de escritura al final del mismo. La función C para escribir a un descriptor de archivo es `fwrite()`.

Reposicionar Tanto la lectura como la escritura se hacen siguiendo a un *apuntador*, que indica cuál fue la última posición del archivo a la que accesó el proceso actual. Al reposicionar el apuntador, se puede *saltar* a otro punto del archivo. La función que reposiciona el apuntador dentro de un descriptor de archivos es `fseek()`.

Hay varias otras operaciones comunes que pueden implementarse con llamadas compuestas a estas operaciones (por ejemplo, *copiar* un archivo puede implementarse como *crear* un archivo nuevo en modo de escritura, abrir en modo de lectura al archivo fuente, e ir *leyendo* de éste y *escribiendo* al nuevo hasta llegar al fin de archivo fuente).

Las operaciones aquí presentadas no son todas las operaciones existentes; dependiendo del sistema operativo, habrá algunas adicionales; estas se presentan como una base general común a los principales sistemas operativos.

Vale la pena mencionar que esta semántica para el manejo de archivos presenta a cada archivo como si fuera una *unidad de cinta*, dentro de la cual la cabeza lectora/escritora simulada puede avanzar o retroceder.

7.2.2 Tablas de archivos abiertos

Tanto el sistema operativo como cada uno de los procesos mantienen normalmente *tablas de archivos abiertos*. Éstas mantienen información acerca de todos los archivos actualmente abiertos, presentándolos hacia el proceso por medio de un *descriptor de archivo*; una vez que un archivo fue abierto, las operaciones que se realizan dentro de éste no son empleando su nombre, sino que su descriptor de archivo.

En un sistema operativo multitareas, más de un proceso podría abrir el mismo archivo a la vez; lo que cada uno de ellos pueda hacer, y cómo esto impacte a lo que vean los demás procesos, depende de la semántica que implemente el sistema; un ejemplo de las diferentes semánticas posibles es el descrito en la sección 7.2.3.

Ahora, ¿por qué estas tablas son mantenidas tanto por el sistema operativo como por cada uno de los procesos? ¿No lleva esto a una situación de información redundante?

La respuesta es que la información que cada uno debe manejar es distinta. El sistema operativo necesita:

Conteo de usuarios del archivo Cuando se solicita, por ejemplo, *desmontar* una partición (por ejemplo, para expulsar una unidad removible) o eliminar un archivo, el sistema debe poder determinar cuándo es momento de declarar la solicitud como *efectuada*. Si algún proceso tiene abierto a un archivo, y particularmente si tiene cambios pendientes de guardar, el sistema debe hacer lo posible por evitar que el archivo *desaparezca* de su visión.

Modos de acceso Aunque un usuario tenga permisos de acceso a determinado recurso, el sistema puede determinar negarlo si llevaría a una inconsistencia. Por ejemplo, si dos procesos abren un mismo archivo en modo de escritura, es probable que los cambios que realice uno sobreesciban a los que haga el otro.

Ubicación en disco El sistema mantiene esta información para evitar que cada proceso tenga que consultar las tablas en disco para encontrar al archivo, o cada uno de sus fragmentos.

Información de bloqueo En caso de que los modos de acceso del archivo requieran protección mutua, puede implementarse por medio de un bloqueo.

Por otro lado, el proceso necesita:

Descriptor de archivo Relación entre el nombre del archivo abierto y el identificador numérico que maneja internamente el proceso. Un archivo abierto por varios procesos tendrá descriptores de archivo distintos en cada uno de ellos.

A nivel implementación, el descriptor de archivo otorgado por el sistema a un proceso es simplemente un número entero, que podría entenderse como *el n-ésimo archivo empleado por el proceso*.⁴

Permisos Los modos válidos de acceso para un archivo. Esto no necesariamente es igual a los permisos que tiene el archivo en cuestión en disco, sino que el *subconjunto* de dichos permisos bajo los cuales está operando para este proceso en particular — Si un archivo fue abierto en modo de sólo lectura, por ejemplo, este campo sólo permitirá la lectura.

7.2.3 Acceso concurrente: Bloqueo de archivos

Dado que los archivos pueden emplearse como mecanismo de comunicación entre procesos que no guarden relación entre sí, incluso a lo largo del tiempo, y para emplear un archivo basta indicar su nombre o ruta, los sistemas operativos multitarea implementan

⁴No sólo los archivos reciben descriptores de archivo. Por ejemplo, en todos los principales sistemas operativos, los descriptores 0, 1 y 2 están relacionados a *flujos de datos*: respectivamente, la entrada estándar (STDIN), la salida estándar (STDOUT) y el error estándar (STDERR); si el usuario no lo indica de otro modo, la terminal desde donde fue ejecutado el proceso.

mecanismos de bloqueo para evitar que varios procesos intentando emplear de forma concurrente a un archivo se corrompan mutuamente.

Algunos sistemas operativos permiten establecer bloqueos sobre determinadas regiones de los archivos, aunque la semántica más común es operar sobre el archivo entero.

En general, la nomenclatura que se sigue para los bloqueos es:

Compartido (*Shared lock*) Podría verse como equivalente a un bloqueo (o *candado*) para realizar lectura — Varios procesos pueden adquirir al mismo tiempo un bloqueo de lectura, e indica que todos los que posean dicho *candado* tienen la expectativa de que el archivo no sufrirá modificaciones.

Exclusivo (*Exclusive lock*) Un bloqueo o *candado* exclusivo puede ser adquirido por un sólo proceso, e indica que realizará operaciones que modifiquen al archivo (o, si la semántica del sistema operativo permite expresarlo, a la *porción* del archivo que indica).

Respecto al *mecanismo* de bloqueo, hay también dos tipos, dependiendo de qué tan explícito tiene que ser su manejo:

Mandatorio u obligatorio (*Mandatory locking*) Una vez que un proceso adquiere un *candado* obligatorio, el sistema operativo se encargará de imponer las restricciones correspondientes de acceso a todos los demás procesos, independientemente de si éstos fueron programados para considerar la existencia de dicho bloqueo o no.

Consultivo o asesor (*Advisory locking*) Este tipo de bloqueos es manejado cooperativamente entre los procesos involucrados, y depende del programador de *cada uno* de los programas en cuestión el solicitar y respetar dicho bloqueo.

Haciendo un paralelo con los mecanismos presentados en el capítulo 4, los mecanismos que emplean mutexes, semáforos o variables de condición serían *consultivos*, y únicamente los que emplean monitores (en que la única manera de llegar a la información es a través del mecanismo que la protege) serían *mandatorios*.

No todos los sistemas operativos implementan las cuatro posibles combinaciones (compartido mandatorio, o compartido compulsivo, exclusivo mandatorio y exclusivo consultivo). Como regla general, en los sistemas Windows se maneja un esquema de bloqueo obligatorio, y en sistemas Unix es de bloqueo consultivo.⁵

Cabe mencionar que el manejo de bloqueos con archivos requiere del mismo cuidado que el de bloqueo por recursos cubierto en la sección 4.4: Dos procesos intentando adquirir un *candado* exclusivo sobre dos archivos pueden caer en un bloqueo mutuo tal como ocurre con cualquier otro recurso.

7.2.4 Tipos de archivo

Si los archivos son la *unidad lógica mínima* con la que se puede guardar información en almacenamiento secundario, naturalmente sigue que existen archivos de diferentes tipos: Cada archivo podría ser un documento de texto, un binario ejecutable, un archivo de audio o video, o un larguísimo etcetera, e intentar emplear a un archivo como uno de un tipo distinto puede resultar desde una frustración al usuario porque el programa no responde como éste quiere, hasta en pérdidas económicas.⁶

Hay tres estrategias principales para que el sistema operativo reconozca al tipo de un archivo:

⁵Esto explica por qué en Windows es tan común que el sistema mismo rechace hacer determinada operación porque {el archivo está abierto por otro programa} (bloqueo mandatorio compartido), mientras que en Unix esta responsabilidad recae en cada uno de los programas de aplicación

⁶Por ejemplo, imprimir un archivo binario resulta en una gran cantidad de hojas inútiles, particularmente tomando en cuenta que hay caracteres de control como el ASCII 12 (avance de forma, *form feed*), que llevan a las impresoras que operan en modo texto a iniciar una nueva página; llevar a un usuario a correr un archivo ejecutable *disfrazado* de un documento inocuo, como se verá a continuación, fue un importante vector de infección de muchos virus.

Extensión En los sistemas CP/M de los 1970, el nombre de cada archivo se dividía en dos porciones, empleando como elemento separador al punto: El nombre del archivo y su extensión. El sistema mantenía una lista de extensiones conocidas, para las cuales sabría cómo actuar, y este diseño se extendería a las aplicaciones, que sólo abrirían a aquellos archivos cuyas extensiones supieran manejar.

Esta estrategia fue heredada por VMS y MS-DOS, de donde la adoptó Windows; ya en el contexto de un entorno gráfico, Windows agrega, más allá de las extensiones directamente ejecutables, la relación de extensiones con los programas capaces de trabajar con ellas, permitiendo invocar a un programa con sólo dar “doble click” en un archivo.

Como nota, este esquema de asociación de tipo de archivo permite ocultar las extensiones toda vez que ya no requieren ser del conocimiento del usuario, sino que son gestionadas por el sistema operativo, abre una vía de ataque automatizado que se popularizó en su momento: El envío de correos con extensiones engañosas duplicadas — Esto es, el programa maligno (un *programa troyano*) se envía a todos los contactos del usuario infectado, presentándose por ejemplo como una imagen, con el nombre `inocente.png.exe`. Por el esquema de ocultamiento mencionado, éste se presenta al usuario como `inocente.png`, pero al abrirlo, el sistema operativo lo reconoce como un ejecutable, y lo ejecuta en vez de abrirlo en un visor de imágenes.

Números mágicos La alternativa que emplean los sistemas Unix es, como siempre, simple y elegante, aunque indudablemente presenta eventuales lagunas: El sistema mantiene una lista compilada de las *huellas digitales* de los principales formatos que debe manejar,⁷ para reconocer el contenido de un archivo basado en sus primeros bytes. Casi todos los formatos de archivo incluyen lo necesario para que se lleve a cabo este reconocimiento, y cuando no es posible hacerlo, se intenta por medio de ciertas reglas *heurísticas*. Por ejemplo, todos los archivos de imagen en *formato de intercambio gráfico* (GIF) inician con la cadena `GIF87a` o `GIF89a`, dependiendo de la versión; los archivos del lenguaje de descripción de páginas PostScript inician con la cadena `% !`, el *Formato de Documentos Portátiles* (PDF) con `%PDF`, etcétera. Un documento en formatos definidos alrededor de XML inicia con `<!DOCTYPE`. Algunos de estos formatos no están *anclados* al inicio, sino que en un punto específico del primer bloque.

Un caso especial de números mágicos es el llamado *hashbang* (`#!`). Esto indica a un sistema Unix que el archivo en cuestión (típicamente un archivo de texto, incluyendo código fuente en algún lenguaje de *script*) debe tratarse como un ejecutable, y empleando como *intérprete* al comando indicado inmediatamente después del *hashbang*. Es por esto que se pueden ejecutar directamente, por ejemplo, los archivos que inician con `#!/usr/bin/bash`: El sistema operativo invoca al programa `/usr/bin/bash`, y le especifica como argumento al archivo en cuestión.

Metadatos externos Los sistemas de archivos empleado por las Apple Macintosh desde 1984 separan en dos *divisiones (forks)* la información de un archivo: Los datos que propiamente constituyen al archivo en cuestión son la *división de datos (data fork)*, y la información *acerca del archivo* se guardan en una estructura independiente llamada *división de recursos (resource fork)*.

Esta idea resultó fundamental para varias de las características *amigables al usuario* que presentó Macintosh desde su introducción — Particularmente, para presentar un entorno gráfico que respondiera ágilmente, sin tener que buscar los datos base de una aplicación dentro de un archivo de mucho mayor tamaño. La *división de recursos* cabe en pocos sectores de disco, y si se toma en cuenta que las primeras Macintosh

⁷Una de las ventajas de este esquema es que cada administrador de sistema puede ampliar la lista con las huellas digitales que requiera localmente

funcionaban únicamente con discos flexibles, el tiempo invertido en leer una lista de iconos podría ser demasiada.

La división de recursos puede contener todo tipo de información; los programas ejecutables son los que le dan un mayor uso, dado que incluyen desde los aspectos gráficos (ícono a mostrar para el archivo, ubicación de la ventana a ser abierta, etc.) hasta aspectos funcionales, como la traducción de sus cadenas al lenguaje particular del sistema en que está instalado. Esta división permite una gran flexibilidad, dado que no es necesario tener acceso al fuente del programa para crear traducciones y temas.

En el tema particular que concierne a esta sección, la división de recursos incluye un campo llamado *creador*, que indica cuál programa fue el que generó al archivo. Si el usuario solicita ejecutar un archivo de datos, el sistema operativo lanzaría al programa *creador*, indicándole que abra al archivo en cuestión.

Las versiones actuales de MacOS ya no emplean esta técnica, sino que una llamada *appDirectory*, para propósitos de esta discusión, la técnica base es la misma.

7.2.5 Estructura de los archivos y métodos de acceso

La razón principal de la existencia del sistema de archivos son *los archivos*. Un archivo almacena información de *algún tipo*, estructurado o no estructurado.

La mayor parte de los sistemas operativos maneja únicamente archivos *sin estructura* — Cada aplicación es responsable de preparar la información de forma congruente, y la responsabilidad del sistema operativo es únicamente entregarlo como un conjunto de bytes. Históricamente, hubo sistemas operativos, como IBM CICS (1968), IBM MVS (1974) o DEC VMS (1977), que administraban ciertos tipos de datos en un formato básico de *base de datos*.

El hecho de que el sistema operativo no imponga estructura a un archivo no significa, claro está, que la aplicación que lo genera no lo haga. La razón por la que los sistemas creados en los últimos 30 años no han implementado este esquema de base de datos es que le *resta* flexibilidad al sistema: El que una aplicación tuviera que ceñirse a los tipos de datos y alineación de campos del sistema operativo impedía su adecuación, y el que la funcionalidad de un archivo tipo base de datos dependiera de la versión del sistema operativo creaba un *acoplamiento* demasiado rígido entre el sistema operativo y las aplicaciones.

Esta práctica ha ido cediendo terreno a dejar esta responsabilidad en manos de procesos independientes en espacio de usuario (como sería un gestor de bases de datos tradicional), o de bibliotecas que ofrezcan la funcionalidad de manejo de archivos estructurados (como en el caso de *SQLite*, empleado tanto por herramientas de adquisición de datos de bajo nivel como *systemtap* como por herramientas tan de escritorio como el gestor de fotografías *shotwell* o el navegador *Firefox*).

En los sistemas derivados de MS-DOS puede verse aún un remanente de los archivos estructurados: En estos sistemas, un archivo puede ser *de texto* o *binario*. Un archivo de texto está compuesto por una serie de caracteres que forman *líneas*, y la separación entre una línea y otra constituye de un *retorno de carro* (CR, carácter ASCII 13) seguido de un *salto de línea* (LF, carácter ASCII 10).⁸

El acceso a los archivos puede realizarse de diferentes maneras:

Acceso secuencial Mantiene la semántica por medio de la cual permite leer de nuestros

⁸Esta lógica es herencia de las máquinas de escribir manuales, en que el *salto de línea* (avanzar el rodillo a la línea siguiente) era una operación distinta a la del *retorno de carro* (devolver la cabeza de escritura al inicio de la línea). En la época de los teletipos, como medida para evitar que se perdieran caracteres mientras la cabeza volvía hasta la izquierda, se decidió separar el inicio de nueva línea en los dos pasos que tienen las máquinas de escribir, para inducir una demora que evitara la pérdida de información.

archivos de forma equivalente a unidad de cinta mencionados en la sección 7.2.1, y como lo ilustra la figura 7.2: El mecanismo principal para leer o escribir es ir avanzando consecutivamente por los bytes que conforman al archivo hasta llegar a su final.

Típicamente se emplea este mecanismo de lectura para leer a memoria código (programas o bibliotecas) o documentos, sea enteros o fracciones de los mismos. Para un contenido estructurado, como una base de datos, resultaría absolutamente ineficiente, dado que no se conoce el punto de inicio o finalización de cada uno de los registros, y probablemente sería necesario que hacer *barridos secuenciales* del archivo completo para cada una de las búsquedas.



Figura 7.2: Archivo de acceso secuencial

Acceso aleatorio El empleo de gestores como *SQLite* u otros muchos motores de base de datos más robustos no exime al usuario de pensar en el archivo como una tabla estructurada, como lo ilustra la figura 7.3. Si la única semántica por medio de la cual el sistema operativo permitiera trabajar con los archivos fuera la equivalente a una unidad de cinta, implementar el acceso a un punto determinado del archivo podría resultar demasiado costoso.

Afortunadamente, que el sistema operativo no imponga registros de longitud fija no impide que el *programa gestor* lo haga. Si en el archivo al cual apunta el descriptor de archivo FD hay 2000 registros de 75 bytes cada uno y el usuario requiere recuperar el registro número 65 hacia el buffer `registro`, puede *repositionar* el apuntador de lectura al byte $65 \times 75 = 4875$ (`seek (FD, 4875)`) y leer los siguientes 75 bytes en `registro` (`read (FD, *registro, 75)`).

	Nombre	Apellido	Teléfono	Correo	ultimaSesion	usuarioDesde
0
4800	José	Chávez	5154-4553	chavez@aqul.no.es	2013.04.05	2012.01.15
4875	Gonzalo	Oliva				
4950	Raquel	Domínguez		rdomgz@aca.si.es		
5025
150000

Figura 7.3: Archivo de acceso aleatorio

Acceso relativo a índice En los últimos años se han popularizado los gestores de base de datos *débilmente estructurados* u *orientados a documentos*, llamados genéricamente *NoSQL*. Estos gestores pueden guardar registros de tamaño variable en disco, por lo que, como lo ilustra la figura 7.4, no pueden encontrar la ubicación correcta por medio de los mecanismos de acceso aleatorio.

Para implementar este acceso, se divide al conjunto de datos en dos secciones (incluso, posiblemente, en dos archivos independientes): La primer sección es una lista corta de identificadores, cada uno con el punto de inicio y término de los datos a los que apunta. Para leer un registro, se emplea acceso aleatorio sobre el índice, y

el apuntador se avanza a la ubicación específica que se solicita.

En el transcurso de un uso intensivo de esta estructura, dado que la porción de índice es muy frecuentemente consultada y relativamente muy pequeña, muy probablemente se mantenga completa en memoria, y el acceso a cada uno de los registros puede resolverse en tiempo muy bajo.

La principal desventaja de este modelo de indexación sobre registros de longitud variable es que sólo resulta eficiente para contenido *mayormente de lectura*: Cada vez que se produce una escritura y cambia la longitud de los datos almacenados, se va generando fragmentación en el archivo, y para resolverla probablemente se hace necesario suspender un tiempo la ejecución de todos los procesos que lo estén empleando (e invalidar, claro, todas las copias en caché de los índices). Ahora bien, para los casos de uso en que el comportamiento predominante sea de lectura, este formato tendrá la ventaja de no desperdiciar espacio en los campos nulos o de valor irrelevante para algunos de los registros, y de permitir la flexibilidad de registrar datos originalmente no contemplados sin tener que modificar la estructura.

Es importante recalcar que la escritura en ambas partes de la base de datos (índice y datos) debe mantenerse con garantías de atomicidad — Si se pierde la sincronía entre ellas, el resultado será una muy probable corrupción de datos.

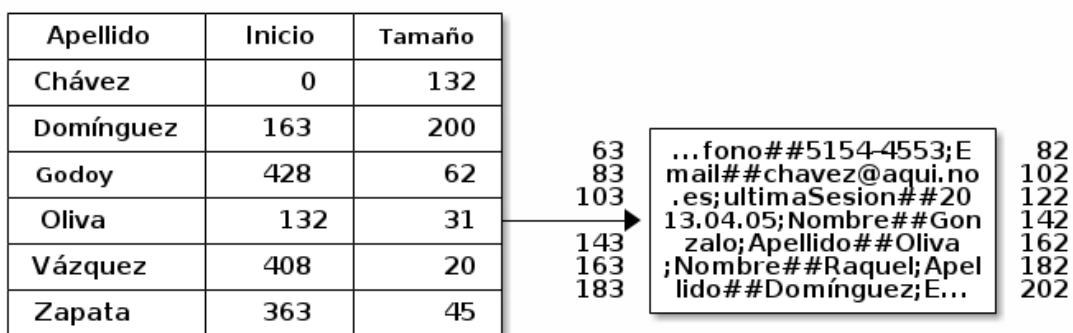


Figura 7.4: Acceso relativo a índice: Un índice apuntando al punto justo de un archivo sin estructura

7.2.6 Transferencias orientadas a bloques

Un sistema de archivos es la representación que se da a un conjunto de archivos y directorios sobre un *dispositivo de bloques*, esto es, un dispositivo que, para cualquier transferencia solicitada desde o hacia él, responderá con un bloque de tamaño predefinido.⁹

Esto es, si bien el sistema operativo presenta una abstracción por medio de la cual la lectura (`read()`) puede ser de un tamaño arbitrario, todas las transferencias de datos desde cualquiera de los discos serán de un múltiplo del tamaño de bloques, definido por el hardware (típicamente 512 bytes).

Al leer, como en el ejemplo anterior, sólamente un registro de 75 bytes, el sistema operativo lee el bloque completo y probablemente lo mantiene en un caché en la memoria principal; si en vez de una lectura, la operación efectuada fue una de escritura (`write()`), y el sector a modificar no ha sido leído aún a memoria (o fue leído hace mucho, y puede haber sido expirado del caché), el sistema tendrá que leerlo nuevamente, modificarlo en memoria, y volver a guardarlo a disco.

⁹La sección 3.8.2 define los dispositivos *de caracteres* y */de bloques*.

7.3 Organización de archivos

Hasta este punto, el enfoque ha sido en qué es y cómo se maneja un archivo. Sin embargo, no tendría sentido hablar de *sistemas de archivos* si no hubiera una gran cantidad de archivos. Es común que un sólo medio de almacenamiento de un equipo de uso doméstico aloje *decenas de miles* de archivos, y en equipos dedicados, no está fuera de lugar tener cientos o miles de veces tanto. Por tanto, se tiene que ver también cómo organizar una gran cantidad de archivos.

7.3.1 Evolución del concepto de *directorio*

El concepto dominante en almacenamiento hoy en día es el de *directorios jerárquicos*. Esto no siempre fue así; conviene revisar brevemente su historia para explicar el por qué de ciertos detalles de implementación del esquema actualmente dominante.

Convenciones de nomenclatura

Cada sistema de archivos puede determinar cuántos y qué caracteres son válidos para designar a uno de sus elementos, y cuáles son separadores válidos. El carácter que se emplea para separar los elementos de un directorio no es un estándar a través de todos los sistemas operativos — Los más comunes en uso hoy en día son la diagonal (/), empleada en sistemas tipo Unix y derivados (incluyendo MacOS X y Android), y la diagonal invertida (\), empleada en CP/M y derivados, incluyendo MS-DOS y Windows.

Diversos sistemas han manejado otros caracteres (por ejemplo, el MacOS histórico empleaba los dos puntos, :), y aunque muchas veces los mantenían ocultos del usuario a través de una interfaz gráfica rica, los programadores siempre tuvieron que manejarlos explícitamente.

A lo largo del presente texto se empleará la diagonal (/) como separador de directorios.

Sistema de archivos *plano*

Los primeros sistemas de archivos limitaban el concepto de directorio a una representación plana de los archivos que lo conformaban, sin ningún concepto de *jerarquía de directorios* como el que hoy resulta natural a los usuarios. Esto se debía, en primer término, a lo limitado del espacio de almacenamiento de las primeras computadoras en implementar esta metáfora (por lo limitado del espacio de almacenamiento, los usuarios no dejaban sus archivos a largo plazo en el disco, sino que los tenían ahí meramente mientras los requerían), y en segundo término, a que no se había aún desarrollado un concepto de separación, permisos y privilegios como el que poco después aparecería.

En las computadoras personales los sistemas de archivos eran también planos en un primer momento, pero por otra razón: En los sistemas *profesionales* ya se había desarrollado el concepto; al aparecer la primera computadora personal en 1975, ya existían incluso las primeras versiones de Unix diseñadas para trabajo en red. La prioridad en los sistemas personales era mantener el código del sistema operativo simple, mínimo. Con unidades de disco capaces de manejar entre 80 y 160KB, no tenía mucho sentido implementar directorios — Si un usuario quisiera llevar a cabo una división temática de su trabajo, lo colocaría en distintos *discos flexibles*. El sistema operativo CP/M nunca soportó jerarquías de directorios, como tampoco lo hizo la primera versión de MS-DOS.¹⁰

El sistema de archivos original de la Apple Macintosh, MFS, estaba construido sobre un modelo plano, pero presentando la *ilusión* de directorios de una forma comparable a las etiquetas: Existían bajo *certas* vistas (pero notoriamente no en los diálogos de abrir y grabar archivos), pero el nombre de cada uno de los archivos tenía que ser único, dado que el directorio al que pertenecía era básicamente sólo un atributo del archivo.

¹⁰El soporte de jerarquías de directorios fue introducido apenas en la versión 2, junto con el soporte a discos duros de 10MB, acompañando al lanzamiento de la IBM PC modelo XT.

Y contrario a lo que dicta la intuición, el modelo de directorio plano no ha desaparecido: El sistema de *almacenamiento en la nube* ofrecido por el servicio Amazon S3 (*Simple Storage Service, Servicio Simple de Almacenamiento*) maneja únicamente *objetos* (comparable con nuestra definición de archivos) y *cubetas* (de cierto modo comparables con las *unidades o volúmenes*), y permite referirse a un objeto o un conjunto de objetos basado en *filtros* sobre el total que conforman a una cubeta.

Conforme se desarrollen nuevas interfaces al programador o al usuario, probablemente se popularicen más ofertas como la que hoy hace Amazon S3. Al día de hoy, sin embargo, el esquema jerárquico sigue, con mucho, siendo el dominante.

Directarios de profundidad fija

Las primeras implementaciones de directorios eran *de un sólo nivel*: El total de archivos en un sistema podía estar dividido en directorios, fuera por tipo de archivo (separando, por ejemplo, programas de sistema, programas de usuario y textos del correo), por usuario (facilitando una separación lógica de los archivos de un usuario de pertenecientes a los demás usuarios del sistema)

El directorio *raíz* (base) se llama en este esquema *MFD* (*Master File Directory, Directorio Maestro de Archivos*), y cada uno de los directorios derivados es un *UFD* (*User File Directory, Directorio de Archivos de Usuario*).

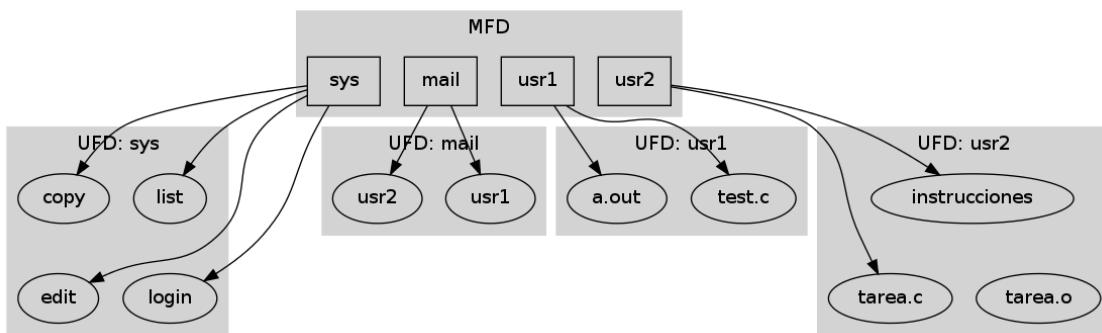


Figura 7.5: Directorio simple, limitado a un sólo nivel de profundidad

Este esquema resuelve el problema principal del nombre global único: Antes de los directorios, cada usuario tenía que cuidar que los nombres de sus archivos fueran únicos en el sistema, y ya teniendo cada uno su propio espacio, se volvió una tarea mucho más simple. La desventaja es que, si el sistema restringe a cada usuario a escribir en su UFD, se vuelve fundamentalmente imposible trabajar en algún proyecto conjunto: No puede haber un directorio que esté tanto dentro de *usr1* como de *usr2*, y los usuarios encontrarán más difícil llevar un proyecto conjunto.

Directarios estructurados en árbol

El siguiente paso natural para este esquema es permitir una *jerarquía ilimitada*: En vez de exigir que exista una capa de directorios, se le puede *dar la vuelta* al argumento, y permitir que cada directorio pueda contener a otros archivos o directorios a niveles arbitrarios. Esto permite que cada usuario (y que el administrador del sistema) estructure su información siguiendo criterios lógicos y piense en el espacio de almacenamiento como un espacio a largo plazo.

Junto con esta estructura nacen las *rutas de búsqueda* (*search path*): Tanto los programas como las bibliotecas de sistema ahora pueden estar en cualquier lugar del sistema de archivos. Al definirle al sistema una *ruta de búsqueda*, el usuario operador puede desentenderse del lugar exacto en el que está determinado programa — El sistema se encargará de buscar

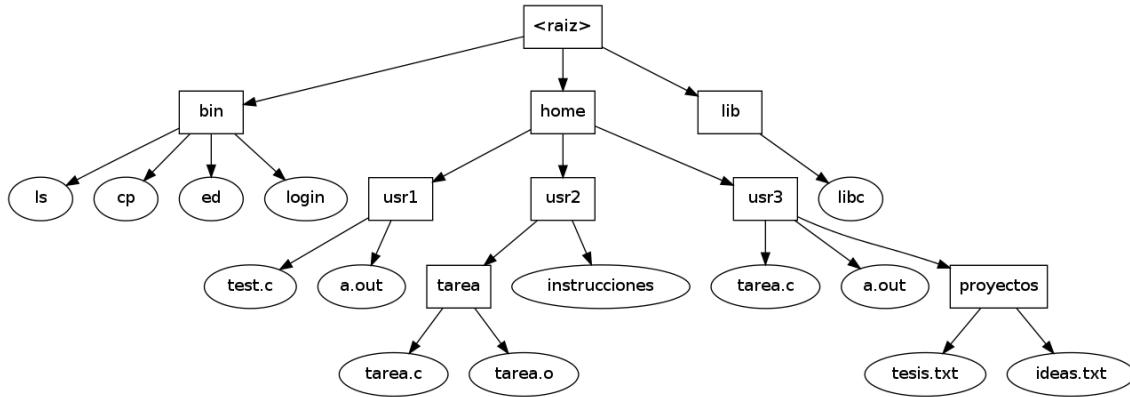


Figura 7.6: Directorio estructurado en árbol

en todos los directorios mencionados los programas o bibliotecas que éste requiera.¹¹

El directorio como un *grafo dirigido*

Si bien parecería que muchos de los sistemas de archivos empleados hoy en día pueden modelarse suficientemente con un árbol, donde hay un sólo nodo raíz, y donde cada uno de los nodos tiene un sólo nodo padre, la semántica que ofrecen es en realidad un *superconjunto estricto* de esta: La de un grafo dirigido.

En un grafo dirigido como el presentado en la figura 7.7, un mismo nodo puede tener varios directorios *padre*, permitiendo por ejemplo que un directorio de trabajo común sea parte del directorio personal de dos usuarios. Esto es, *el mismo objeto* está presente en más de un punto del árbol.

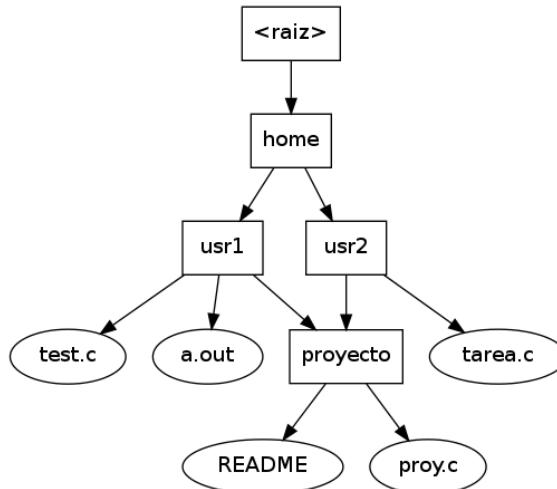


Figura 7.7: Directorio como un *grafo dirigido acíclico*: El directorio *projeto* está tanto en el directorio /home/usr1 como en el directorio /home/usr2

Un sistema de archivos puede permitir la organización como un *grafo dirigido*, aunque es común que la interfaz que presenta al usuario¹² se restrinja a un *grafo dirigido acíclico*:

¹¹La {ruta de búsqueda} refleja la organización del sistema de archivos en el contexto de la instalación específica. Es común que la ruta de búsqueda de un usuario estándar en Unix sea similar a /usr/local/bin:/usr/bin:/bin:~/bin — Esto significa que cualquier comando que sea presentado es buscado, en el orden indicado, en los cuatro directorios presentados (separados por el carácter :, la notación ~ indica el directorio personal del usuario activo). En Windows, es común ver una ruta de búsqueda c:\WINDOWS\system32;c:\WINDOWS

¹²Esta simplificación es simplemente una abstracción, y contiene una pequeña mentira, que será desmentida

Las ligas múltiples son permitidas, siempre y cuando no generen un ciclo.

La semántica de los sistemas Unix implementa directorios como grafos dirigidos por medio de dos mecanismos:

Liga o enlace duro La entrada de un archivo en un directorio Unix es la relación entre la ruta del archivo y el *número de i-nodo* en el sistema de archivos.¹³ Si a partir de un archivo existente se crea una *liga dura* a él, ésta es sencillamente otra entrada en el directorio apuntando al mismo *i-nodo*. Ambas entradas, pues, son el mismo archivo — No hay uno *maestro* y uno *dependiente*.

En un sistema Unix, este mecanismo tiene sólo dos restricciones:

1. Sólo se pueden hacer ligas duras dentro del mismo volumen.
2. No pueden hacerse ligas duras a directorios, sólo a archivos.¹⁴

Liga o enlace simbólico Es un archivo *especial*, que meramente indica a dónde apunta. El encargado de seguir este archivo a su destino (esto es, de *resolver* la liga simbólica) es el sistema operativo mismo; un proceso no tiene que hacer nada especial para seguir la liga.

Una liga simbólica puede *apuntar* a directorios, incluso creando ciclos, o a archivos en otros volúmenes.

Cuando se crea una liga simbólica, la liga y el archivo son dos entidades distintas. Si bien cualquier proceso que abra al archivo destino estará trabajando con la misma entidad, en caso de que éste sea renombrado o eliminado, la liga quedará *rota* (esto es, apuntará a una ubicación inexistente).

Si bien estos dos tipos de liga existen también en los sistemas Windows¹⁵, en dichos sistemas sigue siendo más común emplear los *accesos directos*. Se denomina así a un archivo (identificado por su extensión, .lnk), principalmente creado para poder *apuntar* a los archivos desde el escritorio y los menús — Si un proceso solicita al sistema abrir el *acceso directo*, no obtendrá al archivo destino, sino que al acceso directo mismo.

Ahora, si bien tanto las ligas duras como las ligas simbólicas existen también en Windows, su uso es muy poco frecuente. El API de Win32 ofrece las funciones necesarias, pero éstas no están reflejadas desde la interfaz usuario del sistema — Y son sistemas donde el usuario promedio no emplea una interfaz programador, sino que una interfaz gráfica. Las ligas, pues, no son más empleadas por *questión cultural*: En sus comunidades de usuarios, nunca fueron frecuentes, por lo cual se mantienen como conceptos empleados sólo por los *usuarios avanzados*.

Ya con el conocimiento de las ligas, y reelaborando la figura 7.7 con mayor apego a la realidad: En los sistemas operativos (tanto Unix como Windows), todo directorio tiene dos entradas especiales: Los directorios . y .., que aparecen tan pronto como el directorio es creado, y resultan fundamentales para mantener la *navegabilidad* del árbol.

Como se puede ver en la figura 7.8, en todos los directorios, . es una liga dura al mismo directorio, y .. es una liga al directorio *padre* (de nivel jerárquico inmediatamente superior). Claro está, como sólo puede haber una liga .., un directorio enlazado desde dos lugares distintos sólo apunta hacia uno de ellos con su enlace ..; en este caso, el directorio común *proyecto* está dentro del directorio /home/usr2. La figura representa la *liga simbólica* desde /home/usr1 como una línea punteada.

Hay una excepción a esta regla: El directorio raíz. En este caso, tanto . como .. apuntan al mismo directorio.

en breve.

¹³El significado y la estructura de un i-nodo se abordan en el capítulo 8.

¹⁴Formalmente, puede haberlas, pero sólo el administrador puede crearlas; en la sección 7.3.2 se cubre la razón de esta restricción al hablar de recorrer los directorios.

¹⁵Únicamente en aquellos que emplean el sistema de archivos NTFS, no en los que utilizan alguna de las variantes de FAT

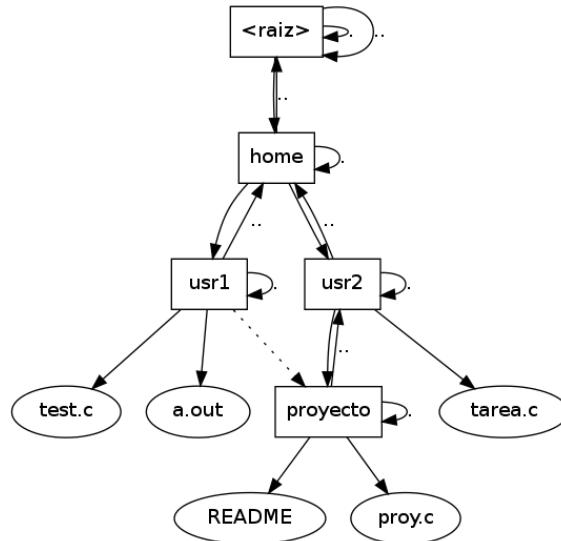


Figura 7.8: Directorio como un *grafo dirigido*, mostrando los *enlaces ocultos* al directorio actual . y al directorio *padre* ..

Esta es la razón por la cual no se puede tomar rigurosamente a un árbol de archivos como a un *grafo dirigido acíclico*, ni en Windows ni en Unix: Tanto las entradas . (al apuntar al mismo directorio donde están contenidas) como las entradas .. (al apuntar al directorio padre) crean ciclos.

7.3.2 Operaciones con directorios

Al igual que los archivos, los directorios tienen una semántica básica de acceso. Los directorios resultan también tipos de datos abstractos con algunas operaciones definidas. Muchas de las operaciones que pueden realizarse con los directorios son análogas a las empleadas para los archivos.¹⁶ Las operaciones básicas a presentar son:

Abrir y cerrar Al igual que los archivos, los directorios deben ser *abiertos* para trabajar con ellos, y *cerrados* cuando ya no se les requiera. Para esto, en C, se emplean las funciones `opendir()` y `closedir()`. Estas funciones trabajan asociadas a un *flujo de directorio* (*directory stream*), que funciona de forma análoga a un descriptor de archivo.

Listado de archivos Para mostrar los archivos que conforman a un directorio, el directorio se *abre* (tal como se haría con un archivo, pero empleando la función `opendir()` en vez de `open()`), y va *leyendo* (con `readdir()`) cada una de sus entradas. Cada uno de los resultados es una estructura `dirent` (*directory entry*, esto es, *entrada de directorio*), que contiene su nombre en `d_name`, un apuntador a su *i-nodo* en `d_ino`, y algunos datos adicionales del archivo en cuestión.

Para presentar al usuario la lista de archivos que conforman un directorio, podría hacerse:

```

#include <stdio.h>
#include <dirent.h>
#include <sys/types.h>

int main(int argc, char *argv[]) {
    struct dirent *archivo;
  
```

¹⁶De hecho, en muchos sistemas de archivos los directorios son meramente archivos de tipo especial, que son presentados al usuario de forma distinta. En la sección 8.1.4 se presenta un ejemplo.

```

DIR *dir;
if (argc != 2) {
    printf("Indique el directorio a mostrar\n");
    return 1;
}
dir = opendir(argv[1]);
while ((archivo = readdir(dir)) != 0) {
    printf("%s\t", archivo->d_name);
}
printf("\n");
closedir(dir);
}

```

Al igual que en al hablar de archivos se puede *repositionar* (`seek()`) al descriptor de archivo, para *rebobinar* el descriptor del directorio al principio del listado se emplea `rewinddir()`.

Buscar un elemento Si bien en el transcurso del uso del sistema operativo resulta una operación frecuente que el usuario solicite el listado de archivos dentro de un directorio, resulta mucho más frecuente buscar a un archivo en particular. La llamada `fopen()` antes descrita efectúa una búsqueda similar a la presentada en el ejemplo de código anterior, claro está, deteniéndose cuando encuentra al archivo en cuestión.

Crear, eliminar o renombrar un elemento Si bien estas operaciones se llevan a cabo sobre el directorio, son invocadas a través de la semántica orientada a archivos: Un archivo es creado con `fopen()`, eliminado con `remove()`, y renombrado con `rename()`.

Recorriendo los directorios

Es frecuente tener que aplicar una operación a todos los archivos dentro de cierto directorio — Por ejemplo, para agrupar a un directorio completo en un archivo comprimido, o para copiar todos sus contenidos a otro medio. Procesar todas las entradas de un directorio, incluyendo las de sus subdirectorios, se denomina *recorrer el directorio* (en inglés, *directory traversal*).

Si se trabaja sobre un sistema de archivos plano, la operación de recorrido completo puede realizarse con un programa tan simple como el presentado en la sección anterior.

Al hablar de un sistema de profundidad fija, e incluso de un directorio estructurado en árbol, la lógica se complica levemente, dado que para recorrer el directorio es necesario revisar, entrada por entrada, si esta es a su vez un directorio (y en caso de que así sea, entrar y procesar a cada uno de sus elementos). Hasta aquí, sin embargo, se puede recorrer el directorio sin requerir de mantener estructuras adicionales en memoria representando el estado.

Sin embargo, al considerar a los grafos dirigidos, se vuelve indispensable mantener en memoria la información de todos los nodos que ya han sido tocados — en caso contrario, al encontrar ciclo (incluso si este es creado por mecanismos como las *ligas simbólicas*), se corre el peligro de entrar en un bucle infinito.

Para recorrer al directorio ilustrado como ejemplo en la figura 7.9, no bastaría tomar nota de las rutas de los archivos conforme son recorridas — Cada vez que los sean procesados, su ruta será distinta. Al intentar respaldar al directorio `/home/jose/proyecto`, por ejemplo, el recorrido resultante podría ser:

- `/home/jose/proyecto`
- `/home/jose/proyecto/miembros`
- `/home/jose/proyecto/miembros/jose`
- `/home/jose/proyecto/miembros/jose/proyectos`

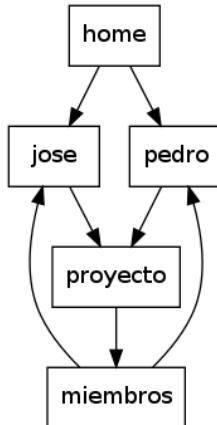


Figura 7.9: Directorio basado en grafo dirigido que incluye ciclos

- /home/jose/proyecto/miembros/jose/proyectos/miembros
- ... Y un etcétera infinito.

Para resolver esta situación, los programas que recorren directorios en los sistemas de archivos *reales* deben emplear un indexado basado en *i-nodo*¹⁷ identifica sin lugar a dudas a cada uno de los archivos. En el caso presentado, si el i-nodo de *jose* fuera 10543, al consultar a los miembros de *miembros* el sistema encontrará que su primer entrada apunta al i-nodo 10543, por lo cual la registraría sólo como un apuntador a datos *ya archivados*, y continuaría con la segunda entrada del directorio, que apunta a *pedro*.

Otros esquemas de organización

Por más que el uso de sistemas de archivos basados en directorios jerárquicos parece universal y es muy ampliamente aceptado, hay cada vez más casos de uso que apuntan a que se pueda estar por dar la bienvenida a una nueva metáfora de organización de archivos.

Hay distintas propuestas, y claro está, es imposible aún saber cuál dirección obtendrá el favor del mercado — O, dado que no necesariamente siga existiendo un modelo apto para todos los usos, de qué segmento del mercado.

7.3.3 Montaje de directorios

Para trabajar con el contenido de un sistema de archivos, el sistema operativo tiene que *montarlo*: Ubicarlo en algún punto del árbol de archivos visible al sistema y al usuario.

Es muy común, especialmente en los entornos derivados de Unix, que un sistema operativo trabaje con distintos sistemas de archivos al mismo tiempo. Esto puede obedecer a varias causas, entre las cuales se encuentran:

Distintos medios físicos Si la computadora tiene más de una unidad de almacenamiento, el espacio dentro de cada uno de los discos se maneja como un sistema de archivos independiente. Esto es especialmente cierto en la presencia de unidades removibles (CDs, unidades USB, discos duros externos, etc.)

Diferentes usos esperados Como se verá más adelante, distintos *esquemas de organización* (esto es, distintos sistemas de archivos) presentan ventajas para distintos patrones de uso. Por ejemplo, tiene sentido que una base de datos resida sobre una organización distinta a la de los programas ejecutables (binarios) del sistema.

Abstracciones de sistemas no-físicos El sistema operativo puede presentar diversas estructuras *con una estructura* de sistema de archivos. El ejemplo más claro de esto es el

¹⁷Que si bien no ha sido definido aún formalmente, para esta discusión bastará saber que es un número único por volumen.

sistema de archivos virtual `/proc`, existente en los sistemas Unix, que permite ver diversos aspectos de los procesos en ejecución (y, en Linux, del sistema en general). Los archivos bajo `/proc` no existen en ningún disco, pero se presentan como si fueran archivos estándar.

Razones administrativas El administrador del sistema puede emplear sistemas de archivos distintos para aislar espacios de usuarios entre sí: Por ejemplo, para evitar que un exceso de mensajes enviados en la bitácora (típicamente bajo `/var/log`) saturen al sistema de archivos principal, o para determinar patrones de uso máximo por grupos de usuarios.

En los sistemas tipo Unix, el mecanismo para montar los archivos es el de un árbol con *puntos de montaje*. Esto es, *todos los archivos y directorios* del sistema operativo están estructurados en *un único árbol*. Cuando se solicita al sistema operativo *montar* un sistema de archivos en determinado lugar, éste se integra al árbol, ocultando todo lo que el directorio en cuestión previamente tuviera.¹⁸

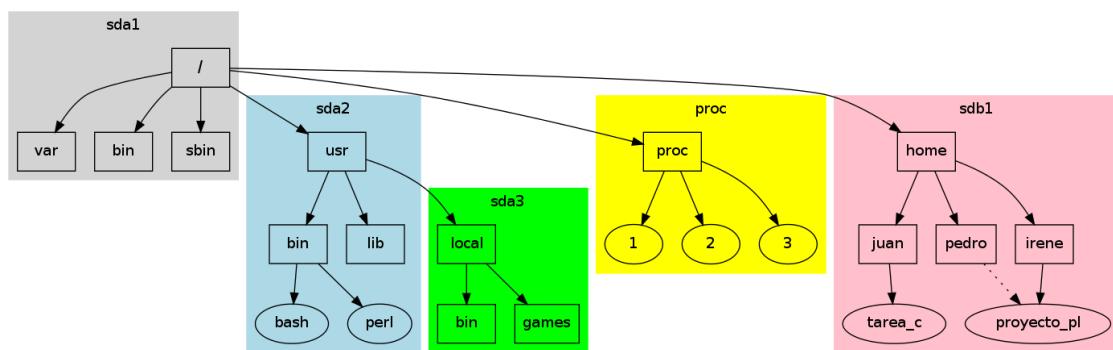


Figura 7.10: Árbol formado del montaje de sda1 en la raíz, sda2 como `/usr`, sdb1 como `/home`, y el directorio virtual `proc`

La manera en que esto se presenta en sistemas Windows es muy distinta. Ahí, cada uno de los volúmenes *detectados* recibe un *identificador de volumen*, y es montado automáticamente en un sistema de directorio estructurado como árbol de un sólo nivel representando a los dispositivos del sistema.¹⁹ Este árbol es presentado a través de la interfaz gráfica (aunque este nombre no significa nada para el API del sistema) como *Mi PC*.

Para especificar la *ruta completa* a determinado archivo, se inicia con el identificador del volumen. De este modo, la especificación absoluta de un archivo es una cadena como `VOL:\Dir1\Dir2\Archivo.ext` — El carácter `:` separa al volumen del árbol del sistema de archivos, y el carácter `\` separa uno de otro a los directorios. Por convención, si no se especifica la unidad, el sistema asumirá que se está haciendo referencia a la *unidad actual* (a grandes rasgos, la última unidad en ser utilizada).

Los identificadores de volumen están preasignados, muchos de ellos según a un esquema heredado desde la época de las primeras PC: Los volúmenes A y B están reservados para las unidades de disco flexible; C se refiere al disco duro de arranque, y las unidades posteriores que va detectando el sistema son D, E, F, etc.

Es posible modificar esta nomenclatura y configurar a los discos para estar en otra ubicación, pero muchas aplicaciones dependen ya de este comportamiento y configuración

¹⁸Hay implementaciones que exigen que el montaje se realice exclusivamente en directorios vacíos; existen otras, como UnionFS, que buscan seguir presentando una interfaz de *lectura* a los objetos que existían en el directorio previo al montaje, pero realizan las escrituras únicamente en el sistema ya montado; estas complican fuertemente algunos aspectos semánticos, por lo cual resultan poco comunes.

¹⁹En realidad, este árbol no sólo incluye a los volúmenes de almacenamiento, sino que a los demás dispositivos del sistema, como los distintos puertos, pero los *oculta* de la interfaz gráfica.

específica.

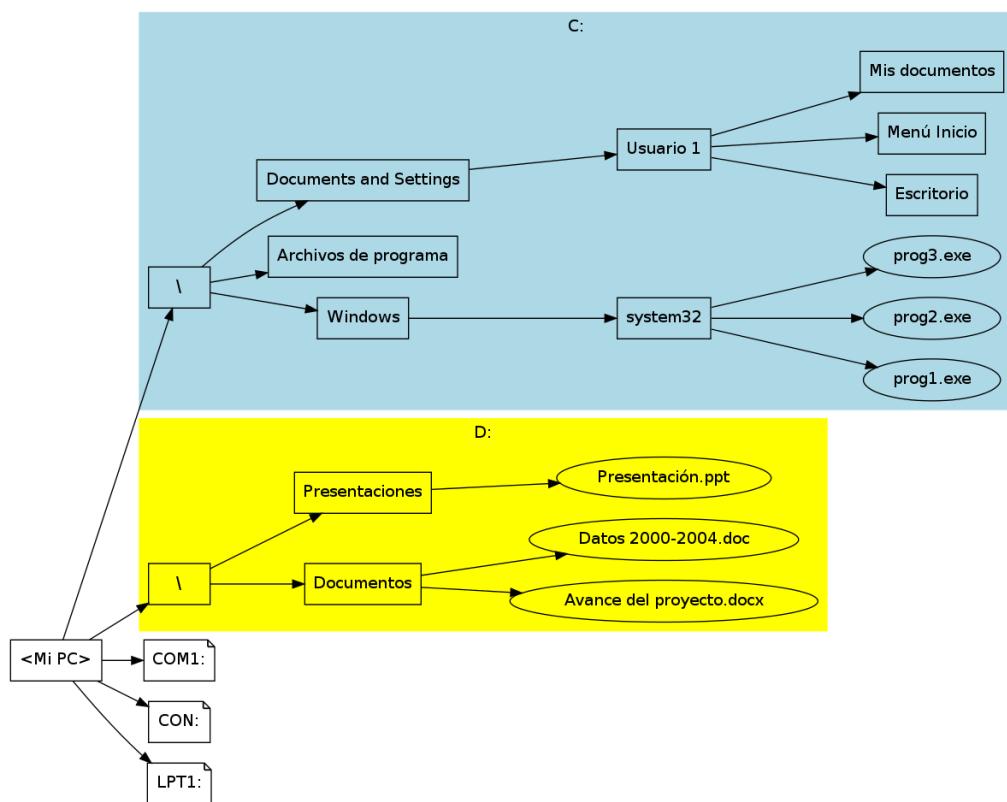


Figura 7.11: Vista de un sistema de archivos Windows

7.4 Sistemas de archivos remotos

Uno de los principales y primeros usos que se dio a la comunicación en red fue el de compartir archivos entre computadoras independientes. En un principio, esto se realizaba de forma *explícita*, con transferencias manuales a través de programas dedicados a ello, como sería hoy en día el FTP.

Por otro lado, desde mediados de los 1980, es posible realizar estas transferencias de forma *implícita* y *automática*, empleando *sistemas de archivos sobre la red* (o lo que es lo mismo, *sistemas de archivos remotos*). Éstos se presentan como caso particular de la *abstracción de sistemas no-físicos* que fueron mencionados en la sección anterior: Si bien el sistema operativo no tiene acceso *directo* a los archivos y directorios que le solicitará el usuario, a través de los módulos de red, sabe cómo obtenerlos y presentarlos *como si fueran locales*.

Al hablar de sistemas de archivos en red, casi siempre se hará siguiendo un *modelo cliente-servidor*. Estos términos no se refieren a las prestaciones relativas de una computadora, sino al rol que ésta juega *dentro de cada conexión* — Esto es, se designa como *cliente* a la computadora que solicita un servicio, y como *servidor* a la que lo provee; es frecuente que dos computadoras sean tanto servidor como cliente la una de la otra en distintos servicios.

7.4.1 Network File System (NFS)

El *Sistema de Archivos en Red* (*Network File System*, mejor conocido por sus siglas, *NFS*) fue creado por Sun Microsystems, y desde 1984 forma parte de su sistema operativo —

Resultó una implementación tan exitosa que a los pocos años formaba parte de todos los sistemas tipo Unix.

NFS está construido sobre el mecanismo *RPC* (*Remote Procedure Call, Llamada a Procedimientos Remotos*), un mecanismo de mensajes y manejo básico de sesión que actúa como una capa superior a TCP/IP, incluyendo facilidades de *descubrimiento de recursos y abstracción*. RPC puede ser comparado con protocolos como *DCE/RPC* de OSF, *DCOM* de Microsoft, y hoy en día, *SOAP* y *XML-RPC*. Estos mecanismos permiten al programador delegar en un *servicio* el manejo de las conexiones de red, particularmente (en el caso particular aquí descrito) la persistencia de sesiones en caso de desconexión, y limitar su atención a una *conexión virtual establecida*.

La motivación de origen para la creación de NFS fue presentar una solución que aprovechara el hardware existente y centralizara la administración: Ofrecer las facilidades para contar con redes donde hubiera un *servidor de archivos*, y donde las estaciones de trabajo tuvieran únicamente una instalación básica,²⁰ y el entorno de usuario completo estuviera disponible en cualquiera de las estaciones.

NFS ofrece sobre la red un sistema de archivos con la semántica Unix completa — Para montar un sistema remoto, basta montarlo²¹ y usarlo como si fuera local. El manejo de permisos, usuarios, e incluso las ligas duras y simbólicas se manejan exactamente como se haría localmente.

NFS es un protocolo muy ligero — No implementa cifrado ni verificaciones adicionales, pero al día de hoy, es uno de los mejores mecanismos para el envío de grandes cantidades de información — Pero siempre en redes que sean *completamente confiables*.

Ahora, NFS se presenta como uno de los componentes de una solución completa. Dado que se espera que la información de usuarios y permisos sea *consistente* en todos los clientes; Sun ofrecía también un esquema llamado *Yellow Pages* (posteriormente renombrado a NIS, *Network Information System*) para compartir la información de autenticación y listas de usuarios.

La desventaja, en entornos sin NIS, es que los permisos se manejan según el ID numérico del usuario. Si en diferentes sistemas el mismo usuario tiene diferentes IDs, los permisos no coincidirán. Es más, dado que el control de acceso principal es únicamente por dirección IP, para tener acceso irrestricto a los archivos de otros usuarios en NFS basta con tener control pleno de una computadora cualquiera en la red para poder *asumir o usurpar la identidad* de cualquier otro usuario.

Por último, para garantizar que las escrituras a archivos se llevaban a cabo cuando eran solicitadas (en contraposición a asumir éxito y continuar), todas las escrituras en un principio sobre NFS eran manejadas de forma *síncrona*, esto es, tras grabar un archivo, el cliente no continuaba con la ejecución hasta no tener confirmación por parte del servidor de que los datos estaban ya guardados en disco. Esto, si bien asegura que el usuario recibirá retroalimentación confiable respecto a la operación que realizó, ocasiona demoras que muchas veces son percibidas como excesivas.

Versiones posteriores del protocolo mejoraron sobre los puntos débiles aquí mencionados. Al día de hoy, casi 30 años después de su presentación, NFS es aún un sistema de archivos en red muy ampliamente empleado.

²⁰Incluso manejando estaciones de trabajo *diskless*, esto es, computadoras sin disco duro, cuyo sistema de arranque tiene la capacidad de solicitar al servidor le envíe incluso el núcleo del sistema operativo que ejecutará

²¹Para montar un sistema remoto, se emplea un comando como `mount archivos.unam.mx:/ext/home /home`, con lo cual el directorio `/ext/home` ubicado en el servidor `archivos.unam.mx` aparecerá montado como directorio `/home` local.

7.4.2 Common Internet File System (CIFS)

El equivalente a NFS en los entornos donde predominan los sistemas Windows es el protocolo CIFS (*Common Internet File System*, Sistema de Archivos Común para Internet). Aparece en los sistemas primarios de Microsoft alrededor de 1990²², originalmente bajo el nombre SMB (*Server Message Block, Bloque de Mensaje del Servidor*).

Las primeras implementaciones estaban basadas en el protocolo *NBF*, frecuentemente conocido como *NetBEUI*, un protocolo no ruteable diseñado para redes pequeñas y entornos sencillos de oficina. A partir de Windows 2000 se ha reimplementado completamente para operar sobre TCP/IP. Es a partir de este momento que se le comienza a denominar *CIFS*, aunque el nombre *SMB* sigue siendo ampliamente utilizado.²³

CIFS se ajusta mucho más a la semántica de los sistemas MS-DOS y Windows, aunque dado el lapso de tiempo que ha existido, ha pasado por varios cambios fundamentales, que al día de hoy complican su uso.

Para tener acceso a un volumen compartido por SMB se introdujo el comando *NET*,²⁴ basta indicar a DOS o Windows (desde la línea de comando) *NET USE W: \\servidor\directorio* para que el recurso compartido bajo el nombre *directorio* dentro del equipo conocido como *servidor* aparezca en el árbol *Mi PC*, y el usuario pueda emplear sus contenidos como si fuera un sistema de archivos local, con un volumen asignado de *W:*.

Cuando LAN Manager fue introducido al mercado, los sistemas Microsoft no manejaban aún el concepto de usuarios, por lo que la única medida de seguridad que implementaba SMB era el manejo de hasta dos contraseñas por directorio compartido: Con una, el usuario obtenía acceso de sólo lectura, y con la otra, de lectura y escritura. Tras la aparición de Windows NT, se agregó un esquema de identificación por usuario/contraseña, que posibilita el otorgamiento de permisos con una *granularidad* mucho menor.²⁵

SMB fue pensado originalmente para una red *pequeña*, con hasta un par de decenas de equipos. La mayor parte de los paquetes eran enviados en modo de *difusión (broadcast)*, por lo que era fácil llegar a la saturación, y no existía un esquema centralizado de resolución de nombres, con lo que era frecuente *no encontrar* a determinado equipo.

Los cambios que CIFS presenta a lo largo de los años son muy profundos. Las primeras implementaciones presentan fuertes problemas de confiabilidad, rendimiento y seguridad, además de estar planteadas para su uso en un sólo tipo de sistema operativo; al día de hoy, estos puntos han todos mejorado fuertemente. En sistemas Unix, la principal implementación, *Samba*, fue creada haciendo ingeniería inversa sobre el protocolo; a lo largo de los años, se ha convertido en un esquema tan robusto que es hoy por hoy tomado como implementación refencia.

7.4.3 Sistemas de archivos distribuidos: Andrew File System (AFS)

Los dos ejemplos de sistema de archivos en red presentados hasta ahora comparten una visión *tradicional* del modelo cliente-servidor: Al ver el comando que inicializa una conexión, e incluso a ver la información que guarda el núcleo del cliente respecto a cualquiera de los archivos, resulta claro cuál es el servidor para cada uno de ellos.

Andrew File System, desarrollado en la Carnegie Mellon University²⁶ y publicado en 1989, plantea presentar un verdadero *sistema de archivos distribuido*, en el cual los *recursos*

²²El desarrollo de SMB nació como {LAN Manager}, originalmente para OS/2

²³Es debido a este nombre que la implementación de CIFS para sistemas Unix, *Samba*, fue llamado de esta manera.

²⁴Este comando es empleado en MS-DOS, pero está también disponible en Windows, y al día de hoy es una de las principales herramientas para administrar usuarios.

²⁵Esto significa, que puede controlarse el acceso permitido más finamente, a nivel archivo individual y usuario individual.

²⁶Como parte del *Proyecto Andrew*, denominado así por el nombre de los fundadores de esta universidad:

compartidos no tengan que estar en un servidor en particular, sino que un conjunto de equipos se repartan la carga (esto es, *agnosticismo a la ubicación*). AFS busca también una *fácil escalabilidad*, la capacidad de agregar tanto espacio de almacenamiento como equipos con rol de servidor. AFS permite inclusive migrar completamente un volumen mientras está siendo empleado, de forma transparente.

Ante la complejidad e inestabilidad adicional que presentan con tanta frecuencia las redes grandes²⁷ (y lo hacían mucho más hace 30 años): AFS debe operar tan confiablemente como sea posible, *incluso sin la certeza de que la red opera correctamente*.

AFS construye fuertemente sobre el modelo de *tickets* y credenciales de *Kerberos*,²⁸ pero se aleja sensiblemente de la semántica de operación de archivos que hasta ahora se han presentado. Muchos eventos, operaciones y estados van ligados al *momento en el tiempo* en que se presentan, a través de un *modelo de consistencia débil* (*weak consistency model*). Muy a grandes rasgos, esto significa que:

- Cuando se abre un archivo, éste se copia completo al cliente. Todas las lecturas *y escrituras* (a diferencia de los esquemas tradicionales, en que éstas son enviadas al servidor *lo antes posible* y de forma síncrona) se dirigen únicamente a la copia local.
- Al cerrar el archivo, éste se copia de vuelta al *servidor de origen*, el cual se *compromete* a notificar a los clientes si un archivo abierto fue modificado (esto es, a *hacer una llamada o callback*). Los clientes pueden entonces intentar incorporar los cambios a su versión de trabajo, o continuar con la copia ya obtenida — Es *esperable* que si un segundo cliente realiza alguna modificación, incorpore los cambios hechos por el primero, pero esta responsabilidad se deja a la implementación del programa en cuestión.

Esto significa en pocas palabras que los cambios a un archivo abierto por un usuario no son visibles a los demás de inmediato; sólo una vez que se cierra un archivo, los cambios hechos a éste son puestos a disposición de las sesiones abiertas actuales, y sólo son enviados como *versión actual* a las sesiones abiertas posteriormente.

Con este cambio semántico, debe quedar claro que AFS no busca ser un sistema para todo uso ni un reemplazo universal de los sistemas de archivos locales, en contraposición de los sistemas de archivos centralizados. AFS no plantea en ningún momento una operación *diskless*. Bajo el esquema aquí descrito, las lecturas y escrituras resultan baratas, porque se realizan exclusivamente sobre el caché local, pero abrir y cerrar un archivo puede ser muy caro, porque debe transferirse el archivo completo.

Hay aplicaciones que verdaderamente sufrirían si tuvieran que implementarse sobre un sistema de archivos distribuído — Por ejemplo, si una base de datos se distribuyera sobre AFS, la carencia de mecanismos de bloqueo sobre *secciones* del archivo, y el requisito de operar sobre *archivos completos* harían impracticable compartir un archivo de uso intensivo y aleatorio.

7.5 Otros recursos

- *File System Interface: Functions for manipulating files*
https://www.gnu.org/software/libc/manual/html_node/File-System-Interface.html
The GNU C Library manual (Free Software Foundation)
- *Disk from the Perspective of a File System*
<http://queue.acm.org/detail.cfm?id=2367378>

Andrew Carnegie y Andrew Mellon

²⁷El uso típico de AFS se planteaba para organizaciones grandes, del orden de decenas de miles de estaciones

²⁸Un sistema de autenticación y autorización centralizado para entornos corporativos.

Marshall Kirk McKusick (2012); ACM Queue

- Traducción al español: *Los discos desde la perspectiva de un sistema de archivos*

<http://cyanezfdz.me/post/los-discos-desde-la-perspectiva-de-un-sistema-de-archivos>
César Yáñez (2013).

- *File-system development with stackable layers*

<https://dl.acm.org/citation.cfm?id=174613.174616>

Heidemann y Popek (1994); ACM Transactions on Computer Systems

- *Serverless network file systems*

<https://dl.acm.org/citation.cfm?doid=225535.225537>

Thomas E. Anderson et. al. (1996); ACM Transactions on Computer Systems

- *OpenPlanets Results Evaluation Framework*

<http://data.openplanetsfoundation.org/ref/>

David Tarrant (2012). Muestra la evolución a lo largo de los años de cómo reconocen archivos de tipos conocidos varias herramientas

- *Finding open files with lsof*

<http://www.ibm.com/developerworks/aix/library/au-lsop.html>)

Sean A. Walberg (2006); IBM DeveloperWorks

8 — Sistemas de archivos

8.1 Plasmando la estructura en el dispositivo

A lo largo del capítulo 7 se presentaron los elementos del sistema de archivos tal como son presentados al usuario final, sin entrar en detalles respecto a cómo organiza toda esta información el sistema operativo en un *dispositivo persistente* — Mencionamos algunas estructuras base, pero dejándolas explícitamente pendientes de definición. En este capítulo se tratarán las principales estructuras y mecanismos empleados para que un sistema de archivos sea ya no sólamente una estructura formal ideal, sino que una entidad almacenada en un dispositivo.

A lo largo de la historia del cómputo, el almacenamiento no siempre se realizó en discos (dispositivos giratorios de acceso aleatorio). En un principio, los medios principales de acceso estrictamente secuencial (tarjetas perforadas, cintas de papel, cintas magnéticas); por más de 30 años, el medio primario de almacenamiento han sido los distintos tipos de discos magnéticos, y desde hace algunos años, estamos viendo una migración a *almacenamiento de estado sólido*, a dispositivos sin partes móviles que guardan la información en un tipo particular de memoria. Volviendo a las categorías presentadas en la sección 3.8, los medios de acceso secuencial son *dispositivos de caracteres*, y tanto discos como unidades de estado sólido son *dispositivos de bloques*.

8.1.1 Conceptos para la organización

Los sistemas de archivo están en general desarrollados pensando en *discos*, y a lo largo de este capítulo, se hará referencia como *el disco* al medio de almacenamiento persistente en el cual esté plasmado el sistema de archivos. En el apéndice C se tocarán algunos de los aspectos que debemos considerar al hablar de sistemas de archivos respaldados en medios *distintos* a un disco.

Mientras tanto, conviene mantener una visión aún bastante idealizada y abstracta: Un *disco* visto desde la perspectiva del sistema operativo será presentado a lo largo del presente capítulo¹ como un arreglo muy grande de *bloques* de tamaño fijo, cada uno de ellos *directamente direccionable*; esto significa que el sistema operativo puede referirse por igual a cualquiera de los bloques del disco a través de una dirección física e inambigua dentro del disco entero. Partiendo de esto, se emplean los siguientes conceptos para almacenar, ubicar o recuperar la información:

Partición Una subdivisión de un disco, por medio de la cual el administrador/usuario del sistema puede definir la forma en que se emplea el espacio del disco, segmentándolo si hace falta según haga falta.

Un disco puede tener varias particiones, y cada una de ellas puede tener un sistema de archivos independiente.

Volumen Colección de bloques *inicializados* con un sistema de archivos que pueden presentarse al usuario como una unidad. Típicamente un volumen coincide con una

¹Para una visión más rigurosa de cómo se relaciona el sistema operativo con los discos y demás mecanismos de almacenamiento, refiérase al apéndice C.

partición (pero no siempre es el caso, como se describirá en las secciones C.2 y C.3). El volumen se *describe* ante el sistema operativo en el *bloque de control de volumen*, también conocido como *superbloque* en Unix, o *Tabla Maestra de Archivos (Master File Table)* en NTFS.

Sistema de archivos Esquema de organización que sigue un determinado *volumen*. Dependiendo del sistema de archivos elegido, cada uno de los componentes aquí presentados ocuparán un distinto lugar en el disco, presentando una semántica propia.

Para poder tener acceso a la información almacenada en determinado volumen, el sistema operativo debe tener soporte para el sistema de archivos particular en que éste esté estructurado.

Directorio raíz La estructura que relaciona cada nombre de archivo con su números de *i-nodo*. Típicamente sólo almacena los archivos que están en el *primer nivel jerárquico* del sistema, y los directorios derivados son únicamente referenciados desde éste. En sistemas de archivos modernos, el directorio normalmente incluye sólo el nombre de cada uno de los archivos y el número de *i-nodo* que lo describe, todos los *metadatos* adicionales están en los respectivos *i-nodos*.

Metadatos Recibe este nombre toda la información *acerca de* un archivo que *no es* el contenido del archivo mismo. Por ejemplo, el nombre, tamaño o tipo del archivo, su propietario, el control de acceso, sus fechas de creación, último acceso y modificación, ubicación en disco, etc.

I-nodo Del inglés *i-node, information node* (nodo de información); en los sistemas tipo Windows, normalmente se le denomina *bloque de control de archivo (FCB)*. Es la estructura en disco que guarda los *metadatos* de cada uno de los archivos, proporcionando un vínculo entre la *entrada en el directorio* y los datos que lo conforman.

La información almacenada incluye todos los metadatos relacionados con el archivo *a excepción del nombre* (mismo que radica únicamente en el *directorío*): Los permisos y propietarios del archivo, sus fechas de creación, última modificación y último acceso, y la *relación de bloques* que ocupa en el disco. Más adelante se abordarán algunos de los esquemas más comunes para presentar esta relación de bloques.

Esta separación entre directorio e *i-nodo* permite a un mismo archivo formar parte de distintos directorios, como se explicó en la sección 7.3.1.

Mapa de bits de espacio libre La función del bitmap es poder gestionar el espacio libre del disco. Recuérdese que el disco se presenta asignado por *bloques*, típicamente de 4096 bytes — En el bitmap cada bloque se representa con un bit, con lo que aquí se puede encontrar de forma compacta el espacio ocupado y disponible, así como el lugar adecuado para crear un nuevo archivo.

El bitmap para un disco de 100GB puede, de esta manera, representarse en 23MB ($\frac{100 \times 10^9}{4096}$), cantidad que puede razonablemente mantener en memoria un sistema de escritorio promedio hoy en día.²

Más adelante se verán algunas estructuras avanzadas que permiten mayor eficiencia en este sentido.

8.1.2 Diferentes sistemas de archivos

Un sistema operativo puede dar soporte a varios distintos sistemas de archivos; un administrador de sistemas puede tener muy diferentes razones que influyan para elegir cuál sistema de archivos empleará para su información — Algunas razones para elegir a uno u otro son que el rendimiento de cada uno puede estar *afinado* para diferentes

²Esto explica por qué, incluso sin estar trabajando activamente con ningún archivo contenido en éste, el sólo hecho de montar un volumen con gran cantidad de datos obliga al sistema a reservarle una cantidad de memoria.

patrones de carga, necesidad de emplear un dispositivo portátil para intercambiar datos con distintos sistemas, e incluso restricciones de hardware.³

A lo largo de esta sección se revisará cómo los principales conceptos a abordar se han implementado en distintos sistemas de archivos; se hará referencia principalmente a una familia de sistema de archivos simple de comprender, aunque muestra claramente su edad: El sistema FAT. La razón de elegir al sistema de archivos FAT es la simplicidad de sus estructuras, que permiten comprender la organización general de la información. Donde sea pertinente, se mencionará en qué puntos principales estiba la diferencia con los principales sistemas de la actualidad.

El sistema FAT fue creado hacia fines de los 1970, y su diseño muestra claras evidencias de haber sido concebido para discos flexibles. Sin embargo, a través de varias extensiones que se han presentado con el paso de los años (algunas con compatibilidad hacia atrás,⁴ otras no), sigue siendo uno de los sistemas más empleados al día de hoy, a pesar de que ya no es recomendado como sistema primario por ningún sistema operativo de escritorio.

Si bien FAT tuvo su mayor difusión con los sistemas operativos de la familia MS-DOS, es un sistema de archivos nativo para una gran cantidad de otras plataformas (muchas de ellas dentro del mercado *embebido*), lo cual se hace obvio al revisar el soporte a atributos extendidos que maneja.

8.1.3 El volumen

Lo primero que requiere saber el sistema operativo para poder montar un volumen es su estructura general: En primer término, de qué *tipo* de sistema de archivos se trata, y acto seguido, la descripción básica del mismo: Su extensión, el tamaño de los *bloques lógicos* que maneja, si tiene alguna *etiqueta* que describa su función ante el usuario, etc. Esta información está contenida en el *bloque de control de volumen*, también conocido como *superbloque* o *tabla maestra de archivos*.⁵

Tras leer la información del superbloque, el sistema operativo determina en primer término si puede proceder — Si no sabe cómo trabajar con el sistema de archivos en cuestión, por ejemplo, no puede presentar información útil alguna al usuario (e incluso arriesgaría destruirla).

Se mencionó ya que el tamaño de bloques (históricamente, 512 bytes; el estándar *Advanced Format* en marzo del 2010 introdujo bloques de 4096 bytes) es establecido por el hardware. Es muy común que, tanto por razones de eficiencia como para alcanzar a direccionar mayor espacio, el sistema de archivos *agrupue* a varios bloques físicos en un bloque lógico. En la sección 8.1.4 se revisará qué factores determinan el tamaño de bloques en cada sistema de archivos.

Dado que describir al volumen es la más fundamental de las operaciones a realizar, muchos sistemas de archivos mantienen *copias adicionales* del superbloque, a veces dispersas a lo largo del sistema de archivos, para poder recuperarlo en caso de que se corrompa.

En el caso de FAT, el volumen indica no sólo la *generación* del sistema de archivos que se está empleando (FAT12, FAT16 o FAT32, en los tres casos denominados así por la cantidad de bits para referenciar a cada uno de los bloques lógicos o *clusters*), sino el tamaño de los *clusters*, que puede ir desde los 2 y hasta los 32 Kb.

³Por ejemplo, los /cargadores de arranque} en algunas plataformas requieren poder leer el volumen donde está alojada la imagen del sistema operativo — Lo cual obliga a que esté en un sistema de archivos nativo a la plataforma.

⁴Se denomina *compatibilidad hacia atrás* a aquellos cambios que permiten interoperar de forma transparente con las versiones anteriores.

⁵Y aquí hay que aclarar: Este bloque {no contiene a los archivos}, ni siquiera a las estructuras que apuntan hacia ellos, sólo describe al volumen para que pueda ser montado

Volúmenes crudos

Si bien una de los principales tareas de un sistema operativo es la organización del espacio de almacenamiento en sistemas de archivos y su gestión para compartirse entre diversos usuarios y procesos, hay algunos casos en que un dispositivo orientado a bloques puede ser puesto a disposición de un proceso en particular para que éste lo gestione directamente. Este modo de uso se denomina el de *dispositivos crudos* o *dispositivos en crudo* (*raw devices*).

Pueden encontrarse dos casos de uso primarios hoy en día para dispositivos orientados a bloques no administrados a través de la abstracción de los sistemas de archivos:

Espacio de intercambio Como se vio en la sección 6.5.2, la gestión de la porción de la memoria virtual que está en disco es mucho más eficiente cuando se hace sin cruzar por la abstracción del sistema operativo — Esto es, cuando se hace en un volumen en crudo. Y si bien el gestor de memoria virtual es parte innegable del sistema operativo, en un sistema *microkernel* puede estar ejecutándose como proceso de usuario.

Bases de datos Las bases de datos relacionales pueden incluir volúmenes muy grandes de datos estrictamente estructurados. Algunos gestores de bases de datos, como Oracle, MaxDB o DB2, recomiendan a sus usuarios el uso de volúmenes crudos, para optimizar los accesos a disco sin tener que cruzar por tantas capas del sistema operativo.

La mayor parte de los gestores de bases de datos desde hace varios años no manejan esta modalidad, por la complejidad adicional que supone para el administrador del sistema y por lo limitado de la ventaja en rendimiento que supone hoy en día, aunque es indudablemente un tema que se presta para discusión e investigación.

8.1.4 El directorio y los i-nodos

El directorio es la estructura que relaciona a los archivos como son presentados al usuario –identificados por una ruta y un nombre– con las estructuras que los describen ante el sistema operativo — Los *i-nodos*.

A lo largo de la historia de los sistemas de archivos, se han implementado muy distintos esquemas de organización. Se presenta a continuación la estructura básica de la popular familia de sistemas de archivos FAT.

El directorio raíz

Una vez que el sistema de archivos está *montado* (ver 7.3.3), todas las referencias a archivos dentro de éste deben pasar a través del directorio. El directorio raíz está siempre en una ubicación *bien conocida* dentro del sistema de archivos — Típicamente al inicio del volumen, en los primeros sectores⁶. Un disco flexible tenía 80 *pistas* (típicamente denominadas *cilindros* al hablar de discos duros), con lo que, al ubicar al directorio en la pista 40, el tiempo promedio de movimiento de cabezas para llegar a él se reducía a la mitad. Si todas las operaciones de abrir un archivo tienen que pasar por el directorio, esto resultaba en una mejoría muy significativa.

El directorio es la estructura que determina el formato que debe seguir el nombre de cada uno de los archivos y directorios: Es común que en un sistema moderno, el nombre de un archivo pueda tener hasta 256 caracteres, incluyendo espacios, caracteres internacionales, etc. Algunos sistemas de archivos son *sensibles a mayúsculas*, como es el caso de los sistemas nativos a Unix (el archivo `ejemplo.txt` es distinto de `Ejemplo.TXT`), mientras que otros no lo son, como es el caso de NTFS y VFAT (`ejemplo.txt` y `Ejemplo.TXT` son idénticos ante el sistema operativo).

⁶Una excepción a esta lógica se presentó en la década de los 1980, cuando los diseñadores del sistema AmigaOS decidieron ubicar al directorio en el sector *central* de los volúmenes, para reducir a la mitad el tiempo promedio de acceso a la parte más frecuentemente referida del disco

Todas las versiones de FAT siguen para los nombres de archivos un esquema claramente arcáico: Los nombres de archivo pueden medir hasta 8 caracteres, con una extensión opcional de 3 caracteres más, dando un total de 11. El sistema no sólo no es sensible a mayúsculas y minúsculas, sino que todos los nombres deben guardarse completamente en mayúsculas, y permite sólo ciertos caracteres no alfanuméricos. Este sistema de archivos no implementa la separación entre directorio e i-nodo, que hoy es la norma, por lo que cada una de las entradas en el directorio mide exactamente 32 bytes. Como es de esperarse en un formato que ha vivido tanto tiempo y ha pasado por tantas generaciones como FAT, algunos de estos campos han cambiado substancialmente sus significados. La figura 8.1 muestra los campos de una entrada del directorio bajo FAT32.



Figura 8.1: Formato de la entrada del directorio bajo FAT (Mohammed, 2007)

La extensión VFAT fue agregada con el lanzamiento de Windows 95. Esta extensión permitía que, si bien el nombre *real* de un archivo seguiría estando limitada al formato presentado, pudieran agregarse entradas adicionales al directorio utilizando el atributo de *etiqueta de volumen* de maneras que un sistema MS-DOS debiera ignorar.⁷

Esto presenta una complicación adicional al hablar del directorio *raiz* de una unidad: Si bien los directorios derivados no tienen este límite, al estar el directorio raiz ubicado en una sección fija del disco, tiene una longitud límite máxima: En un disco flexible (que hasta el día de hoy, por su limitada capacidad, se formatea bajo FAT12), desde el bloque 20 y hasta el 33, esto es, 14 bloques. Con un tamaño de sector de 512 bytes, el directorio raiz mide $512 \times 14 = 7168$ bytes, esto es, $\frac{7168}{32=224}$ entradas como máximo. Y si bien esto puede no parecer muy limitado, ocupar cuatro entradas por archivo cuando, empleando VFAT, se tiene un nombre medianamente largo reduce fuertemente el panorama.

El problema no resulta tan severo como podría parecer: Para FAT32, el directorio raiz ya no está ubicado en un espacio reservado, sino que como parte del espacio de datos, por lo cual es extensible en caso de requerirse.

Los primeros sistemas de archivos estaban pensados para unidades de muy baja capacidad; por mucho tiempo, las implementaciones del directorio eran simplemente listas lineales con los archivos que estaban alojados en el volumen. En muchos de estos primeros sistemas no se contemplaban directorios jerárquicos, sino que presentaban un único espacio *plano* de nombres; cuando estos sistemas fueron evolucionando para soportar directorios anidados, por compatibilidad hacia atrás (y por consideraciones de rendimiento que se abordan a continuación) siguieron almacenando únicamente al directorio raiz en esta posición privilegiada, manejando a todos los directorios que derivaran de éste como si fueran archivos, repartidos por el disco.

En un sistema que implementa los directorios como listas lineales, entre más archivos haya, el tiempo que toma casi cualquier operación se incrementa linealmente (dado que potencialmente se tiene que leer al directorio completo para encontrar a un archivo). Y las listas lineales presentan un segundo problema: Cómo reaccionar cuando se *llena* el espacio que tienen asignado.

⁷La *etiqueta de volumen* estaba definida para ser empleada exclusivamente a la *cabeza* del directorio, dando una etiqueta global al sistema de archivos completo; el significado de una entrada de directorio con este atributo hasta antes de la incorporación de VFAT {no estaba definida}.

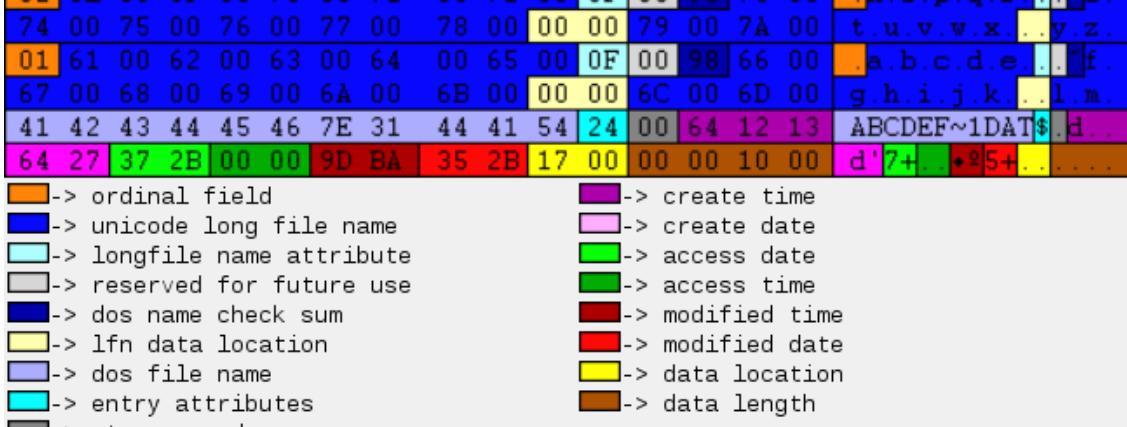
0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
43	2E	00	64	00	61	00	74	00	00	00	0F	00	98	FF	FF
FF	FF	FF	FF	FF	FF	FF	FF	FF	FF	00	00	FF	FF	FF	FF
02	6E	00	6F	00	70	00	71	00	72	00	0F	00	98	73	00
74	00	75	00	76	00	77	00	78	00	00	00	79	00	7A	00
01	61	00	62	00	63	00	64	00	65	00	0F	00	98	66	00
67	00	68	00	69	00	6A	00	6B	00	00	00	6C	00	6D	00
41	42	43	44	45	46	7E	31	44	41	54	24	00	64	12	13
64	27	37	2B	00	00	9D	BA	35	2B	17	00	00	00	10	00
															

Figura 8.2: Entradas representando archivo con nombre largo bajo VFAT (Imagen: Peter Clark, ver otros recursos)

Como ya se presentó, FAT asigna un espacio fijo al directorio raíz, pero los subdirectorios pueden crecer arbitrariamente. Un subdirectorio es básicamente una entrada con un tipo especial de archivo — Si el doceavo byte de una entrada de directorio, que indica los *atributos del archivo* (ver figura 8.1 y cuadro 7.2) tiene al bit 4 activado, la región de datos correspondientes a dicho archivo será interpretada como un subdirectorio.

La tabla de asignación de archivos

Queda claro que FAT es un sistema heredado, y que exhibe muchas prácticas que ya se han abandonado en diseños modernos de sistemas de archivos. Se vio que dentro de la entrada de directorio de cada archivo está prácticamente su *i-nodo* completo: La información de permisos, atributos, fechas de creación — Y muy particularmente, el apuntador al *cluster* de inicio (bytes 26-28, mas los 20-22 para FAT32). Esto resulta en una de las grandes debilidades de FAT: La tendencia a la fragmentación.

La familia FAT obtiene su nombre de la Tabla de Asignación de Archivos (*File Allocation Table*), que aparece antes del directorio, en los primeros sectores del disco.⁸ Cada byte de la FAT representa a un *cluster* en el área de datos; cada entrada en el directorio indica, en su campo *cluster*, cuál es el primer *cluster* que conforma al archivo. Ahora bien, conforme se usa un disco, y los archivos crecen y se eliminan, y van llenando los espacios vacíos que van dejando, FAT va asignando espacio *conforme encuentra nuevos clusters libres*, sin cuidar que sea espacio continuo. Los apuntadores al *siguiente cluster* se van marcando en la tabla, *cluster por cluster*, y el último *cluster* de cada archivo recibe el valor especial (dependiendo de la versión de FAT) 0xFFFF, 0xFFFFFFF o 0xFFFFFFFF.

Ahora bien, si los directorios son sencillamente archivos que reciben un tratamiento especial, estos son también susceptibles a la fragmentación. Dentro de un sistema Windows 95 o superior (empleando VFAT), con directorios anidados a cuatro o cinco niveles como lo establece su jerarquía estándar⁹, la simple tarea de recorrerlos para encontrar determinado archivo puede resultar muy penalizado por la fragmentación.

⁸ Esta tabla es tan importante que, dependiendo de la versión de FAT, se guarda por duplicado, o incluso por triplicado.

⁹ Por ejemplo, C:\Documents and Settings\Usuario\Menú Inicio\Programa Ejemplo\Programa Ejemplo.lnk

Cuadro 8.1: Significado de cada uno de los bits del byte de *atributos del archivo* en el directorio FAT. La semántica que se presenta es la empleada por los sistemas MS-DOS y Windows; otros sistemas pueden presentar comportamientos adicionales.

Bit	Nombre	Descripción
0	Sólo lectura	El sistema no permitirá que sea modificado.
1	Oculto	No se muestra en listados de directorio.
2	Sistema	El archivo pertenece al sistema y no debe moverse de sus <i>clusters</i> (empleado, por ejemplo, para los componentes a cargar para iniciar al sistema operativo)
3	Etiqueta	Indica el nombre del volumen, no un archivo. En VFAT, expresa la <i>continuación</i> de un nombre largo.
4	Subdirectorio	Los <i>clusters</i> que componen a este archivo son interpretados como un subdirectorio, no como un archivo.
5	Archivado	Empleado por algunos sistemas de respaldo para indicar si un archivo fue modificado desde su última copia.
6	Dispositivo	Para uso interno del sistema operativo, no fue adoptado para los archivos.
7	Reservado	Reservado, no debe ser manipulado.

La eliminación de entradas del directorio

Sólo unos pocos sistemas de archivos guardan sus directorios ordenados — Si bien esto facilitaría las operaciones más frecuentes que se realizan sobre de ellos (en particular, la búsqueda: Cada vez que un directorio es recorrido hasta encontrar un archivo tiene que leerse potencialmente completo), mantenerlo ordenado ante cualquier modificación resultaría mucho más *caro*, dado que tendría que reescribirse el directorio completo al crearse o eliminarse un archivo dentro de éste, y lo que es más importante, más *peligroso*, dado que aumentaría el tiempo que los datos del directorio están en un estado inconsistente, aumentando la probabilidad de que ante una interrupción repentina (fallo de sistema, corte eléctrico, desconexión del dispositivo, etc.) se presentara corrupción de la información llevando a pérdida de datos. Al almacenar las entradas del directorio sin ordenar, las escrituras que modifican esta crítica estructura se mantienen atómicas: Un sólo sector de 512 bytes puede almacenar 16 entradas básicas de FAT, de 32 bytes cada una.¹⁰

Ordenar las entradas del directorio teniendo sus contenidos ya en memoria y, en general, diferir las modificaciones al directorio resulta mucho más conveniente en el caso general. Esto vale también para la eliminación de archivos — A continuación se abordará la estrategia que sigue FAT. Cabe recordar que FAT fue diseñado cuando el medio principal de almacenamiento era el disco flexible, decenas de veces más lento que el disco duro, y con mucha menor confiabilidad.

Cuando se le solicita a un sistema de archivos FAT eliminar un archivo, éste no se borra del directorio, ni su información se libera de la tabla de asignación de archivos, sino que se *marca* para ser ignorado, reemplazando el primer carácter de su nombre por

¹⁰Aunque en el caso de VFAT, las diferentes entradas que componen un sólo nombre de archivo pueden quedar separadas en diferentes sectores.

0xE5. Ni la entrada de directorio, ni la *cadena de clusters* correspondiente en las tablas de asignación,¹¹ son eliminadas — Sólo son marcadas como *disponibles*. El espacio de almacenamiento que el archivo eliminado ocupa debe, entonces, ser *sumado* al espacio libre que tiene el volumen. Es sólo cuando se crea un nuevo archivo empleando esa misma entrada, o cuando otro archivo ocupa el espacio físico que ocupaba el que fue eliminado, que el sistema operativo marca *realmente* como desocupados los *clusters* en la tabla de asignación.

Es por esto que desde los primeros días de las PC existen tantas herramientas de recuperación (o *des-borramiento, undeletion*) de archivos: Siempre que no haya sido creado un archivo nuevo que ocupe la entrada de directorio en cuestión, recuperar un archivo es tan simple como volver a ponerle el primer carácter a su nombre.

Este es un ejemplo de un *mecanismo flojo* (en contraposición de los *mecanismos ansiosos*, como los vistos en la sección 6.5.1). Eliminar un archivo requiere de un trabajo mínimo, mismo que es *diferido* al momento de reutilización de la entrada de directorio.

8.1.5 Compresión y desduplicación

Los archivos almacenados en un área dada del disco tienden a presentar patrones comunes. Algunas situaciones ejemplo que llevarían a estos patrones comunes son:

- Dentro del directorio de trabajo de cada uno de los usuarios hay típicamente archivos creados con los mismos programas, compartiendo encabezado, estructura, y ocasionalmente incluso parte importante del contenido.
- Dentro de los directorios de binarios del sistema, habrá muchos archivos ejecutables compartiendo el mismo *formato binario*.
- Es muy común también que un usuario almacene versiones distintas del mismo documento.
- Dentro de un mismo documento, es frecuente que el autor repita en numerosas ocasiones las palabras que describen sus conceptos principales.

Conforme las computadoras aumentaron su poder de cómputo, desde fines de los 1980 se presentaron varios mecanismos que permitían aprovechar las regularidades en los datos almacenados en el disco para comprimir el espacio utilizable en un mismo medio. La compresión típicamente se hace por medio de mecanismos estimativos derivados del análisis del contenido¹², que tienen por resultado un nivel variable de compresión: Con tipos de contenido altamente regulares (como podría ser texto, código fuente, o audio e imágenes *en crudo*), un volumen puede almacenar frecuentemente mucho más del 200% de su espacio real.

Con contenido poco predecible o con muy baja redundancia (como la mayor parte de formatos de imágenes o audio, que incluyen ya una fase de compresión, o empleando cualquier esquema de cifrado) la compresión no ayuda, y sí reduce el rendimiento global del sistema en que es empleada.

Compresión de volumen completo

El primer sistema de archivos que se popularizó fue *Stacker*, comercializado a partir de 1990 por *Stac Electronics*. Stacker operaba bajo MS-DOS, creando un dispositivo de bloques virtual alojado en un disco estándar¹³. Varias implementaciones posteriores de esta misma época se basaron en este mismo principio.

¹¹Este tema será abordado en breve, en la sección 8.2.4, al hablar de las tablas de asignación de archivos.

¹²Uno de los algoritmos más frecuentemente utilizados y fáciles de entender es la *Codificación Huffman*; este y la familia de algoritmos *Lempel-Ziv* sirven de base para prácticamente la totalidad de implementaciones.

¹³Esto significa que, al solicitarle la creación de una unidad comprimida de 30MB dentro del volumen C (disco duro primario), esta aparecería disponible como un volumen adicional. El nuevo volúmen requería de la carga de un *controlador* especial para ser *montado* por el sistema operativo.

Ahora, sumando la variabilidad derivada del enfoque probabilístico al uso del espacio con el ubicarse como una compresión orientada al volumen entero, resulta natural encontrar una de las dificultades resultantes del uso de estas herramientas: Dado que el sistema operativo estructura las operaciones de lectura y escritura por bloques de dimensiones regulares (por ejemplo, el tamaño típico de sector hardware de 512 bytes), al poder estos traducirse a más o menos bloques reales al pasar por una capa de compresión, es posible que el sistema de archivos tenga que reacomodar constantemente la información al *crecer* alguno de los bloques previos. Conforme mayor era el tiempo de uso de una unidad comprimida por Stacker, se notaba más degradación en su rendimiento.

Además, dado que bajo este esquema se empleaba el sistema de archivos estándar, las tablas de directorio y asignación de espacio resultaban también comprimidas. Estas tablas, como ya se ha expuesto, contienen la información fundamental del sistema de archivos; al comprimirlos y reescribirlos constantemente, la probabilidad de que resulten dañadas en alguna falla (eléctrica o lógica) aumenta. Y sí, si bien los discos comprimidos por Stacker y otras herramientas fueron populares principalmente durante la primera mitad de los 1990, conforme aumentó la capacidad de los discos duros fue declinando su utilización.

Compresión archivo por archivo

Dado el éxito del que gozó Stacker en sus primeros años, Microsoft anunció como parte de las características de la versión 6.0 de MS-DOS (publicada en 1993) que incluiría *DoubleSpace*, una tecnología muy similar a la de Stacker. Microsoft incluyó en sus sistemas operativos el soporte para DoubleSpace por siete años, cubriendo las últimas versiones de MS-DOS y las de Windows 95, 98 y Millenium, pero como ya se vio, la compresión de volumen completo presentaba importantes desventajas.

Para el entonces nuevo sistemas de archivos NTFS, Microsoft decidió incluir una característica distinta, más segura y más modular: Mantener el sistema de archivos funcionando de forma normal, sin compresión, y habilitar la compresión *archivo por archivo* de forma transparente al usuario.

Este esquema permite al administrador del sistema elegir, por archivos o carpetas, qué áreas del sistema de archivos desea almacenar comprimidas; esta característica viene como parte de todos los sistemas operativos Windows a partir de la versión XP, liberada en el año 2003.

Si bien la compresión transparente a nivel archivo se muestra mucho más atractiva que la compresión de volumen completo, no es una panacea y es frecuente encontrar en foros en línea la recomendación de deshabilitarla. En primer término, comprimir un archivo implica que un cambio pequeño puede tener un impacto mucho mayor: Modificar un bloque puede implicar que el tamaño final de los datos cambie, lo cual se traduciría a la reescritura del archivo desde ese punto en adelante; esto podría mitigarse insertando espacios para preservar el espacio ya ocupado, pero agrega complejidad al proceso (y abona en contra de la compresión). Los archivos comprimidos son además mucho más sensibles a la corrupción de datos, particularmente en casos de fallo de sistema o de energía: Dado que un cambio menor puede resultar en la necesidad de reescribir al archivo completo, la ventana de tiempo para que se produzca un fallo se incrementa.

En archivos estructurados para permitir el acceso aleatorio, como podrían ser las tablas de bases de datos, la compresión implicaría que los registros no estarán ya alineados al tamaño que el programa gestor espera, lo cual acarrearía necesariamente una penalización en el rendimiento y en la confiabilidad.

Por otro lado, los formatos nativos en que se expresan los datos que típicamente más espacio ocupan en el almacenamiento de los usuarios finales implican ya un alto grado de compresión: Los archivos de fotografías, audio o video están codificados empleando diversos esquemas de compresión aptos para sus particularidades. Y comprimir un archivo

que de suyo está ya comprimido no sólo no reporta ningún beneficio, sino que resulta en desperdicio de trabajo por el esfuerzo invertido en descomprimirlo cada vez que es empleado.

La compresión transparente archivo por archivo tiene innegables ventajas, sin embargo, por las desventajas que implica, no puede tomarse como el modo de operación por omisión.

Desduplicación

Hay una estrategia fundamentalmente distinta para optimizar el uso del espacio de almacenamiento, logrando muy altos niveles de *sobreuso*: Guardar sólo una copia de cada cosa.

Ha habido sistemas implementando distintos tipos de desduplicación desde fines de los 1980, aunque su efectividad era muy limitada y, por tanto, su uso se mantuvo como muy marginal hasta recientemente.

El que se retomara la desduplicación se debe en buena medida a la *consolidación* de servidores ante la adopción a gran escala de mecanismos de virtualización (ver apéndice B, y en particular la sección B.5). Dado que un mismo servidor puede estar alojando a decenas o centenas de *máquinas virtuales*, muchas de ellas con el mismo sistema operativo y programas base, los mismos archivos se repiten muchas veces; si el sistema de archivos puede determinar que determinado archivo o bloque está ya almacenado, podría almacenarse sólamente una copia.

La principal diferencia entre la desduplicación y las *ligas duras* mencionados en la sección 7.3.1 es que, en caso de que cualquiera de estos archivos (o bloques) sea modificado, el sistema de archivos tomará el espacio necesario para representar estos cambios y evitará que esto afecte a los demás archivos. Además, si dos archivos inicialmente distintos se hacen iguales, se liberará el espacio empleado por uno de ellos de forma automática.

Para identificar qué datos están duplicados, el mecanismo más utilizado es calcular el *hash criptográfico* de los datos¹⁴; este mecanismo permite una búsqueda rápida y confiable de coincidencias, sea a nivel archivo o a nivel bloque.

La desduplicación puede realizarse *en línea* o *fuera de línea* — Esto es, analizar los datos buscando duplicidades al momento que estos llegan al sistema, o, dado que es una tarea intensiva tanto en uso de procesador como de entrada/salida, realizarla como una tarea posterior de mantenimiento, en un momento de menor ocupación del sistema.

Desduplicar a nivel archivo es mucho más ligero para el sistema que hacerlo a nivel bloque, pero hacerlo a nivel bloque lleva típicamente a una optimización del uso de espacio mucho mayor.

Al día de hoy, el principal sistema de archivos que implementa desduplicación es ZFS¹⁵, desarrollado por Sun Microsystems (hoy Oracle). En Linux, esta característica forma parte de BTRFS, aunque no ha alcanzado los niveles de estabilidad como para recomendarse para su uso en entornos de producción.

En Windows, esta funcionalidad se conoce como *Single Instance Storage (Almacenamiento de Instancia Única)*. Esta característica apareció a nivel de archivo, implementada en espacio de usuario, como una de las características del servidor de correo *Exchange Server* entre los años 2000 y 2010. A partir de Windows Server 2003, la funcionalidad de desduplicación existe para NTFS, pero su uso es poco frecuente.

El uso de desduplicación, particularmente cuando se efectúa a nivel bloques, tiene un alto costo en memoria: Para mantener buenos niveles de rendimiento, la tabla que relaciona el hash de datos con el sector en el cual está almacenado debe mantenerse en

¹⁴Por ejemplo, empleando el algoritmo SHA-256, el cual brinda una probabilidad de *colisión* de 1 en 2^{256} , suficientemente confiable para que la pérdida de datos sea mucho menos probable que la falla del disco.

¹⁵Las características básicas de ZFS serán presentadas en la sección C.3.2

memoria. En el caso de la implementación de ZFS en FreeBSD, la documentación sugiere dedicar 5GB de memoria por cada TB de almacenamiento (0.5% de la capacidad total).

8.2 Esquemas de asignación de espacio

Hasta este punto, la presentación de la *entrada de directorio* se ha limitado a indicar que ésta apunta al punto donde inicia el espacio empleado por el archivo. No se ha detallado en cómo se implementa la asignación y administración de dicho espacio. En esta sección se hará un breve repaso de los tres principales mecanismos, para después de ésta explicar cómo es la implementación de FAT, abordando sus principales virtudes y debilidades.

8.2.1 Asignación contigua

Los primeros sistemas de archivos en disco empleaban un esquema de *asignación contigua*. Para implementar un sistema de archivos de este tipo, no haría falta el contar con una *tabla de asignación de archivos*: Bastaría con la información que forma parte del directorio de FAT — La extensión del archivo y la dirección de su primer *cluster*.

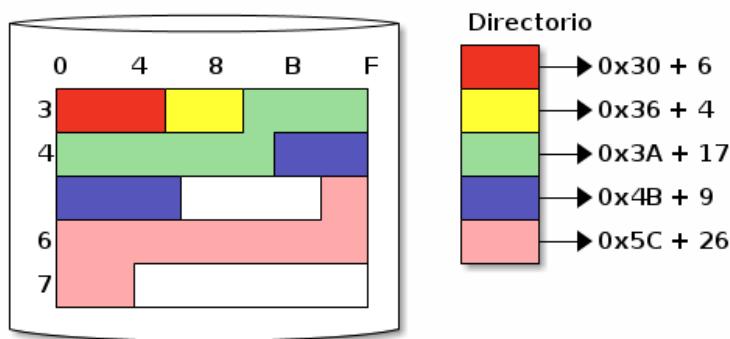


Figura 8.3: Asignación contigua de archivos: Directorio con inicio y longitud

La principal ventaja de este mecanismo de asignación, claro está, es la simplicidad de su implementación. Brinda además la mejor velocidad de transferencia del archivo, dado que al estar cada uno de los archivos en espacio contiguo en el disco, el movimiento de cabezas se mantiene al mínimo. Sin embargo, este mecanismo se vuelve sumamente inconveniente en medios que soporten lectura y escritura: Es muy sensible a la *fragmentación externa*; si un archivo requiere crecer, debe ser movido íntegramente a un bloque más grande (lo cual toma demasiado tiempo), y el espacio que libera un archivo en caso de reducirse su necesidad de espacio queda *atrapado* entre bloques asignados; podemos tener mucho más espacio disponible que el que podamos asignar a un nuevo archivo.

Los esquemas de asignación contigua se emplean hoy en día principalmente en sistemas de archivo de sólo lectura — Por ejemplo, lo emplea el sistema principal que utilizan los CD-ROMs, el ISO-9660, pensado para aprovechar al máximo un espacio que, una vez grabado, sólo podrá abrirse en modo de sólo lectura. Esto explica por qué, a diferencia de como ocurre en cualquier otro medio de almacenamiento, al *quemar* un CD-ROM es necesario preparar primero una *imagen* en la que los archivos ocupen sus posiciones definitivas, y esta imagen debe grabarse al disco en una sola operación.

8.2.2 Asignación ligada

Un enfoque completamente distinto sería el de *asignación ligada*. Este esquema brinda mucho mayor flexibilidad al usuario, sacrificando la simplicidad y la velocidad: Cada

entrada en el directorio apunta a un primer *grupo* de sectores (o *cluster*), y éste contiene un apuntador que indica cuál es el siguiente.

Para hacer esto, hay dos mecanismos: El primero, reservar un espacio al final de cada *cluster* para guardar el apuntador, y el segundo, crear una tabla independiente, que guarde únicamente los apuntadores.

En el primer caso, si se manejan *clusters* de 2048 bytes, y se reservan los 4 bytes (32 bits) finales de cada uno, el resultado sería de gran incomodidad al programador: Frequentemente, los programadores buscan alinear sus operaciones con las fronteras de las estructuras subyacentes, para optimizar los accesos (por ejemplo, evitar que un sólo registro de base de datos requiera ser leído de dos distintos bloques en disco). El programador tendría que diseñar sus estructuras para ajustarse a la poco ortodoxa cantidad de 2044 bytes.

Y más allá de esta inconveniencia, guardar los apuntadores al final de cada *cluster* hace mucho más lento el manejo de todos los archivos: Al no tener en una sola ubicación la relación de *clusters* que conforman a un archivo, todas las transferencias se convierten en *secuenciales*: Para llegar directamente a determinado bloque del archivo, habrá que atravesar todos los bloques previos para encontrar su ubicación.

Particularmente por este segundo punto es mucho más común el empleo de una *tabla de asignación de archivos* — Y precisamente así es como opera FAT (de hecho, esta tabla es la que le da su nombre). La tabla de asignación es un mapa de los *clusters*, representando a cada uno por el espacio necesario para guardar un apuntador.

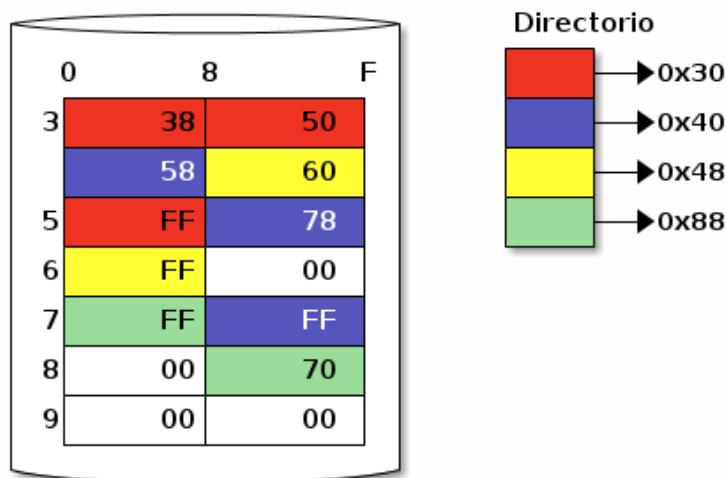


Figura 8.4: Asignación ligada de archivos: Directorio con apuntador sólo al primer *cluster*

La principal ventaja del empleo de asignación ligada es que desaparece la *fragmentación interna*.¹⁶ Al ya no requerir la *pre-asignación* de un espacio contiguo, cada uno de los archivos puede crecer o reducirse según sea necesario.

Ahora, la *asignación ligada* no sólo resulta más lenta que la contigua, sino que presenta una mayor *sobrecarga administrativa*: El espacio desperdiciado para almacenar los apuntadores típicamente es cercano al 1% del disponible en el medio.

Este esquema también presenta *fragilidad de metadatos*: Si alguno de los apuntadores se

¹⁶Con *fragmentación interna* se hace aquí referencia al concepto presentado en la sección 6.2.1. El fenómeno generalmente conocido como *fragmentación* se refiere a la necesidad de *compactación*; es muy distinto, y sí se presenta bajo este esquema: Cada archivo se separa en pequeños *fragmentos* que pueden terminar espaciados por todo el disco, impactando fuertemente en el rendimiento del sistema

pierde o corrompe, lleva a que se pierda el archivo *completo* desde ese punto y hasta su final (y abre la puerta a la corrupción de otro archivo, si el apuntador queda apuntando hacia un bloque empleado por éste; el tema de fallos y recuperación bajo estos esquemas es abordado en la sección 8.3).

Hay dos mecanismos de mitigación para este problema: El empleado por FAT es guardar una (o, bajo FAT12, dos) copias adicionales de la tabla de asignación, entre las cuales que el sistema puede verificar si se mantengan consistentes y buscar corregirlas en caso contrario. Por otro lado, puede manejarse una estructura de *lista doblemente ligada* (en vez de una *lista ligada* sencilla) en que cada elemento apunte tanto al siguiente como al anterior, con lo cual, en caso de detectarse una inconsistencia en la información, esta pueda ser recorrida de *atrás hacia adelante* para confirmar los datos correctos. En ambos casos, sin embargo, la sobrecarga administrativa se duplica.

8.2.3 Asignación indexada

El tercer modelo es la *asignación indexada*, el mecanismo empleado por casi todos los sistemas de archivos modernos. En este esquema, se crea una estructura intermedia entre el directorio y los datos, única para cada archivo: el *i-nodo* (o *nodo de información*). Cada i-nodo guarda los metadatos y la lista de bloques del archivo, reduciendo la probabilidad de que se presente la *corrupción de apuntadores* mencionada en la sección anterior.

La sobrecarga administrativa bajo este esquema potencialmente es mucho mayor: Al asignar el i-nodo, éste se crea ocupando como mínimo un *cluster* completo. En el caso de un archivo pequeño, que quepa en un sólo *cluster*, esto representa un desperdicio del 100% de espacio (un *cluster* para el i-nodo y otro para los datos);¹⁷ para archivos más grandes, la sobrecarga relativa disminuye, pero se mantiene siempre superior a la de la asignación ligada.

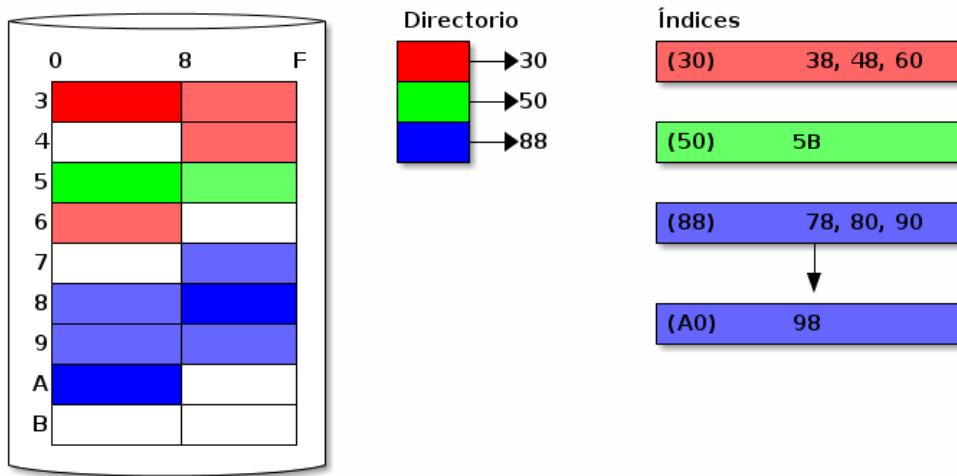


Figura 8.5: Asignación indexada de archivos: Directorio con apuntador al i-nodo (llevado a un i-nodo de tamaño extremadamente ineficiente)

Un esquema de asignación indexada nos da una mayor eficiencia de caché que la asignación ligada: Si bien en dicho caso es común guardar copia de la tabla de asignación en memoria para mayor agilidad, con la asignación indexada bastará hacer caché únicamente

¹⁷Algunos sistemas de archivos, como ReiserFS, BTRFS o UFS, presentan esquemas de asignación *sub-cluster*. Estos denominan *colas* (*tails*) a los archivos muy pequeños, y pueden ubicarlos ya sea dentro de su mismo i-nodo o compartiendo un mismo cluster con un *desplazamiento* dentro de éste. Esta práctica no ha sido adoptada por sistemas de archivos de uso mayoritario por su complejidad relativa.

de la información importante, esto es, únicamente de los archivos que se emplean en un momento dado. El *mapa de asignación* para los archivos y directorios que no hayan sido empleados recientemente no requerirán estar en memoria.

Claro está, mientras que en los esquemas anteriores la tabla central de asignación de archivos puede emplearse directamente como el *bitmap* del volumen, en los sistemas de archivos de asignación indexada se vuelve necesario contar con un *bitmap* independiente — Pero al sólo requerir representar si cada *cluster* está vacío u ocupado (y ya no apuntar al siguiente), resulta de mucho menor tamaño.

Ahora, ¿qué pasa cuando la lista de *clusters* de un archivo no cabe en un i-nodo? Este ejemplo se ilustra en el tercer archivo de la figura 8.6: En el caso ilustrado, cada i-nodo puede guardar únicamente tres apuntadores.¹⁸ Al tener un archivo con cuatro *clusters*, se hace necesario extender al i-nodo con una lista adicional. La implementación más común de este mecanismo es que, dependiendo del tamaño del archivo, se empleen apuntadores con los niveles de indirección que *vayan haciendo falta*.

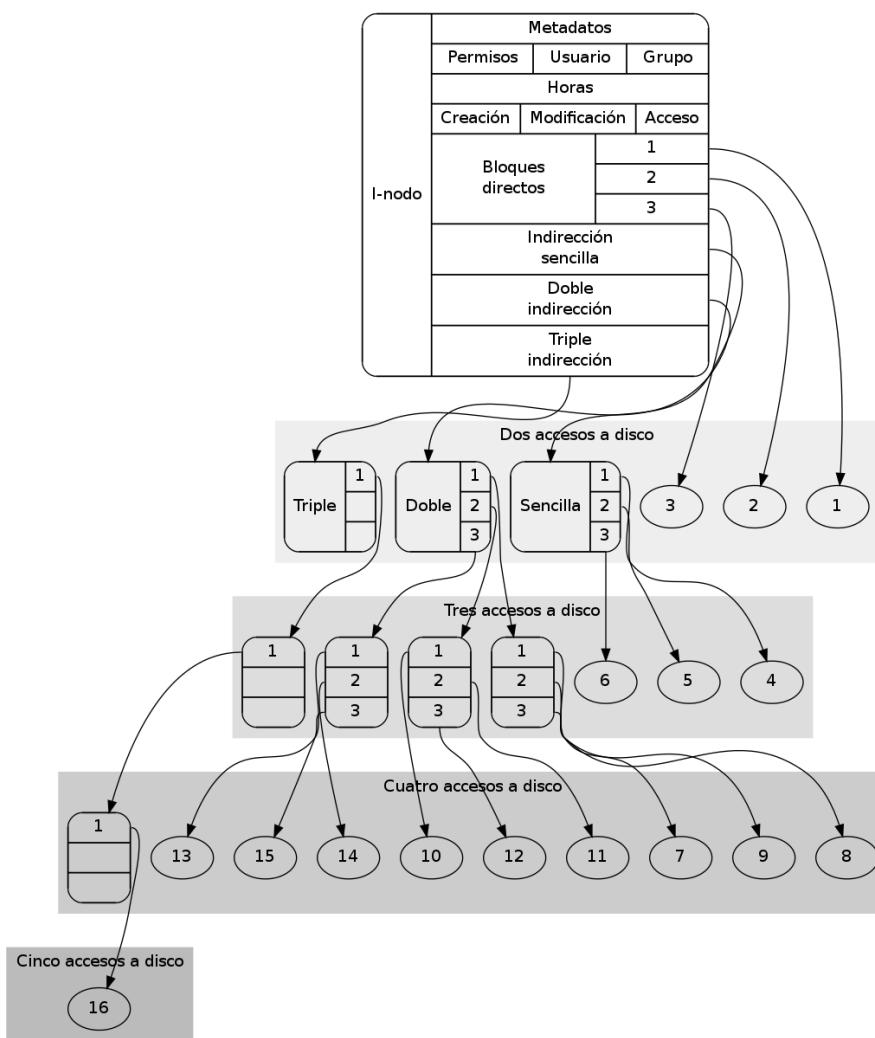


Figura 8.6: Estructura típica de un i-nodo en Unix, mostrando además el número de accesos a disco necesarios para llegar a cada cluster (con sólo tres cluster por lista)

¿Qué tan grande sería el archivo máximo direccionable bajo este esquema y únicamente tres indirecciones? Suponiendo magnitudes que típicas hoy en día (*clusters* de 4KB y direcciones de 32 bits), en un *cluster* vacío caben 128 apuntadores ($\frac{4096}{32}$). Si los metadatos

¹⁸Esto resulta un límite demasiado bajo, y fue elegido meramente para ilustrar el presente punto.

ocupan 224 bytes en el i-nodo, dejando espacio para 100 apuntadores:

- Un archivo de hasta $(100 - 3) \times 4KB = 388KB$ puede ser referido por completo directamente en el i-nodo, y es necesario un sólo acceso a disco para obtener su lista de *clusters*.
- Un archivo de hasta $(100 - 3 + 128) \times 4KB = 900KB$ puede representarse con el bloque de indirección sencilla, y obtener su lista de *clusters* significa dos accesos a disco adicionales.
- Con el bloque de doble indirección, puede hacerse referencia a archivos mucho más grandes:
 $(100 - 3 + 128 + (128 \times 128)) \times 4KB = 66436KB \approx 65MB$
 Sin embargo, aquí ya llama la atención otro importante punto: Para acceder a estos 65MB, es necesario que realizar hasta 131 accesos a disco. A partir de este punto, resulta importante que el sistema operativo asigne *clusters* cercanos para los metadatos (y, de ser posible, para los datos), pues la diferencia en tiempo de acceso puede ser muy grande.
- Empleando triple indirección, se puede llegar hasta:
 $(100 - 3 + 128 + (128 \times 128) + (128 \times 128 \times 128)) \times 2KB = 8455044KB \approx 8GB$
 Esto es ya más de lo que puede representarse en un sistema de 32 bits. La cantidad de bloques a leerse, sin embargo, para encontrar todos los *clusters* significarían hasta 16516 accesos a disco (en el peor de los casos).

8.2.4 Las tablas en FAT

Volviendo al caso que se presenta como ejemplo, el sistema de archivos FAT: en este sistema, cada entrada del directorio apunta al primer *cluster* que ocupa cada uno de los archivos, y se emplea un esquema de asignación ligada. El directorio tiene también un campo indicando la *longitud total* del archivo, pero esto no es empleado para leer la información, sino para poderla presentar más ágilmente al usuario (y para que el sistema operativo sepa dónde indicar *fin de archivo* al leer el último *cluster* que compone a determinado archivo).

La estructura fundamental de este sistema de archivos es la tabla de asignación de archivos (*File Allocation Table*) — Tanto que de ella toma su nombre FAT.

Cada entrada de la FAT mide lo que la longitud correspondiente a su versión (12, 16 o 32 bits), y puede tener cualquiera de los valores descritos en el cuadro 8.2.

Cuadro 8.2: Valores especiales que puede almacenar FAT; cualquier otro valor indica la dirección del siguiente *cluster* que forma parte del archivo al cual pertenece el registro en cuestión.

FAT12	FAT16	FAT32	Significado
0x000	0x0000	0x00000000	Disponible, puede ser asignado
0xFF7	0xFFFF7	0xFFFFFFFF7	Cluster dañado, no debe utilizarse
0xFFF	0xFFFF	0xFFFFFFFFFF	Último cluster de un archivo

Llama la atención que haya un valor especial para indicar que un *cluster* tiene sectores dañados. Esto remite de vuelta al momento histórico de la creación de la familia FAT: Siendo el medio predominante de almacenamiento el disco flexible, los errores en la superficie eran mucho más frecuentes de lo que lo son hoy en día.

Una característica que puede llamar la atención de FAT es que parecería permitir la fragmentación de archivos *por diseño*: Dado que el descriptor de cada *cluster* debe apuntar al siguiente, puede asumirse que el *caso común* es que los *clusters* no ocuparán contiguos en el disco. Claro está, la tabla puede apuntar a varios *clusters* adyacentes, pero el sistema

de archivos mismo no hace nada para que así sea.

The figure displays two tables representing the File Allocation Table (FAT) for 15 clusters. Each table has columns labeled 0 through F. The first table shows a standard assignment where each cluster is a different color. The second table shows a corrupted assignment where some clusters are red (bad clusters) or yellow (10 cluster entry).

0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
F8 FF	FF 7F	03 00	04 00	05 00	06 00	07 00	08 00	yy							
09 00	0A 00	0B 00	0C 00	0D 00	0E 00	0F 00	10 00								

Legend for the first table:

- 0x00 -> media type
- 0x01 -> partition state
- 0x02 -> cluster 2
- 0x03 -> cluster 3
- 0x04 -> cluster 4
- 0x05 -> cluster 5
- 0x06 -> cluster 6
- 0x07 -> cluster 7
- 0x08 -> cluster 8
- 0x09 -> cluster 9
- 0x0A -> cluster 10
- 0x0B -> cluster 11
- 0x0C -> cluster 12
- 0x0D -> cluster 13
- 0x0E -> cluster 14
- 0x0F -> cluster 15

0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
F8 FF	FF 7F	06 00	04 00	05 00	FF FF	07 00	08 00	yy							
09 00	14 00	0B 00	0C 00	0D 00	0E 00	0F 00	10 00								
11 00	12 00	13 00	FF FF	FF FF	00 00	F7 FF	F7 FF								

Legend for the second table:

- 0x00 -> media type
- 0x01 -> partition state
- 0x02 -> 1st of a 6 cluster entry
- 0x03 -> 3 cluster entry
- 0x04 -> 2 to 5 of a 6 cluster entry
- 0x05 -> 10 cluster entry
- 0x06 -> last cluster of 6 entry
- 0x07 -> available cluster
- 0x08 -> 2 bad clusters

Figura 8.7: Ejemplo de entradas en la tabla de asignación de archivos (Imagen: Peter Clark, ver otros recursos)

En la sección 8.1.4, al presentar el formato del directorio de FAT, se mencionó que los subdirectorios son en realidad archivos de un tipo especial: una suerte de archivos estructurados (ver sección 7.2.5), gestionados por el sistema operativo. Lo único que distingue a un directorio de un archivo normal es que, en la entrada que lo describe en su directorio padre, el doceavo byte de la entrada (que indica los *atributos del archivo*, ver figura 8.1 y cuadro 7.1) tiene activado el bit 4.

Un directorio es almacenado en disco *exactamente* como cualquier otro archivo. Si se le asigna únicamente un *cluster*, y el tamaño del *cluster* es pequeño (2KB), podrá almacenar sólo 64 entradas ($\frac{2048}{32}$) y cada *cluster* adicional le dará 64 entradas más. Y como tal, está sujeto también a la fragmentación: Conforme se agregan entradas al directorio, éste crece. Llegado el momento, requiere *clusters* adicionales. Y si un directorio termina disperso por todo el disco, resultará –como cualquier otro archivo– más lento leerlo y trabajar con él. Siempre que se abra un archivo dentro de un directorio grande, o que se lo recorra para abrir algún archivo en algún subdirectorio suyo, el sistema tendrá que buscar todos sus fragmentos a lo largo del disco.

Ante estos dos aspectos, no puede perderse de vista la edad que tiene FAT. Otros sistemas de archivos más modernos han resuelto este problema a través de los *grupos de asignación*: Los directorios del sistema son *esparcidos* a lo largo del volumen, y se intenta ubicar a los archivos cerca de los directorios desde donde son referidos¹⁹. Esto tiene por consecuencia que los archivos que presentan *cercanía temática* (que pertenecen al mismo usuario o proyecto, o que por alguna razón están en la misma parte de la jerarquía del sistema) quedan ubicados en disco cerca unos de otros (y cerca de sus directorios). Y dado que es probable que sean empleados aproximadamente al mismo tiempo, esto reduce las distancias que recorrerán las cabezas. Además, al esparcir los archivos, se distribuye también mejor el espacio libre, por lo cual el impacto de los cambios de tamaño de un

¹⁹Claro está, en el caso de los archivos que están como *ligas duras* desde varios directorios, pueden ubicarse sólo cerca de uno de ellos

archivo en lo relativo a la fragmentación se limita a los que forman parte del mismo bloque de asignación.

Los sistemas de archivos que están estructurados siguiendo esta lógica de grupos de asignación no evitan la fragmentación, pero sí la mayor parte de sus consecuencias negativas. Para mantener este esquema operando confiablemente, eso sí, requieren de mantener disponibilidad de espacio — Al presentarse saturación, esta estrategia pierde efectividad. Para evitar que esto ocurra, es muy frecuente en los sistemas Unix que haya un cierto porcentaje (típicamente cercano al 5%) del disco que esté disponible únicamente para el administrador — En caso de que el sistema de archivos pase del 95%, los usuarios no podrán escribir ya a él, pero el administrador puede efectuar tareas de mantenimiento para volver a un rango operacional.

8.3 Fallos y recuperación

El sistema de archivos FAT es *relativamente frágil*: No es difícil que se presente una situación de *corrupción de metadatos*, y muy particularmente, de la estructura de las tablas de asignación. Los usuarios de sistemas basados en FAT en Windows sin duda conocen a los programas `CHKDSK` y `SCANDISK`, dos programas que implementan la misma funcionalidad base, y difieren principalmente en su interfaz al usuario: `CHKDSK` existe desde los primeros años de MS-DOS, y está pensado para su uso interactivo en línea de comando; `SCANDISK` se ejecuta desde el entorno gráfico, y presenta la particularidad de que no requiere (aunque sí recomienda fuertemente) *acceso exclusivo* al sistema de archivos mientras se ejecuta.

¿Cómo es el funcionamiento de estos programas?

A lo largo de la vida de un sistema de archivos, conforme los archivos se van asignando y liberando, van cambiando su tamaño, y conforme el sistema de archivos se monta y des-monta, pueden ir apareciendo *inconsistencias* en su estructura. En los sistemas tipo FAT, las principales inconsistencias²⁰ son:

Archivos cruzados En inglés, *cross-linked file*. Recuérdese que la entrada en el directorio de un archivo incluye un apuntador al primer *cluster* de una *cadena*. Cada cadena debe ser única, esto es, ningún *cluster* debe pertenecer a más de un archivo. Si dos archivos incluyen al mismo *cluster*, esto indica una inconsistencia, y la única forma de resolverla es *truncar* a uno de los archivos en el punto inmediato anterior a este cruce.

Cadenas perdidas o huérfanas En inglés, *lost clusters*. Cuando hay espacio marcado como ocupado en la tabla de asignación, pero no hay ninguna entrada de directorio haciendo referencia a ella, el espacio está efectivamente bloqueado y, desde la perspectiva del usuario, inutilizado; además, estas cadenas pueden ser un archivo que el usuario aún requiera.

Este problema resultó tan frecuente en versiones históricas de Unix que al día de hoy es muy común tener un directorio llamado `lost+found` en la raíz de todos los sistemas de archivos, donde `fsck` (el equivalente en Unix de `CHKDSK`) creaba ligas a los archivos perdidos por corrupción de metadatos.

Cada sistema de archivos podrá presentar un distinto conjunto de inconsistencias, dependiendo de sus estructuras básicas y de la manera en que cada sistema operativo las maneja.

En la década de los 1980 comenzaron a entrar a mercado los *controladores de disco inteligentes*, y en menos de diez años dominaban ya el mercado. Estos controladores, con buses tan disímiles como SCSI, IDE, o los más modernos, SAS y SATA, introdujeron

²⁰Que no las únicas. Estas y otras más están brevemente descritas en la página de manual de `dosfsck` (ver sección 8.4).

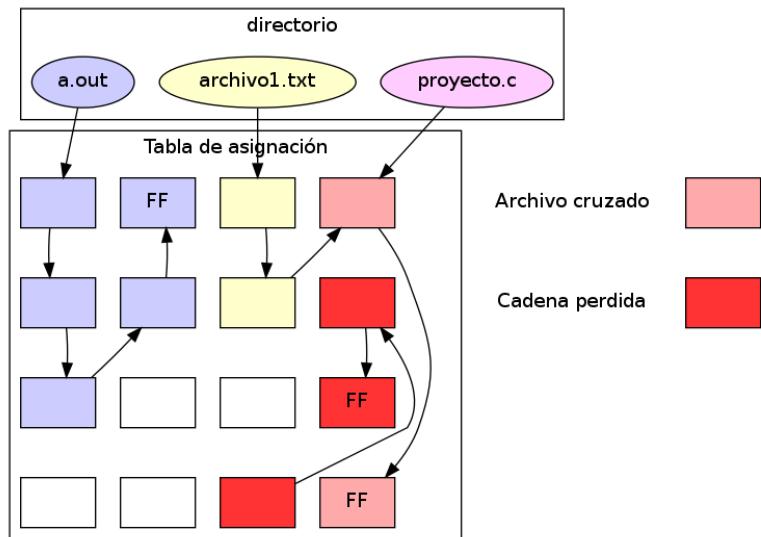


Figura 8.8: Inconsistencias en un sistema de archivos tipo FAT

muchos cambios que fueron disociando cada vez más al sistema operativo de la gestión física directa de los dispositivos; en el apéndice C se presenta a mayor profundidad lo que esto ha significado para el desarrollo de sistemas de archivos y algoritmos relacionados. Sin embargo, para el tema en discusión, los *controladores inteligentes* resultan relevantes porque, si bien antes el sistema operativo podía determinar con toda certeza si una operación se había realizado o no, hoy en día los controladores dan un *acuse de recibo* a la información en el momento en que la colocan en el caché incorporado del dispositivo — En caso de un fallo de corriente, esta información puede no haber sido escrita por completo al disco.

Es importante recordar que las operaciones con los metadatos que conforman al sistema de archivos no son atómicas. Por poner un ejemplo, crear un archivo en un volumen FAT requiere:

1. Encontrar una lista de *clusters* disponibles suficiente para almacenar la información que conformará al archivo
2. Encontrar el siguiente espacio disponible en el directorio
3. Marcar en la tabla de asignación la secuencia de *clusters* que ocupará el archivo
4. Crear en el espacio encontrado una entrada con el nombre del archivo, apuntando al primero de sus *clusters*
5. Guardar los datos del archivo en cuestión en los *clusters* correspondientes

Cualquier fallo que se presente después del tercer paso (cuando hacemos la primera modificación) tendrá como consecuencia que el archivo resulte corrupto, y muy probablemente que el sistema de archivos todo *presente inconsistencias o esté en un estado inconsistente*.

8.3.1 Datos y metadatos

En el ejemplo recién presentado, el sistema de archivos estará consistente siempre que se terminen los pasos 3 y 4 — La consistencia del sistema de archivos es independiente de la validez de los datos del archivo. Lo que busca el sistema de archivos, más que asegurar la integridad de los *datos* de uno de los archivos, es asegurar la de los *metadatos*: Los datos que describen la estructura.

En caso de que un usuario desconecte una unidad a media operación, es muy probable que se presentará pérdida de información, pero el sistema de archivos debe buscar no presentar ningún problema que ponga en riesgo *operaciones posteriores* o *archivos no relacionados*.

nados. La corrupción y recuperación de datos en archivos corruptos y truncados, si bien es también de gran importancia, cae más bien en el terreno de las aplicaciones del usuario.

8.3.2 Verificación de la integridad

Cada sistema operativo incluye programas para realizar verificación (y, en su caso, corrección) de la integridad de sus sistemas de archivo. En el caso de MS-DOS y Windows, como ya se vio, estos programas son `CHKDSK` y `SCANDISK`; en los sistemas Unix, el programa general se llama `fsck`, y típicamente emplea a asistentes según el tipo de sistema a revisar — `fsck.vfat`, `fsck.ext2`, etc.

Estos programas hacen un *barrido* del sistema de archivos, buscando evidencias de inconsistencia. Esto lo hacen, en líneas generales:

- Siguiendo todas las cadenas de *clusters* de archivos o tablas de i-nodos, y verificando que no haya archivos cruzados (compartiendo erróneamente bloques)
- Verificando que todas las cadenas de *clusters*, así como todos los directorios, sean alcanzables y sigan una estructura válida
- Recalculando la correspondencia entre las estructuras encontradas y los diferentes bitmaps y totales de espacio vacío

Estas operaciones son siempre procesos intensivos y complejos. Como requieren una revisión profunda del volumen entero, es frecuente que duren entre decenas de minutos y horas. Además, para poder llevar a cabo su tarea deben ejecutarse teniendo acceso exclusivo al volumen a revisar, lo cual típicamente significa colocar al sistema completo en modo de mantenimiento.

Dado el elevado costo que tiene verificar el volumen entero, en la década de 1990 surgieron varios esquemas orientados a evitar la necesidad de invocar a estos programas de verificación: Las *actualizaciones suaves*, los *sistemas de archivos con bitácora*, y los *sistemas de archivos estructurados en bitácora*.

8.3.3 Actualizaciones suaves (*soft updates*)

Este esquema aparentemente es el más simple de los que presentaremos, pero su implementación resultó mucho más compleja de lo inicialmente estimado, y en buena medida por esta causa hoy en día no ha sido empleado más ampliamente. La idea básica detrás de este esquema es estructurar el sistema de archivos de una forma más simple y organizar las escrituras al mismo de modo que el estado resultante *no pueda* ser inconsistente, ni siquiera en caso de fallo, y de exigir que todas las operaciones de actualización de metadatos se realicen de forma *síncrona*.²¹

Ante la imposibilidad de tener un sistema *siempre consistente*, esta exigencia se relajó para permitir inconsistencias *no destructivas*: Pueden presentarse *cadenas perdidas*, dado que esto no pone en riesgo a ningún archivo, sólo disminuye el espacio total disponible.

Esto, aunado a una reestructuración del programa de verificación (`fsck`) como una tarea *ejecutable en el fondo*²² y en una tarea de *recolector de basura*, que no requiere intervención humana (dado que no pueden presentarse inconsistencias destructivas), permite que un sistema de archivos que no fue *limpiamente desmontado* pueda ser montado y utilizado de inmediato, sin peligro de pérdida de información o de corrupción.

Al requerir que todas las operaciones sean síncronas, parecería que el rendimiento global del sistema de archivos tendría que verse afectado, pero por ciertos patrones de acceso muy frecuentes, resulta incluso beneficioso. Al mantenerse un ordenamiento lógico

²¹Esto es, no se le reporta éxito en alguna operación de archivos al usuario sino hasta que ésta es completada y grabada a disco.

²²Una tarea que no requiere de intervención manual por parte del operador, y se efectúa de forma automática como parte de las tareas de mantenimiento del sistema.

entre las dependencias de todas las operaciones pendientes, el sistema operativo puede *combinar* a muchas de estas y reducir de forma global las escrituras a disco.

A modo de ejemplos: si varios archivos son creados en el mismo directorio de forma consecutiva, cada uno de ellos a través de una llamada `open()` independiente, el sistema operativo combinará a todos estos accesos en uno sólo, reduciendo el número de llamadas. Por otro lado, un patrón frecuente en sistemas Unix es, al crear un archivo de uso temporal, solicitar al sistema la creación de un archivo, abrir el archivo recién creado, y ya teniendo al descriptor de archivo, eliminarlo — En este caso, con estas tres operaciones seguidas, *soft updates* podría ahorrarse por completo la escritura a disco.

Esta estrategia se vio impactada por los controladores inteligentes: Si un disco está sometido a carga intensa, no hay garantía para el sistema operativo del orden que seguirán *en verdad* sus solicitudes, que se *forman* en el caché propio del disco. Dado que las actualizaciones suaves dependen tan profundamente de confiar en el ordenamiento, esto rompe por completo la confiabilidad del proceso.

Las actualizaciones suaves fueron implementadas hacia 2002 en el sistema operativo FreeBSD, y fueron adoptadas por los principales sistemas de la familia *BSD, aunque NetBSD lo retiró en 2012, prefiriendo el empleo de sistemas con bitácora — Muy probablemente, la lógica detrás de esta decisión sea la cantidad de sistemas que emplean esta segunda estrategia que se abordará a continuación, y lo complejo de mantener dentro del núcleo dos estrategias tan distintas.

8.3.4 Sistemas de archivo con bitácora (*journaling file systems*)

Este esquema tiene su origen en el ámbito de las bases de datos distribuídas. Consiste en separar un área del volumen y dedicarla a llevar una bitácora con todas las *transacciones* de metadatos.²³ Una *transacción* es un conjunto de operaciones que deben aparecer como atómicas.

La bitácora es generalmente implementada como una *lista ligada circular*, con un apuntador que indica cuál fue la última operación realizada exitosamente. Periódicamente, o cuando la carga de transferencia de datos disminuye, el sistema verifica qué operaciones quedaron pendientes, y *avanza* sobre la bitácora, marcando cada una de las transacciones conforme las realiza.

En caso de tener que recuperarse de una condición de fallo, el sistema operativo sólo tiene que leer la bitácora, encontrar cuál fue la última operación efectuada, y aplicar las restantes.

Una restricción de este esquema es que las transacciones guardadas en la bitácora deben ser *idempotentes* — Esto es, si una de ellas es efectuada dos veces, el efecto debe ser exactamente el mismo que si hubiera sido efectuada una sola vez. Por poner un ejemplo, no sería válido que una transacción indicara “Aregar al directorio *x* un archivo llamado *y*”, dado que si la falla se produce después de procesar esta transacción pero antes de avanzar al apuntador de la bitácora, el directorio *x* quedaría con dos archivos *y* — Una situación inconsistente. En todo caso, tendríamos que indicar “registrar al archivo *y* en la posición *z* del directorio *x*”; de esta manera, incluso si el archivo ya había sido registrado, puede volverse a registrar sin peligro.

Este esquema es el más implementado hoy en día, y está presente en casi todos los sistemas de archivos modernos. Si bien con un sistema con bitácora no hace falta verificar el sistema de archivos completo tras una detención abrupta, esta no nos exime de que, de tiempo en tiempo, el sistema de archivos sea verificado: Es altamente recomendado hacer

²³Existen implementaciones que registran también los datos en la bitácora, pero tanto por el tamaño que ésta requiere como por el impacto en velocidad que significa, son poco utilizadas. La sección 8.3.5 presenta una idea que elabora sobre una bitácora que almacena tanto datos como metadatos.

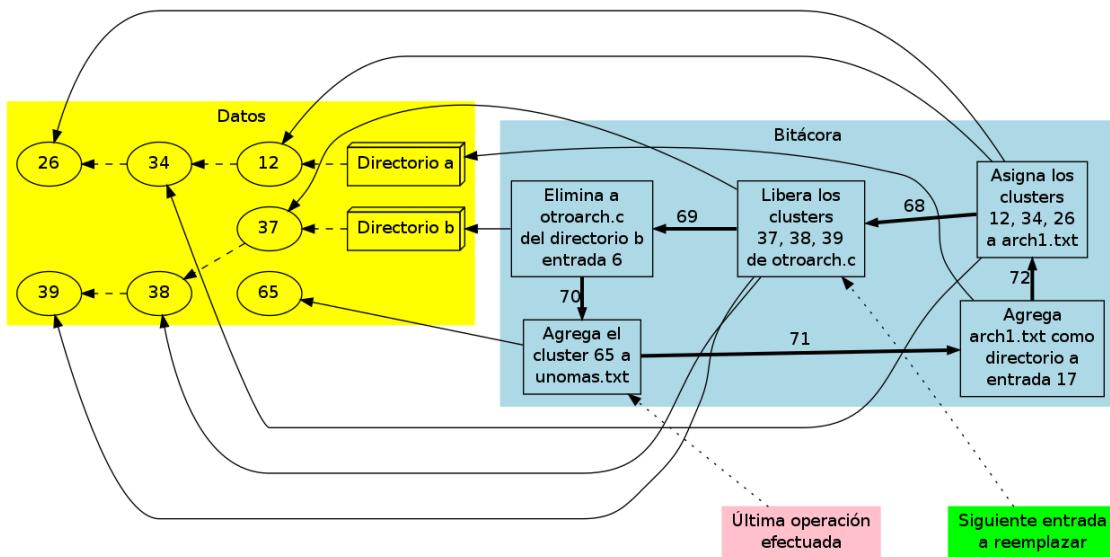


Figura 8.9: Sistema de archivos con bitácora

una verificación periódica en caso de presentarse alguna corrupción, sea por algún bug en la implementación, fallos en el medio físico, o factores similarmente poco frecuentes.

La mayor parte de sistemas de archivos incluyen contadores de *cantidad de montajes* y de *fecha del último montaje*, que permiten que el sistema operativo determine, de forma automática, si corresponde hacer una verificación preventiva.

8.3.5 Sistemas de archivos estructurados en bitácora (*log-structured file systems*)

Si se lleva el concepto del sistema de archivos con bitácora a su límite, y se designa a la *totalidad* del espacio en el volumen como la bitácora, el resultado es un sistema de archivos *estructurado en bitácora*. Obviamente, este tipo de sistemas de archivos presenta una organización completamente diferente de los sistemas de archivo tradicionales.

Las ideas básicas detrás de la primera implementación de un sistema de archivos de esta naturaleza (Ousterhout y Rosenblum, 1992) apuntan al empleo agresivo de caché de gran capacidad, y con un fuerte mecanismo de *recolección de basura*, reacomodando la información que esté más cerca de la *cola* de la bitácora (y liberando toda aquella que resulte redundante).

Este tipo de sistemas de archivos facilita las escrituras, haciéndolas siempre secuenciales, y buscan –a través del empleo del caché ya mencionado– evitar que las cabezas tengan que desplazarse con demasiada frecuencia para recuperar fragmentos de archivos.

Una consecuencia directa de esto es que los sistemas de archivos estructurados en bitácora fueron los primeros en ofrecer *fotografías (snapshots)* del sistema de archivos: Es posible apuntar a un momento en el tiempo, y –con el sistema de archivos aún en operación– montar una copia de sólo lectura con la información del sistema de archivos *completa* (incluyendo los datos de los archivos).

Los sistemas de archivo estructurados en bitácora, sin embargo, no están optimizados para cualquier carga de trabajo. Por ejemplo, una base de datos relacional, en que cada uno de los registros es típicamente actualizados de forma independiente de los demás, y ocupan apenas fracciones de un bloque, resultaría tremadamente ineficiente. La implementación referencial de Ousterhout y Rosenblum fue parte de los sistemas *BSD, pero dada su tendencia a la *extrema fragmentación*, fue eliminado de ellos.

Este tipo de sistemas de archivo ofrece características muy interesantes, aunque es un campo que aún requiere de mucha investigación e implementaciones ejemplos. Muchas de

las implementaciones en sistemas libres han llegado a niveles de funcionalidad aceptables, pero por diversas causas han ido perdiendo el interés o el empuje de sus desarrolladores, y su ritmo de desarrollo ha decrecido. Sin embargo, varios conceptos muy importantes han nacido bajo este tipo de sistemas de archivos, algunos de los cuales (como el de las *fotografías*) se han ido aplicando a sistemas de archivo estándar.

Por otro lado, dado el fuerte crecimiento que están registrando los medios de almacenamiento de estado sólido (en la sección C.1.2 se abordarán algunas de sus particularidades), y dado que estos sistemas aprovechan mejor varias de sus características, es probable que el interés en estos sistemas de archivos resurja.

8.4 Otros recursos

- *Practical File System Design*
<http://www.nobius.org/~dbg/>
Dominic Giampaolo (1999). El autor fue parte del equipo que implementó el sistema operativo BeOS, un sistema de alto rendimiento pensado para correr en estaciones de alto rendimiento, particularmente enfocado al video. El proyecto fracasó a la larga, y BeOS (así como BeFS, el sistema que describe) ya no se utilizan. Este libro, descargable desde el sitio Web del autor, tiene una muy buena descripción de varios sistemas de archivos, y aborda a profundidad técnicas que hace 15 años eran verdaderamente novedosas, y hoy forman parte de casi todos los sistemas de archivos con uso amplio, e incluso algunas que no se han logrado implementar y que BeFS sí ofrecía.
- *A method for the construction of Minimum-Redundancy Codes*
http://compression.graphicon.ru/download/articles/huff/huffman_1952_minimum-redundancy-codes.pdf
David A. Huffman (1952); Proceedings of the I. R. E
- *FAT Root Directory Structure on Floppy Disk and File Information*
http://www.codeguru.com/cpp/cpp/cpp_mfc/files/article.php/c13831
Mufti Mohammed (2007); Codeguru
- *File Allocation Table: 16bit*
<http://www.beginningtoseethelight.org/fat16/index.htm>
Peter Clark (2001)
- *Dosfsck: check and repair MS-DOS filesystems*
http://www.linuxcommand.org/man_pages/dosfsck8.html
Werner Almesberger (1997)
- *A Fast File System for UNIX*
<http://www.cs.berkeley.edu/~brewer/cs262/FFS.pdf>
Marshall Kirk Mckusick, William N. Joy, Samuel J. Leffler, Robert S. Fabry (1984); ACM Transactions on Computer Systems
- *The Design and Implementation of a Log-Structured File System*
<http://www.cs.berkeley.edu/~brewer/cs262/LFS.pdf>
Mendel Rosenblum, J. K. Ousterhout (1992); ACM Transactions on Computer Systems
- *The Second Extended File System: Internal Layout*
<http://www.nongnu.org/ext2-doc/>
Dave Poirier (2001-2011)
- *NILFS2 en Linux*
<http://cyanezfdz.me/articles/2012/08/nilfs2.html>
César Yáñez (2012)
- *Disk from the Perspective of a File System*

<http://queue.acm.org/detail.cfm?id=2367378>

Marshall Kirk McKusick (2012); ACM Queue

- Traducción al español: *Los discos desde la perspectiva de un sistema de archivos*

<http://cyanezfdz.me/post/los-discos-desde-la-perspectiva-de-un-sistema-de-archivos>

César Yáñez (2013)

- *A hash-based DoS attack on Btrfs*

<http://lwn.net/Articles/529077/>

Jonathan Corbet (2012); Linux Weekly News

- *Default /etc/apache2/mods-available/disk_cache.conf is incompatible with ext3*

<http://bugs.debian.org/682840>

Christoph Berg (2012). Reporte de fallo de Debian ilustrando los límites en números de archivos para Ext3.

- *File-system development with stackable layers*

<https://dl.acm.org/citation.cfm?id=174613.174616>

Heidemann y Popek (1994); ACM Transactions on Computer Systems

- *Serverless network file systems*

<https://dl.acm.org/citation.cfm?id=225535.225537>

Thomas E. Anderson et. al. (1996); ACM Transactions on Computer Systems

- *Log-structured file systems: There's one in every SSD*

<https://lwn.net/Articles/353411/>

Valerie Aurora (2009); Linux Weekly News

A — Software libre y licenciamiento

A.1 Software libre

Este apéndice, a diferencia del resto del libro, no se enfoca a asuntos técnicos, sino que a un aspecto *social*: A la construcción del conocimiento de forma colectiva, colaborativa, que ha resultado en el movimiento del software libre.

Si bien este tema es meramente tangente al que desarrolla el libro, los autores consideran importante incluirlo no sólo por la importancia que el software libre ha tenido para el desarrollo y estudio de los sistemas operativos, sino que directamente –como se explica en la sección A.2– para el presente libro en particular.

A.1.1 *Free as in Freedom*: El proyecto GNU

Citando la definición que brinda la Fundación del Software Libre (Free Software Foundation), el software libre es todo programa en el cual los usuarios tienen la libertad para ejecutar, copiar, distribuir, estudiar, modificar y mejorar el software. Esto es, todo programa cuyo modelo de licenciamiento respete las *cuatro libertades del software*:

- La libertad de ejecutar el programa para cualquier propósito.
- La libertad de estudiar cómo funciona el programa, y cambiarlo para que haga lo que usted quiera. El acceso al código fuente es una condición necesaria para ello.
- La libertad de redistribuir copias para ayudar a su prójimo.
- La libertad de distribuir copias de sus versiones modificadas a terceros. Esto le permite ofrecer a toda la comunidad la oportunidad de beneficiarse de las modificaciones. El acceso al código fuente es una condición necesaria para ello.

El software libre en tanto *movimiento ideológico* tiene bien identificados sus orígenes y génesis: En septiembre de 1983, Richard M. Stallman anunció el nacimiento del /Proyecto GNU/, que buscaba crear un sistema operativo tipo Unix, junto con todas las herramientas básicas del sistema (partiendo naturalmente desde la creación de un entorno de edición y un compilador). Tras sistematizar los fundamentos del proyecto, en marzo de 1985 publicó el *Manifiesto GNU*, documento que al día de hoy es lectura obligada para comprender al fenómeno que nació en ese momento.

Algunos meses más tarde, en octubre de 1985, creó la Fundación de Software Libre (FSF, Free Software Foundation), enfocada en la consecución de fondos para impulsar la creación de dicho sistema, en dar a conocer su trabajo, tanto para lograr que fuera ampliamente utilizado como para reclutar a más programadores y acelerar su ritmo de desarrollo.

El trabajo de la FSF es desde cualquier óptica impresionante por su magnitud y por su envergadura técnica. Sin embargo, probablemente su mayor contribución sea la *Licencia Pública General* (GPL), que será abordada en la sección A.1.4.

A.1.2 El software libre antes de GNU

El software libre como hoy se conoce existió mucho antes del proyecto GNU: Era la norma prácticamente hasta la aparición de las computadoras personales.

Los sistemas operativos, las herramientas de sistema y los compiladores eran, en un principio, entregadas por los fabricantes junto con el equipo de cómputo no sólo como *objetos binarios*, sino que en código fuente. Esto era natural: Los operadores de las computadoras no limitaban su uso a adecuar el software, sino que era común que adecuaran también el hardware: Cada equipo instalado era, hasta cierto punto, único.

Para hacerlo, claro, casi siempre era necesario modificar al software de forma correspondiente. Esto requería el acceso al código fuente, e implícitamente pasaba por las *cuatro libertades* ya enunciadas.

Durante las primeras décadas, prácticamente la totalidad del desarrollo del cómputo se llevó a cabo siguiendo la *tradición académica*: Los programas eran distribuidos, sea en cintas o incluso en listados impresos, requiriendo únicamente –como en un artículo científico– que se preserve la *atribución de autoría*. Sólo de este modo puede entenderse el desarrollo (y la supervivencia hasta el día de hoy) de sistemas con la relevancia de CP-CMS, creado por la muy pragmática y corporativa empresa IBM y cuya progenie sigue empleándose como núcleo de su arquitectura de virtualización z/VM (ver sección B.3) o Unix.

Unix nació como una reacción al sistema operativo *Multics*, desarrollado principalmente entre 1965 y 1970, y en el que participaban de forma conjunta varias empresas y el Instituto de Tecnología de Massachusetts (MIT). Multics resultó un proyecto demasiado grande, y AT&T lo abandonó en 1969; del equipo de AT&T que trabajaba en Unix, dos de los desarrolladores (Ken Thompson y Dennis Ritchie) comenzaron a escribir en 1969 un sistema mucho menos ambicioso, tomando algunos de los principales criterios de diseño, pero simplificando fuertemente el modelo de usuario y los requisitos en hardware. El nombre de Unix (originalmente *Unics*) es, de hecho, una broma sobre el nombre *Multics*.

Citando a Dennis Ritchie:¹

Lo que queríamos preservar no sólo era un buen ambiente en el cual programar, sino que un sistema alrededor del cual pudiera formarse una cofradía. Sabíamos por experiencia propia que la esencia del cómputo comunal, provisto por computadoras de acceso remoto y tiempo compartido, no se limita únicamente a escribir programas en una terminal en vez de emplear tarjetas perforadas, sino que favorece la comunicación cercana.

El párrafo inicial de este apéndice, que hace referencia a la *naturaleza social* del software libre, resuena con esta cita. El desarrollo de software va mucho más allá de su impacto técnico: Es una actividad tan social como cualquier otro desarrollo intelectual.

A lo largo de sus primeros diez años de vida, Unix pasó rápidamente de ser un sistema de *juguete* a ser, sin proponérselo, la base de desarrollo tecnológico sobre la cual se tendieron las bases de Internet. Decenas de empresas y universidades obtuvieron copias de Unix y lo modificaron, agregando funcionalidad —Y *compartiendo* esta nueva funcionalidad con el resto de la *comunidad* que se formó alrededor de Unix.

A.1.3 El software propietario como anomalía histórica

La *anomalía histórica* resulta, más bien, el auge que tuvo el software *proprietario* o *privativo*.² Una de las probables razones para éste puede ser, paradójicamente, el nacimiento del segmento del cómputo *aficionado*, como los presentados en la sección 2.4: Las primeras computadoras personales carecían del almacenamiento y poder de cómputo suficiente

¹/The Evolution of the Unix Time-sharing System; Language Design and Programming Methodology conference; Sydney, Australia, 1979. <http://cm.bell-labs.com/cm/cs/who/dmr/hist.html>

²Se designa de esta forma al software *no-libre*.

para siquiera compilar sus propios entornos operativos, razón por la cual las empresas productoras recurrieron a una *distribución exclusivamente binaria*.

El inicio de la masificación del cómputo llevó a que varias empresas nacientes identificaran un nicho de mercado donde podrían vender *licencias de uso* de los programas que produjeran, cobrando relativamente poco por cada licencia, pero aspirando a vender un gran volumen.

En este sentido, vale mucho la pena leer la [carta abierta a los entusiastas](#) que Bill Gates, socio de la entonces naciente y pequeña empresa Micro-Soft publicó en varias revistas de cómputo personal; la publicación original fue en el *Homebrew Computer Club Newsletter* (*periódico del club de cómputo casero*) en enero de 1976, y fue replicado en varias otras revistas.

Esta carta abierta tuvo amplias repercusiones, y desató un interesante debate que los lectores interesados podrán encontrar (y seguir en copias de los textos originales) desde el [artículo de Wikipedia respecto a esta /carta abierta/](#).

A.1.4 Esquemas libres de licenciamiento

Las licencias resultan fundamentales para comprender al software libre, tanto en su planteamiento ideológico primigenio como en el tipo de comunidad de desarrollo que aglutinan. Lo que es más, sólo se puede hablar de software libre en tanto esté asociado a un esquema de licenciamiento, dado que es éste el que determina las condiciones de uso a que estará sujeto un programa.³

A continuación, se abordan los dos principales enfoques del licenciamiento libre.

Licenciamiento académico/permisivo: MIT, BSD, X11, etc.

Las licencias derivadas del *primer momento* del software libre descrito son, en su conjunto, como licencias *académicas* o *permisivas*. Esto es porque, sin dejar de cubrir las cuatro libertades presentadas al principio del presente apéndice, el único requisito que imponen ante el usuario o distribuidor es el de la *atribución*.

De ahí el símil con la *academia*. No es de sorprender que algunas de las licencias más frecuentemente referidas de este tipo provengan directamente del ámbito Universitario: La licencia MIT proviene del Instituto de Tecnología de Massachusetts (ampliamente conocido bajo dicha sigla), y la licencia BSD hace referencia a la *Distribución de Software de Berkeley*, una de las principales ramas del desarrollo histórico de Unix, liderada por la Universidad de California en Berkeley.

Hay decenas de licencias que caben en esta categoría, con variaciones relativamente muy menores entre ellas. Los principales puntos que tienen en común son:

- Son licencias *muy cortas*. Siendo documentos legales, son muy sencillas y no dejan espacio a interpretaciones ambiguas.
- Se limitan a autorizar expresamente el uso del software, en fuente o en binario, y a *rechazar cualquier reclamo de garantía o responsabilidad por su uso*.
- Permiten la derivación en proyectos propietarios.

Una crítica reiterada al uso de estos esquemas de licenciamiento por parte de la FSF es que permiten la *privatización* de mejorías hechas al software libre — Pero al mismo tiempo, este punto constituye una de las principales fortalezas de este licenciamiento.

La masificación de Internet, y su adopción en los sistemas operativos más variados, se debió en gran parte a que el desarrollo de la *pila TCP/IP* fue liberado bajo un licenciamiento BSD. Al día de hoy, muchos de los componentes fundamentales de conectividad en

³Todos los países firmantes de la Convención de Berna garantizan la protección del derecho de autor *sin necesidad de registro*, de donde deriva que todo programa que sea publicado sin una licencia que *expresamente* lo haga libre, estará sujeto a *todos los derechos reservados*: Prohibición a todo tipo de uso sin autorización expresa y explícita del autor.

prácticamente la totalidad de sistemas operativos siguen incluyendo la nota de que los derechos de autor de determinados componentes pertenecen a *los regentes de la Universidad de California*.

Dado que empresas tan dispares como Sun Microsystems, Digital Research, IBM, Hewlett-Packard, Microsoft y Apple (por mencionar sólo a las que han dirigido distintos aspectos del mercado del cómputo) pudieron adoptar esta pila ya desarrollada, y que había una masa crítica de *sistemas abiertos* empleando TCP/IP, este protocolo de red creció hasta eclipsar a las diferentes apuestas propietarias de las diferentes empresas. Posteriormente, con el auge de los sistemas operativos libres, estos pudieron también adoptar esta base tecnológica en igualdad de condiciones.

Licenciamiento *Copyleft*: GPL, LGPL, MPL, CDDL, etc.

Para la FSF, el desarrollo de software es explícitamente un hecho social, y la creación de un sistema libre es un imperativo ético. La principal herramienta que emplearon para difundir y exigir la libertad del software fue el conjunto de licencias *Copyleft*.⁴ Y como se vio, si bien esto podría no ser compartido por los diferentes actores (personas y empresas), el desarrollo de Unix partió desde este mismo punto de vista.

Como se mencionó al inicio del presente apéndice, una de las principales obras de la FSF fue la creación de un modelo de licenciamiento que expresa este imperativo ético: Una familia de licencias cuyos principales exponentes son la *Licencia Pública General* (*General Public License*, GPL) y la *Licencia Pública General para Bibliotecas* (*Library General Public License*, LGPL, hoy renombrada a *Licencia Pública General Disminuida*, *Lesser General Public License*).

Existen varios ejemplos de licenciamiento que siguen estas ideas básicas; probablemente los más importantes sean la *Licencia Pública de Mozilla* (MPL) o la *Licencia Común de Distribución y Desarrollo* (CDDL, desarrollada por Sun Microsystems), y su principal diferencia con las presentadas por la FSF es que fueron propuestas no por grupos idealistas para el desarrollo de software aún inexistente, sino que por empresas que tenían ya un cuerpo de software, y encontraron este modelo como el más *sustentable* para continuar su desarrollo.

La principal característica de estos esquemas es que permiten el uso del software para cualquier fin, imponiendo como única condición que, *en caso de redistribución* (ya sea en su estado original o con modificaciones), el destinatario no sólo reciba el objeto binario ejecutable sino que el código fuente del cual éste provino, *bajo las mismas condiciones de licenciamiento original*.

Este esquema asegura que lo que una vez fue software libre *Copyleft* siempre lo siga siendo. El licenciamiento GPL ha llevado a que muchas empresas empleen al sistema operativo Linux como base para su desarrollo *contribuyan* sus cambios de vuelta a la comunidad — Convirtiendo a Linux al paso de los años de un sistema relativamente aficionado y con mediocre soporte a hardware en un sistema verdaderamente sólido y universal.

Muchos han criticado a este *espíritu viral* de las licencias *Copyleft*: Una vez que un proyecto incorpora componentes GPL, esta licencia podría *infectar* al proyecto entero obligándolo a adoptar esta licencia, resultando en graves perjuicios para las empresas que invierten en desarrollo. Si bien esto se ha demostrado falso repetidamente, sigue siendo un punto de propaganda frecuentemente empleado para evitar el empleo de software libre.

El objetivo del presente apéndice no es entrar a desmenuzar las diferencias entre estos esquemas o resolver las controversias, sino únicamente presentarlos de forma descriptiva.

⁴Término empleado para contraponerse a la noción de *Copyright*, Derecho de autor.

A.2 Obras culturales libres

Los distintos esquemas de software libre fueron logrando una masa crítica y poco a poco rompieron las predicciones de fracaso. 1998 fue un año crítico, en que varios importantes proyectos propietarios decidieron migrar a un licenciamiento libre por resultar más conveniente y sustentable.

Ya con esta experiencia previa, y conforme el acceso a Internet se masificaba cada vez más, comenzó a verse la necesidad de crear con esquemas similares de licenciamiento libre para otros productos de la creatividad humana, no únicamente para el desarrollo del software. Si bien las licencias académicas podrían aplicarse sin demasiado problema a productos que no fueran software, las licencias *Copyleft* llevan demasiadas referencias al código fuente y al binario como parte de su definición.

Del mismo modo que hay diferentes *escuelas de pensamiento* y puntos de vista ideológicos que han llevado al surgimiento de diversas licencias de software libre, respondiendo a distintas necesidades y matices ideológicos.

El proyecto *Wikipedia* fue anunciado en enero del 2001. Al convocar a *todo mundo* y no sólo a un manojo de especialistas, a crear contenido enciclopédico, este experimento iniciado por Jimmy Wales y Larry Sanger demostró que la creación es un acto profundamente social. Miles de voluntarios de todo el mundo han contribuido para hacer de la Wikipedia el compendio de conocimiento humano más grande de la historia. Al nacer, la Wikipedia adoptó el modelo de licenciamiento recomendado por la FSF para manuales y libros de texto: La *Licencia de Documentación Libre de GNU* (GFDL).

El modelo de la GFDL resulta, sin embargo, de difícil comprensión y aplicación para muchos autores, y la licencia no resulta apta para obras creativas más allá de lo que puede constituir *documentación*.

El marco regulatorio de la Convención de Berna, que rige al derecho de autor, estipula (como ya se mencionó) que toda creación plasmada en un medio físico está protegida, y todo uso no expresamente autorizado por una licencia expresa está prohibido. La tarea de crear esquemas de licenciamiento aptos para lo que se fue definiendo como *obras culturales libres* resultó más compleja por la riqueza de su expresión. En pocos años hubo una proliferación de licencias que buscaban ayudar a los autores de obras creativas de todo tipo — No se abordarán los distintos intentos, sino que —aprovechando que la distancia en tiempo permiten simplificar— se tocará sólo el esquema de licenciamiento que más impacto ha tenido.

A.2.1 La familia de licencias *Creative Commons*

En el año 2001, el abogado estadounidense Larry Lessig inició el proyecto *Creative Commons* (en adelante, CC). Citando del libro *Construcción Colaborativa del Conocimiento* (Wolf, Miranda 2011):

Pero no sólo el conocimiento formalizado puede compartirse. En 2001 nació Creative Commons (CC), impulsada por el abogado estadounidense Larry Lessig. Esta organización liderada localmente en una gran cantidad de países por personalidades versadas en temas legales, fue creada para servir como punto de referencia para quien quiera crear obras artísticas, intelectuales y científicas libres. Asimismo, ofrece un marco legal para que gente no experta en estos temas pueda elegir los términos de licenciamiento que juzgue más adecuados para su creación, sin tener que ahondar de más en las áridas estepas legales; se mantiene asesorada y liderada por un grupo de abogados, cuya principal labor es traducir y adecuar las licencias base de CC para cada una de las jurisdicciones en que sean aplicables. Alrededor de este modelo ha surgido un grupo de creadores, y una gran cantidad de sitios de alto perfil en la red

han acogido su propuesta. Si bien no todas las licencias de CC califican como cultura libre, algunas que claramente sí lo son han ayudado fuertemente a llevar estas ideas a la conciencia general.

CC creó un conjunto de licencias, permitiendo a los autores expresar distintos *grados* de libertad para sus obras. Uno de los principales elementos para su éxito y adopción masiva fue simplificar la explicación de estos distintos elementos, y la presentación de las alternativas bajo siglas mnemotécnicas.

Las licencias CC han pasado, al momento de edición del presente material, por cuatro versiones mayores, que han ido corrigiendo defectos en el lenguaje legal, y agregando o clarificando conceptos. Las opciones de las licencias CC son:⁵

CC0 (Dominio Público) La rigidez del convenio de Berna hace muy difícil en la mayor parte de las jurisdicciones el liberar una obra renunciando expresamente a todos los derechos patrimoniales que conlleva. La licencia *cero* o *dedicación al dominio público* explica esta renuncia expresa de derechos.

BY (Atribución) Todas las combinaciones de licencias CC a excepción de CC0 incluyen la cláusula de *atribución*: La obra puede emplearse para cualquier fin, pero toda redistribución debe reconocer el crédito de manera adecuada, proporcionar un enlace a la licencia, e indicar si se han realizado cambios. Puede hacerlo en cualquier forma razonable, pero no de forma tal que sugiera que tiene el apoyo del licenciante o lo recibe por el uso que hace.

SA (Compartir Igual) Si un usuario del la obra en cuestión decide mezclar, transformar o crear nuevo material a partir de ella, puede distribuir su contribución siempre que utilice la misma licencia que la obra original. Esto es, la cláusula *Compartir Igual* le confiere un carácter *Copyleft* al licenciamiento elegido.

NC (No Comercial) La obra puede ser utilizada, reproducida o modificada según lo permitido por los otros componentes elegidos de la licencia siempre y cuando esto no se considere o dirija hacia una ganancia comercial o monetaria.

ND (No Derivadas) La obra puede ser redistribuida acorde con los otros componentes elegidos de la licencia, pero debe ser redistribuida sólo si no se afecta su integridad: no puede ser modificada sin autorización expresa del autor.

Las licencias CC han sido empleadas para todo tipo de creaciones: Libros, música, películas, artes plásticas — Incluso, si bien no era su fin original, para licenciamiento de software. Y su gran éxito estaba no sólo en su uso, sino en que han llevado la noción del licenciamiento permisivo y de las obras culturales libres a una gran cantidad de creadores que, sin CC, probablemente habrían publicado sus creaciones bajo la tradicional modalidad *todos los derechos reservados*.

Creative Commons y las obras culturales libres

No todas las licencias CC califican de *obras culturales libres*: En 2005, Benjamin Mako Hill exploró el paralelismo entre CC y el movimiento del software libre en su texto *Towards a Standard of Freedom: Creative Commons and the Free Software Movement*; este trabajo sirvió como semilla para la definición de *Obras culturales libres*, publicada en 2006. De forma paralela a las cuatro libertades del software, esta definición acepta como obras libres a aquellas que garantizan:

- La libertad de usar el trabajo y disfrutar de los beneficios de su uso.
- La libertad de estudiar el trabajo y aplicar el conocimiento adquirido de él.
- La libertad de hacer y redistribuir copias, totales o parciales, de la información o expresión.

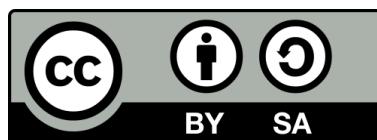
⁵Parte del texto aquí presentado ha sido tomado del [asistente para la elección de licencias](#) de Creative Commons; dicho texto está licenciado bajo un esquema CC-BY (atribución) 4.0.

- La libertad de hacer cambios y mejoras, y distribuir los trabajos derivados.

De las opciones de licenciamiento CC, las que están aprobados como obras culturales libres son CC0 (Dominio Público), BY (Atribución) y SA (Compartir Igual). Las variedades NC (No Comercial) y ND (No Derivadas), si bien permiten una mayor divulgación y circulación de la obra, restringen demasiado la apropiación que puede realizar un usuario, por lo que no constituyen obras culturales libres.

A.3 El licenciamiento empleado para la presente obra

Los autores de este libro buscaron contribuir con material de calidad libremente apropiable y reutilizable para la enseñanza superior en países hispanoparlantes. Para lograr este fin, todo el material contenido en el libro (texto, código fuente e imágenes) está licenciado bajo *Creative Commons Atribución 4.0 Internacional* (CC BY 4.0)⁶, salvo si se menciona explícitamente de otra manera.



Esto significa que usted es libre para:

- *Compartir* — copiar y redistribuir el material en cualquier medio o formato
- *Adaptar* — remezclar, transformar y crear a partir del material
- Para cualquier propósito, *incluso comercialmente*
- El licenciatario no puede revocar estas libertades en tanto usted siga los términos de la licencia

Bajo los siguientes términos:

- *Atribución* — Usted debe reconocer el crédito de una obra de manera adecuada, proporcionar un enlace a la licencia, e indicar si se han realizado cambios. Puede hacerlo en cualquier forma razonable, pero no de forma tal que sugiera que tiene el apoyo del licenciatario o lo recibe por el uso que hace.
- *CompartirIgual* — Si usted mezcla, transforma o crea nuevo material a partir de esta obra, usted podrá distribuir su contribución siempre que utilice la *misma licencia* que la obra original.

No hay restricciones adicionales — Usted no puede aplicar términos legales ni medidas tecnológicas que restrinjan legalmente a otros hacer cualquier uso permitido por la licencia.

A.4 Otros recursos

- *The GNU Manifesto*
<https://www.gnu.org/gnu/manifesto.html>
Richard Stallman (1985); Dr. Dobb's Journal
- *GNU General Public License version 3*
<https://gnu.org/licenses/gpl.html>
Free Software Foundation (2007)
- *GPL, BSD, and NetBSD - why the GPL rocketed Linux to success*
<http://www.dwheeler.com/blog/2006/09/01/>
David A. Wheeler (2006)
- *The Grumpy Editor's guide to free documentation licenses*
<https://lwn.net/Articles/108250/>
Jonathan Corbet (2004); Linux Weekly News

⁶<https://creativecommons.org/licenses/by/4.0/deed.es>

- *GNU Free Documentation License version 1.3*
<https://gnu.org/licenses/fdl.html>
Free Software Foundation (2008)
- *Construcción Colaborativa del Conocimiento*
<http://seminario.edusol.info/>
Gunnar Wolf, Alejandro Miranda (2011)
- *Towards a Standard of Freedom: Creative Commons and the Free Software Movement*
<http://www.advogato.org/article/851.html>
Benjamin Mako Hill (2005)
- *Obras culturales libres*
<http://freedomdefined.org/Definition/Espanol>
Freedom Defined (2006)

B — Virtualización

B.1 Introducción

La *virtualización* no es un concepto nuevo. Sin embargo, tras largos años de estar relegado a un segundo plano, en la actualidad se torna fundamental en referencia a los sistemas operativos, particularmente en rol de servidores. Este tema se abordará de momento desde una óptica más bien descriptiva, y posteriormente se profundizará en algunos de sus aspectos.

En primer término, es importante aclarar que el concepto de *virtualización* no se refiere a una única tecnología o metodología, es un término que agrupa a muy distintas tecnologías que existen – de diversas formas – desde hace décadas. Cada una de ellas tiene su lugar, con diferentes usos y propósitos, algunos de los cuales se usan de forma transparente para el usuario promedio.

Del mismo modo, aunque se abordarán diversas tecnologías que pueden clasificarse como virtualización, la línea divisoria entre cada una de ellas no siempre es clara. Una implementación específica puede caer en más de una categoría, o puede ir migrando naturalmente de un tipo hacia otro.

A nivel general, *virtualizar* consiste en proveer algo que no está ahí, aunque parece estarlo. Más específicamente, presentar a un sistema elementos que se comporten de la misma forma que un componente físico (*hardware*), sin que exista en realidad — Un acto de ilusionismo o de magia, en el cual se busca presentar el elemento de forma tan convincente que la ilusión se mantenga tanto como sea posible.¹

La naturaleza de dichos elementos, y el cómo se implementan, dependen del tipo de virtualización.

Para casi todos los casos que se presentan, se emplearán los siguientes términos:
Anfitrión El hardware o sistema *real*, que ofrece el mecanismo de virtualización. En inglés se le denomina *host*.

Huésped El sistema o las aplicaciones que se ejecutan en el entorno virtualizado. En inglés se les denomina *guest*.

B.2 Emulación

La técnica de virtualización más sencilla, y que hace más tiempo existe en las computadoras personales, es la emulación. Emular consiste en implementar *en software* algo que se presente como el hardware de un sistema de cómputo completo, típicamente de una arquitectura hardware distinta a la del anfitrión (la arquitectura *nativa*).² El emulador puede ser visto (de una forma tremadamente simplificada) como una lista de equivalencias,

¹Una aproximación inicial a este concepto puede ser un archivo con la imagen de un disco en formato ISO: mediante determinados mecanismos, es posible “engaños” a un sistema operativo de forma que “piense” que al acceder al archivo ISO está efectivamente leyendo un CD o DVD de una unidad que no existe físicamente.

²A lo largo de esta discusión, se hará referencia a la *arquitectura hardware* como al juego de instrucciones que puede ejecutar *nativamente* un procesador. Por ejemplo, un procesador x86 moderno puede ejecutar nativamente código i386 y x86_64, pero no ARM.

de cada una de las instrucciones en la arquitectura *huésped* a la arquitectura del sistema *anfitrión*.

Vale la pena recalcar que una emulación no se limita con traducir del lenguaje y la estructura de un procesador a otro — Para que una computadora pueda ser utilizada, requiere de una serie de chips de apoyo — Desde los controladores de cada uno de los *buses* hasta los periféricos básicos (teclado, video). Casi todas las emulaciones incluirán un paso más allá: Los periféricos mismos (discos, interfaces de red, puertos). Todo esto tiene que ser implementado por el emulador.

Resulta obvio que emular un sistema completo es *altamente* ineficiente. Los sistemas huéspedes resultantes típicamente tendrán un rendimiento cientos o miles de veces menor al del anfitrión.

Ahora bien, ¿qué pasa cuando hay dos arquitecturas de cómputo que emplean al mismo procesador? Este caso fue relativamente común en la década de los 80 y 90; si bien en general las computadoras de 8 bits no tenían el poder de cómputo necesario para implementar la emulación de arquitecturas similares, al aparecer tres líneas de computadoras basadas en el CPU Motorola 68000 (Apple Macintosh, Atari ST y Commodore Amiga), diferenciadas principalmente por sus *chipsets*, aparecieron emuladores que permitían ejecutar programas de una línea en la otra, prácticamente a la misma velocidad que en el sistema nativo.

Hoy en día, la emulación se emplea para hacer *desarrollos cruzados*, más que para emplear software *ya escrito y compilado*. La mayor parte de la emulación tradicional hoy se emplea para el *desarrollo de software*. Hoy en día, la mayor parte de las computadoras vendidas son sistemas *embebidos*³ o dispositivos móviles, que hacen imposible (o, por lo menos, muy difícil) desarrollar software directamente en ellos. Los programadores desarrollan en equipos de escritorio, corren entornos de prueba en emuladores del equipo destino. A pesar del costo computacional de realizar la emulación, la diferencia de velocidad entre los equipo de escritorio de gama alta y los embebidos permiten que frecuentemente la velocidad del emulador sea muy similar –incluso superior– a la del hardware emulado.

B.2.1 Emulando arquitecturas inexistentes

Pero la emulación no se limita a hardware existente, y no sólo se emplea por la comodidad de no depender de la velocidad de equipos específicos. Es posible crear emuladores para arquitecturas que *nunca han sido implementadas* en hardware real.

Esta idea viene de los 1970, cuando comenzó la explosión de arquitecturas. La Universidad de California en San Diego propuso una arquitectura llamada *p-system*, o *sistema-p*, la cual definiría una serie de instrucciones a las que hoy se clasificarían como *código intermedio* o *bytecode*, a ser ejecutado en una *máquina-p*, o *p-machine*. El lenguaje base para este sistema fue el *Pascal*, mismo que fue adoptado muy ampliamente de manera principal en entornos académicos a lo largo de los 1970 y 1980 por su limpieza y claridad estructural. Todo programa compilado para correr en un *sistema-p* correría sin modificaciones en cualquier arquitectura hardware que lo implementara.

Los *sistemas-p* gozaron de relativa popularidad hasta mediados de los 1980, logrando implementaciones para las arquitecturas de microcomputadoras más populares — El MOS 6502, el Zilog Z80 y el Intel 80x86.

Hay una diferencia muy importante entre la emulación de una arquitectura real y la de una arquitectura inexistente: Emular una computadora entera requiere implementar no sólo las instrucciones de su procesador, sino que *todos los chips de apoyo*, ¡incluso hay que convertir la entrada del teclado en las interrupciones que generaría un controlador

³Computadoras pequeñas, limitadas en recursos, y típicamente carentes de una interfaz usuario — Desde puntos de acceso y ruteadores hasta los controladores de cámaras, equipos de sonido, automóviles, y un larguísimo etcétera

de teclado! Emular una arquitectura hipotética permite manejar diversos componentes de forma abstracta, y permite definir estructuras de mucho más alto nivel que las que se encuentran implementadas en hardware. Por ejemplo, si bien resultaría impráctico crear como tipo de datos nativo para una arquitectura en hardware una abstracción como las cadenas de caracteres, estas sí existen como *ciudadanos de primera clase* en casi todas las arquitecturas meramente virtuales.

Esta idea ha sido ampliamente adoptada y forma parte de la vida diaria. En la década de los 1990, *Sun Microsystems* desarrolló e impulsó la arquitectura *Java*, actualizando la idea de las *máquinas-p* a los paradigmas de desarrollo que aparecieron a lo largo de veinte años, y dado que el cómputo había dejado de ser un campo especializado y escaso para masificarse, invirtiendo fuertemente en publicidad para impulsar su adopción.

Uno de los slogans que mejor describen la intención de Sun fue *WORA: Write Once, Run Anywhere* (Escribe una vez, corre donde sea). El equivalente a una *máquina-p* (rebautizada como *JVM: Máquina Virtual Java*) se implementaría para las arquitecturas hardware más limitadas y más poderosas. Sun creó también el lenguaje *Java*, diseñado para aprovechar la arquitectura de la *JVM*, enfatizando en la orientación a objetos e incorporando facilidades multi-hilos. Al día de hoy existen distintas implementaciones de la *JVM*, de diferentes empresas y grupos de desarrolladores y con diferentes focos de especialización, pero todas ellas deben poder ejecutar el *bytecode* de *Java*.

A principios de los años 2000, y como resultado del litigio con Sun que imposibilitó a Microsoft a desarrollar extensiones propietarias a *Java* (esto es, desarrollar máquinas virtuales que se salieran del estándar de la *JVM*), Microsoft desarrolló la arquitectura *.NET*. Su principal aporte en este campo es la separación definitiva entre lenguaje de desarrollo y código intermedio producido: La máquina virtual de *.NET* está centrada en el *CLI (Common Language Infrastructure, Infraestructura de Lenguajes Comunes)*, compuesta a su vez por el *CIL (Common Intermediate Language, Lenguaje Intermedio Común)*, que es la especificación del *bytecode* o código intermedio) y el *CLR (Common Language Runtime, Ejecutor del Lenguaje Común)*, que es la implementación de la máquina virtual sobre la arquitectura hardware nativa).

En los años 90, una de las principales críticas a *Java* (y esta crítica podría ampliarse hacia cualquier otra plataforma comparable) era el desperdicio de recursos de procesamiento al tener que traducir, una y otra vez, el código intermedio para su ejecución en el procesador. Hacia el 2010, el panorama había ya cambiado fuertemente. Hoy en día las máquinas virtuales implementan varias técnicas para reducir el tiempo que se desperdicia emulando:

Traducción dinámica Compilación parcial del código a ejecutar a formatos nativos, de modo que sólo la primera vez que se ejecuta el código intermedio tiene que ser traducido

Traducción predictiva Anticipar cuáles serán las siguientes secciones de código que tendrán que ser ejecutadas para, paralelamente al avance del programa, traducirlas a código nativo de forma preventiva

Compilación justo a tiempo (JIT) Almacenar copia del código ya traducido de un programa, de modo que no tenga que traducirse ni siquiera a cada ejecución, sino que sólo una vez en la vida de la máquina virtual

A través de estas estrategias, el rendimiento de las arquitecturas emuladas es ya prácticamente idéntico al del código compilado nativamente.

B.2.2 De lo abstracto a lo concreto

Si bien las arquitecturas de máquinas virtuales planteadas en el apartado anterior se plantearon directamente para no ser implementadas en hardware, el éxito comercial de la plataforma llevó a crear una línea de chips que ejecutara *nativamente* código intermedio *Java*, con lo cual podrían ahorrarse pasos y obtener mejor rendimiento de los sistemas

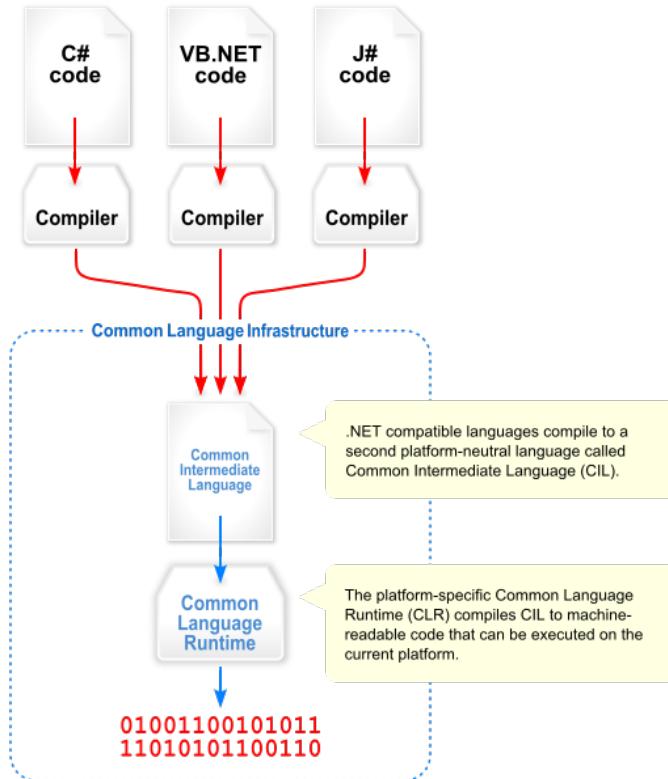


Figura B.1: Arquitectura de la infraestructura de lenguajes comunes (CLI) de .NET (Imagen de la Wikipedia: *Common Language Infrastructure*)

destino. Sun definió la arquitectura MAJC (*Microprocessor Architecture for Java Computing*, Arquitectura de microprocesadores para el cómputo con Java) en la segunda mitad de los 1990, e incluso produjo un chip de esta arquitectura, el MAJC 5200.

La arquitectura MAJC introdujo conceptos importantes que han sido retomados para el diseño de procesadores posteriores, pero la complejidad llevó a un rendimiento deficiente, y el chip resultó un fracaso comercial.

Es importante mencionar otra aproximación. Transitando en el sentido inverso al de Sun con MAJC, Transmeta, una empresa hasta entonces desconocida, anunció en el 2000 el procesador Crusoe, orientado al mercado de bajo consumo energético. Este procesador, en vez de implementar una arquitectura ya existente para entrar a un mercado ya muy competido y dinámico, centró su oferta en que Crusoe trabajaría mano a mano con un módulo llamado CMS (*Code Morphing Software*, Software de Transformación de Código), siendo así el primer procesador diseñado para *emular por hardware* a otras arquitecturas. Crusoe fue lanzado al mercado con el CMS para la arquitectura x86 de Intel, y efectivamente, la emulación era completamente transparente al usuario.⁴ El procesador mismo, además, no implementaba algunas características que hoy en día se consideran fundamentales, como una unidad de manejo de memoria, dado que eso podía ser implementado por software en el CMS. Separando de esta manera las características complejas a una segunda capa, podían mantenerse más bajos tanto el número de transistores (y, por tanto, el gasto energético) y los costos de producción.

La segunda generación de chips Transmeta (*Efficeon*) estaba basada en una arquitectura

⁴Empleando Transmeta, se podían observar ciertos comportamientos curiosos: Por ejemplo, dado el amplio espacio de caché que implementaba el CMS, el código ejecutable se mantenía *ya traducido* listo para el procesador, por lo cual la primera vez que se ejecutaba una función era notablemente más lenta que en ejecuciones posteriores. Sin embargo, si bien estas diferencias son medibles y no deben escapar a la vista de quien está analizando a conciencia estos procesadores, resultaban invisibles para el usuario final.

muy distinta, buscando un rendimiento mejorado. Pero, gracias al CMS, esto resulta imperceptible al usuario.

A pesar de estas ideas interesantes y novedosas, Transmeta no pudo mantener el dinamismo necesario para despegar, y cesó sus operaciones en 2009.

B.2.3 ¿Emulación o simulación?

Una pregunta frecuente que se presenta al hablar de este tema es acerca de la diferencia entre la *emulación* y la *simulación*. Todos los casos presentados anteriormente se tratan de *emulación*.

Emular significa *imitar las acciones de otro, procurando igualarlas e incluso excederlas* (Diccionario de la Real Academia Española, 23^a edición). Esto significa que un emulador reproduce todos los procesos internos que realizaría el sistema nativo, y busca cubrir todos los comportamientos respectivos implementando los mismos mecanismos.

Simular, por otra parte y según este mismo diccionario, significa *Representar algo, fingiendo o imitando lo que no es*. Un sistema simulador simula o finge las áreas de determinado sistema que interesan al usuario; puede emplear datos pre-cargados para generar ciertas respuestas, obviando los procesos que los generarían.

A diferencia de los ejemplos presentados a lo largo de esta sección, que llevan a ejecutar software arbitrario para la plataforma destino buscando idealmente que éstos no detecten siquiera una diferencia en comportamiento, un simulador puede presentar mucho mayor detalle en determinadas áreas, pero no realiza las funciones substantivas del sistema simulado. Por ejemplo, es muy común (incluso para el entrenamiento de pilotos reales) el uso de simuladores de vuelo; estos programas pueden representar una cabina equivalente a la de un avión real, con todos sus monitores y controles, pero nadie esperaría que lo trasladen de un lugar a otro. Muchos de los lectores habrán empleado software de simulación de circuitos electrónicos, que permiten el diseño y pruebas simples de circuitos, pero no esperarán que simular en la computadora un núcleo de ferrita rodeado por una bobina resulte en un receptor de radio.

B.3 Virtualización asistida por hardware

Actualmente se usa la virtualización como una herramienta para la consolidación de servicios, de gran ayuda para los administradores de sistemas. Este uso se refiere principalmente a lo que se presentará en este apartado, así como en las secciones B.4 (*Paravirtualización*) y B.5 (*Contenedores*). Y si bien este *zumbido* de la virtualización se ha producido mayormente a partir del 2006-2007, no se trata de tecnologías o ideas novedosas — Existe desde fines de los 1960. Hasta hace algunos años, sin embargo, se mantenía dentro del ámbito de los servidores a gran escala, fuera del alcance de la mayor parte de los usuarios. Es necesario estudiar la génesis de esta herramienta, para poder comprender mejor cómo opera y cómo se implementa.

En 1964, IBM creó la primer *familia de computadoras*, la serie 360. Presentaron la entonces novedosa idea de que una organización podía adquirir un modelo sencillo y, si sus necesidades se ajustaban al modelo de cómputo, podrían migrar fácilmente hacia modelos más poderosos dado que tendrían *compatibilidad binaria*.

Uno de los modelos de esta familia fue la S-360-67, con la característica distintiva en ser la única de la serie 360 en ofrecer una unidad de manejo de memoria (MMU), con lo cual permitía la reubicación de programas en memoria. Esto, sin embargo, creaba un problema: El software desarrollado para los equipos más pequeños de la familia estaba creado bajo un paradigma de usuario único, y si bien podría ser ejecutado en este modelo, eso llevaría a un desperdicio de recursos (dado que el modelo 67 tenía todo lo necesario para operar en modo multitarea).

La respuesta de IBM fue muy ingeniosa: Desarrollar un sistema operativo mínimo, CP (*Control Program*, Programa de Control) con el único propósito de crear y gestionar máquinas virtuales dentro del hardware S/360-67, dentro de *cada una de las cuales* pudiera ejecutarse *sin requerir modificaciones* un sistema operativo estándar de la serie 360. De entre los varios sistemas operativos disponibles para la S/360, el que más frecuentemente se utilizó fue el CMS,⁵ un sistema sencillo, interactivo y monousuario. La combinación CP/CMS proporcionaba un sistema operativo multiusuario, con plena protección entre procesos, y con compatibilidad con los modelos más modestos de la serie 360.

Aún después de la vida útil de la serie 360 original, IBM mantuvo compatibilidad con este modelo hacia la serie 370, e incluso hoy, 50 años más tarde, se encuentra aún como z/VM en la línea de *Sistemas z*.

Vale la pena mencionar que tanto CP como CMS fueron distribuidos desde el principio de forma consistente con lo que en la actualidad se conoce como *software libre*: IBM los distribuía en fuentes, con permiso de modificación y redistribución, y sus diferentes usuarios fueron enviando las mejoras que realizaban de vuelta a IBM, de modo que hoy en día incorpora el trabajo de 50 años de desarrolladores.

B.3.1 El hipervisor

El modelo CP/CMS lleva a una separación bastante limpia entre un *multiplexador de hardware* (CP) y el sistema operativo propiamente dicho (CMS). Y si bien la dupla puede ser vista como un sólo sistema operativo, conforme se fueron ejecutando en máquinas virtuales sistemas operativos más complejos se hizo claro que el CP tendría que ser *otra cosa*. Partiendo del concepto de que el sistema operativo es el *supervisor* de la actividad de los usuarios, yendo un paso más hacia arriba, se fue popularizando el nombre de *hipervisor* para el programa que administra y virtualiza a los supervisores. Algunas características primarias que definen qué es un hipervisor son:

- Es únicamente un *micro-sistema operativo*, dado que no cubre muchas de las áreas clásicas ni presenta las interfaces abstractas al usuario final — Sistemas de archivos, mecanismos de comunicación entre procesos, gestión de memoria virtual, evasión de bloqueos, etcétera.
- Se limita a gestionar bloques de memoria física contiguos y fijos, asignación de dispositivos y *poco* más que eso.
- Normalmente no tiene una interfaz usuario directa, sino que es administrado a través de llamadas privilegiadas desde alguno de los sistemas operativos huésped.

Estas líneas se han ido haciendo borrosas con el tiempo. Ahora, por ejemplo, muchos hipervisores entienden a los sistemas de archivos, permitiendo que los espacios de almacenamiento ofrecidos a sus sistemas operativos huésped sean simples archivos para el sistema anfitrión (y no particiones o dispositivos enteros). Algunos hipervisores, como KVM bajo Linux se presentan integrados como un componente más de un sistema operativo estándar.

B.3.2 Virtualización asistida por hardware en x86

Hasta alrededor del año 2005, la virtualización no se mencionaba muy frecuentemente. Si bien había hardware virtualizable 40 años atrás, era hardware bastante especializado — y caro. Ese año, Intel sacó al mercado los procesadores con las extensiones necesarias para la virtualización, bajo el nombre *Vanderpool Technology* (o VT-x). Al año siguiente, AMD hizo lo propio, denominándolas *extensiones Pacifica*. Hoy en día, casi todas las computadoras de escritorio de rango medio-alto tienen el soporte necesario para llevar

⁵Originalmente, las siglas CMS eran por el *Cambridge Monitor System*, por haber sido desarrollado en la división de investigación de IBM en Cambridge, pero posteriormente fue renombrado a {Conversational Monitor System}, *Sistema de Monitoreo Conversacional*

a cabo virtualización asistida por hardware. Y si bien en un principio el tema tardó en tomar tracción, llevó a un replanteamiento completo de la metodología de trabajo tanto de administradores de sistemas como de programadores.

En contraste con las arquitecturas diseñadas desde un principio para la virtualización, los usuarios de computadoras personales (inclusive cuando estas son servidores en centros de datos — Siguen estando basadas en la misma arquitectura básica) se enfrentan a una mayor variedad de dispositivos para todo tipo de tareas.⁶ Y si bien la virtualización permite aparentar varias computadoras distintas corriendo sobre el mismo procesador, esta no incluye a los dispositivos. Al presentarse una máquina virtual, el sistema anfitrión esta casi siempre⁷ emulando hardware. Claro está, lo más frecuente es que el hipervisor ofrezca a los huéspedes la emulación de dispositivos relativamente viejos y simples.⁸ Esto no significa que estén limitados a las prestaciones del equipo emulado (por ejemplo, a los 10Mbps para los que estaba diseñada una tarjeta de red NE2000), sino que la interfaz del núcleo para enviar datos a dicho dispositivo es una sencilla y que ha sido empleada tanto tiempo que presenta muy poca inestabilidad.

Y este último punto permite un acercamiento mayor a una de las ventajas que ofrecen los sistemas operativos virtualizados — La estabilidad. Los controladores de dispositivos provistos por fabricante han sido responsabilizados una y otra vez, y con justa razón, de la inestabilidad de los sistemas operativos de escritorio. En particular, son en buena medida culpables de la fama de inestabilidad que obtuvo Windows. Los fabricantes de hardware no siempre gozan de suficiente conocimiento acerca del sistema operativo como para escribir controladores suficientemente seguros y de calidad, y por muchos años, los sistemas Windows no implementaban mayor verificación al comportamiento de los controladores — que, siendo un sistema monolítico, eran código ejecutado con privilegios de núcleo.

Al emplear el sistema operativo huésped únicamente controladores ampliamente probados y estabilizados a lo largo de muchos años, la estabilidad que ofrece una máquina virtualizada muchas veces supera a la que obtendría ejecutándose de forma nativa. Claro, el conjunto de máquinas virtuales que se ejecute dentro de un sistema anfitrión sigue siendo susceptible a cualquier inestabilidad del mismo sistema anfitrión, sin embargo, es mucho menos probable que un programa mal diseñado logre congelarse esperando respuesta del hardware (emulado), y mucho menos afectar a los demás huéspedes.

B.4 Paravirtualización

La virtualización asistida por hardware, por conveniente que resulte, sigue presentando algunas desventajas:

- No todos los procesadores cuentan con las extensiones de virtualización. Si bien cada vez es más común encontrarlas, es aún en líneas generales un factor de diferenciación entre las líneas económicas y de lujo.
- La capa de emulación, si bien es delgada, conlleva un cierto peso.
- Si bien es posible virtualizar arquitecturas como la x86, hay muchas arquitecturas para las cuales no existen las extensiones hardware necesarias.

⁶Una descripción completa de la complejidad a la que debe enfrentarse un hipervisor bajo arquitectura x86 excede con mucho el ámbito del presente texto; se sugiere a los lectores interesados referirse al excelente artículo de [Bugnion et. al. \(2012\)](#) detallando la implementación de VMWare.

⁷Hay mecanismos para reservar y dirigir un dispositivo físico existente a una máquina virtual específica, pero hacerlo implica que éste dispositivo no será *multiplexado* hacia las demás máquinas virtuales que se ejecuten paralelamente.

⁸Por ejemplo, KVM bajo Linux emula tarjetas de red tipo NE2000, tarjetas de sonido tipo Soundblaster16 y tarjetas de video Cirrus Logic, todos ellos de la década de los 1990.

La *paravirtualización*, o *virtualización asistida por el sistema operativo*, parte de un planteamiento distinto: En vez de engañar al sistema operativo para que funcione sobre un sistema que parece real pero no lo es, la paravirtualización busca hacerlo *con pleno conocimiento y cooperación* por parte de los sistemas huéspedes. Esto es, la paravirtualización consiste en alojar a sistemas operativos huésped que, a sabiendas de que están corriendo en hardware virtualizado, *no hacen llamadas directas a hardware* sino que las traducen a llamadas al sistema operativo anfitrión.

Vale la pena reiterar en este punto: Los sistemas operativos huésped bajo un entorno paravirtualizado saben que no están corriendo sobre hardware real, por lo que en vez de enviar las instrucciones que controlen al hardware, envían llamadas al sistema a su hipervisor. Hasta cierto punto, el proceso de adecuación de un sistema para que permita ser paravirtualizado puede ser equivalente a adecuar al sistema operativo para que corra en una arquitectura nueva — Muy parecida a la del hardware *real*, sí, pero con diferencias fundamentales en aspectos profundos.

Y si bien ya se explicó en la sección anterior que la virtualización puede ayudar a presentar un sistema idealizado que reduzca la inestabilidad en un sistema operativo, al hablar de paravirtualización este beneficio naturalmente crece: Los controladores de hardware sencillos y bien comprendidos que se usaban para gestionar los dispositivos emulados se convierten casi en simples pasarelas de llamadas al sistema, brindando además de una sobrecarga mínima, aún mayor estabilidad por simplicidad del código.

B.4.1 Paravirtualización y software libre

La paravirtualización resulta muy atractiva, presentando muy obvias ventajas. Pero a pesar de que es posible emplearla en cualquier arquitectura hardware, no siempre es posible emplearla.

Como se mencionó anteriormente, incorporar dentro de un sistema operativo el soporte para una arquitectura de paravirtualización es casi equivalente a traducirlo a una nueva arquitectura hardware. Para que los autores de un entorno que implemente paravirtualización logren que un sistema operativo nuevo pueda ser ejecutado en su arquitectura, deben poder manipular y modificar su código fuente: De otra manera, ¿cómo se le podría adecuar para que supiera desenvolverse en un entorno no nativo?

El proyecto de gestión de virtualización y paravirtualización *Xen*, hoy impulsado por la empresa *XenSource*, nació como un proyecto académico de la Universidad de Cambridge, presentando su versión 1.x a través de un artículo en 2003 (ver [Xen and the Art of Virtualization](#)). Este artículo presenta su experiencia paravirtualizando a una versión entonces actual de Linux y de Windows XP. Sin embargo, Xen sólo pudo ser empleado por muchos años como plataforma de paravirtualización de Linux porque, dado que la adaptación de Windows se realizó bajo los términos del *Academic Licensing Program*, que permitía a los investigadores acceso y modificación al código fuente, pero no su redistribución — La versión paravirtualizable de Windows XP existe, pero no puede distribuirse fuera de XenSource.

En tanto, el trabajo necesario para lograr la paravirtualización de un sistema operativo libre, como Linux, FreeBSD u otros, puede ser libremente redistribuído. No sólo eso, sino que el esfuerzo de realizar la adaptación pudo compartirse entre desarrolladores de todo el mundo, dado que esta entonces novedosa tecnología resultaba de gran interés.

B.4.2 Paravirtualización de dispositivos

Las ideas derivadas de la paravirtualización pueden emplearse también bajo entornos basados en virtualización plena: Si el sistema operativo está estructurado de una forma modular (sin que esto necesariamente signifique que es un sistema *microkernel*, sino que permite la carga dinámica de controladores o *drivers* para el hardware, como prácticamente

la totalidad de sistemas disponibles comercialmente hoy en día), no hace falta modificar al sistema operativo completo para gozar de los beneficios de la paravirtualización en algunas áreas.

De esta manera, si bien es posible ejecutar un sistema operativo *sin modificaciones* que espera ser ejecutado en hardware real, los dispositivos que típicamente generan más actividad de entrada y salida⁹ pueden ser atendidos por drivers paravirtuales. Por supuesto, varios aspectos que son parte del núcleo *duro* del sistema, como la administración de memoria o el manejo de interrupciones (incluyendo al temporizador) tendrán que seguirse manejando a través de una emulación, aunque mucho más delgada.

Según mediciones empíricas realizadas en 2007 por Qumranet (quienes liderearon el desarrollo del módulo de virtualización asistido por hardware KVM en Linux), las clases de dispositivos `virtio` y `pv` resultaron entre 5 y 10 veces más rápidas que la emulación de dispositivos reales.

Mediante esta estrategia es posible ejecutar sistemas operativos propietarios, como los de la familia Windows, con buena parte de las ventajas de la paravirtualización, sobre entornos de virtualización asistida por hardware.

B.5 Contenedores, o virtualización a nivel sistema operativo

Una estrategia completamente distinta para la creación de máquinas virtuales es la de *contenedores*. A diferencia de emulación, virtualización asistida por hardware y paravirtualización, al emplear contenedores *sólo se ejecuta un sistema operativo*, que es el mismo para los sistemas anfitrión y huesped. El anfitrión implementará una serie de medidas para *aumentar el grado de separación* que mantiene entre procesos, agregando la noción de *contextos* o *grupos* que se describirán en breve. Dado que el sistema operativo es el único autorizado para tener acceso directo al hardware, no hace falta ejecutar un hipervisor.

Podría presentarse un símil: Las tecnologías antes descritas de virtualización implementan *hardware virtual* para cada sistema operativo, mientras que los contenedores más bien presentan un *sistema operativo virtual* para el conjunto de procesos que definen el comportamiento de cada máquina virtual — Muchos autores presentan a la virtualización por contenedores bajo el nombre *virtualización a nivel sistema operativo*. Y si bien el efecto a ojos del usuario puede ser comparable, este método más que una multiplexación de máquinas virtuales sobre hardware real opera a través de restricciones adicionales sobre los procesos de usuario.

Al operar a un nivel más alto, un contenedor presenta algunas limitantes adicionales (principalmente, se pierde la flexibilidad de ejecutar sistemas operativos distintos), pero obtiene también importantes ventajas.

El desarrollo histórico de los contenedores puede rastrearse a la llamada al sistema `chroot()`, que restringe la visión del sistema de archivos de un proceso a sólo el directorio hacia el cual ésta fue invocada.¹⁰ Esto es, si dentro de un proceso se invoca `chroot('/usr/local')` y posteriormente se le pide abrir el archivo `/boot.img`, a pesar de que éste indique una ruta absoluta, el archivo que se abrirá será `/usr/local/boot.img`.

Ahora bien, `chroot()` no es (ni busca ser) un verdadero aislamiento, sólo proporciona

⁹Medios de almacenamiento, interfaz de red y salida de video

¹⁰La llamada `chroot()` fue creada por Bill Joy en 1982 para ayudarse en el desarrollo del sistema Unix 4.2BSD. Joy buscaba probar los cambios que iba haciendo en los componentes en espacio de usuario del sistema sin modificar su sistema *vivo* y en producción, esto es, sin tener que reinstalar y reiniciar cada vez, y con esta llamada le fue posible instalar los cambios dentro de un directorio específico y probarlos como si fueran en la raíz.

un inicio¹¹ — Pero conforme más usuarios comenzaban a utilizarlo para servicios en producción, se hizo claro que resultaría útil ampliar la conveniencia de `chroot()` a un verdadero aislamiento.

El primer sistema en incorporar esta funcionalidad fue *FreeBSD*, creando el subsistema *Jails* a partir de su versión 4.0, del año 2000. No tardaron mucho en aparecer implementaciones comparables en los distintos sistemas Unix. Hay incluso un producto propietario, el *Parallels Virtuozzo Containers*, que implementa esta funcionalidad para sistemas Windows.

Un punto importante a mencionar cuando se habla de contenedores es que se pierde buena parte de la universalidad mencionada en las secciones anteriores. Si bien las diferentes implementaciones comparten principios básicos de operación, la manera en que implementan la separación e incluso la nomenclatura que emplean difieren fuertemente.

El núcleo del sistema crea un *grupo* para cada *contenedor* (también conocido como *contexto de seguridad*), aislandolos entre sí por lo menos en los siguientes áreas:

Tablas de procesos Los procesos en un sistema Unix se presentan como un árbol, en cuya raíz está siempre el proceso 1, `init`. Cada contenedor inicia su existencia ejecutando un `init` propio y enmascarando su identificador de proceso real por el número 1

Señales, comunicación entre procesos Ningún proceso de un contenedor debe poder interferir con la ejecución de uno en otro contenedor. El núcleo restringe toda comunicación entre procesos, regiones de memoria compartida y envío de señales entre procesos de distintos grupos.

Interfaces de red Varía según cada sistema operativo e implementación, pero en líneas generales, cada contenedor tendrá una interfaz de red con una *dirección de acceso a medio* (MAC) distinta.¹² Claro está, cada una de ellas recibirá una diferente dirección IP, y el núcleo ruteará e incluso aplicará reglas de firewall entre ellas.

Dispositivos de hardware Normalmente los sistemas huesped no tienen acceso directo a ningún dispositivo en hardware. En algunos casos, el acceso a dispositivos será multiplexado, y en otros, un dispositivo puede especificarse a través de su configuración. Cabe mencionar que, dado que esta multiplexión no requiere *emulación* sino que únicamente una cuidadosa *planificación*, no resulta tan oneroso como la emulación.

Límites en consumo de recursos Casi todas las implementaciones permiten asignar cotas máximas para el consumo de recursos compartidos, como espacio de memoria o disco o tiempo de CPU empleados por cada uno de los contenedores.

Nombre del equipo Aunque parezca trivial, el nombre con el que una computadora *se designa a sí misma* debe también ser aislado. Cada contenedor debe poder tener un nombre único e independiente.

Una de las principales características que atrae a muchos administradores a elegir la virtualización por medio de contenedores es un consumo de recursos óptimo: Bajo los demás métodos de virtualización (y particularmente al hablar de emulación y de virtualización asistida por hardware), una máquina virtual siempre ocupará algunos recursos, así esté inactiva. El hipervisor tendrá que estar notificando a los temporizadores, enviando los paquetes de red recibidos, etcétera. Bajo un esquema de contenedores, una máquina virtual que no tiene trabajo se convierte sencillamente en un grupo de procesos *dormidos*, probables candidatos a ser *paganados* a disco.

¹¹Como referencia a por qué no es un verdadero aislamiento, puede referirse al artículo {How to break out of a =chroot()= jail} (Simes, 2002)

¹²Es común referirse a las direcciones MAC como direcciones físicas, sin embargo, todas las tarjetas de red permiten configurar su dirección, por lo cual la apelación *física* resulta engañosa.

B.6 Otros recursos

- *Bringing Virtualization to the x86 Architecture with the Original VMware Workstation*
<https://dl.acm.org/citation.cfm?doid=2382553.2382554>
Bugnion et. al. (2012); ACM Transactions on Computer Systems
- *Performance Evaluation of Intel EPT Hardware Assist*
http://www.vmware.com/pdf/Perf_ESX_Intel-EPT-eval.pdf
VMWare Inc. (2006-2009)
- *Performance Aspects of x86 Virtualization*
http://communities.vmware.com/servlet/JiveServlet/download/1147092-17964/PS_TA68_288534_166-1_FIN_v5.pdf
Ole Agesen (2007); VMWare
- *Xen and the Art of Virtualization*
<http://www.cl.cam.ac.uk/netos/papers/2003-xensosp.pdf>
Paul Barham, Boris Dragovic et. al. (2003)
- *KVM: The Linux Virtual Machine Monitor*
<http://kernel.org/doc/ols/2007/ols2007v1-pages-225-230.pdf>
Avi Kivity, Yaniv Kamay, Dor Laor, Uri Lublin, Anthony Liguori (2007) Qumranet / IBM)
- *KVM PV devices*
[http://www.linux-kvm.org/wiki/images/d/dd/KvmForum2007\\$kvm_pv_drv.pdf](http://www.linux-kvm.org/wiki/images/d/dd/KvmForum2007$kvm_pv_drv.pdf)
Dor Laor (2007); Qumranet
- *How to break out of a =chroot()= jail*
<http://www.bpfh.net/computing/docs/chroot-break.html>
Simes (2002)
- *Notes from a container*
<http://lwn.net/Articles/256389/>
Jonathan Corbet (2007); Linux Weekly News
- *CGROUPS*
<https://www.kernel.org/doc/Documentation/cgroups/cgroups.txt>
Paul Menage (Google) (2004-2006), kernel.org

C— El medio físico y el almacenamiento

C.1 El medio físico

A lo largo del presente texto, particularmente de los capítulos 7 y 8 y siguiendo las prácticas a que ha impuesto la realidad de los últimos 40 años, el término genérico de *disco* se ha empleado prácticamente como sinónimo de *medio de almacenamiento a largo plazo*.

En este apéndice se abordan en primer término las características principales del medio aún prevalente, los discos duros magnéticos rotativos, y una introducción a las diferencias que presentan respecto a otros medios, como los discos ópticos y los de estado sólido, así como las implicaciones que éstos tienen sobre el material presentado en el capítulo 8.

Cabe mencionar que la razón de separar este contenido hacia un apéndice es que, si bien estas funciones resultan relevantes para los sistemas operativos y estos cada vez más van asumiendo las funciones que aquí serán descritas, estas comenzaron siendo implementadas por hardware especializado; fue apenas hasta la aparición de los esquemas de manejo avanzado de volúmenes (que serán cubiertos en la sección C.3) que entran al ámbito del sistema operativo.

C.1.1 Discos magnéticos rotativos

El principal medio de almacenamiento empleado en los últimos 40 años es el *disco magnético*. Hay dos tipos diferentes de disco, aunque la lógica de su funcionamiento es la misma: Los *discos duros* y los *discos flexibles* (o *floppies*).

La principal diferencia entre estos es que los primeros son típicamente almacenamiento *interno* en los equipos de cómputo, y los segundos fueron pensados para ser almacenamiento *transportable*. Los discos duros tienen mucha mayor capacidad y son mucho más rápidos, pero a cambio de ello, son correspondientemente más sensibles a la contaminación por partículas de polvo y a daños mecánicos, razón por la cual hoy en día se venden, junto con el mecanismo lector e incluso la electrónica de control, en empaque sellado.

Un disco flexible es una hoja de material plástico, muy similar al empleado en las cintas magnéticas, resguardado por un estuche plástico. Al insertarse el disco en la unidad lectora, esta lo hace girar sujetándolo por el centro, y las cabezas lectoras (en un principio una sola; posteriormente aparecieron las unidades de doble cara, con dos cabezas lectoras) se deslizan por una ventana que tiene el estuche.

La mayor parte de los discos flexibles presentaban velocidades de rotación de entre 300 y 400 revoluciones por minuto — Presentaban, pues, una *demora rotacional* de entre 0.15 y 0.2 segundos. La *demora rotacional* es el tiempo que toma la cabeza lectora en volver a posicionarse sobre un mismo sector del disco. (Ver figura C.1)

A lo largo de más de 20 años se presentaron muy diferentes formatos físicos siguiendo esta misma lógica, designándose principalmente por su tamaño (en pulgadas). La capacidad de los discos, claro está, fue creciendo con el paso de los años — Esto explica la aparente contradicción de que los discos (físicamente) más chicos tenían más capacidad que los más grandes.

El nombre de *disco duro* o *disco flexible* se debe al medio empleado para el almacenamiento.

Cuadro C.1: Principales formatos de disco flexible que se popularizaron en el mercado

	8 pulgadas	5.25 pulgadas	3.5 pulgadas
Fecha de introducción	1971	1976	1982
Capacidad	150KB-1.2MB	110KB-1.2MB	264KB-2.88MB
Velocidad (kbit/s)	33	125-500	250-1000
Pistas por pulgada	48	48-96	135

miento de la información (y no a la rigidez de su *estuche*, como mucha gente erróneamente cree): Mientras que los discos flexibles emplean una hoja plástica flexible, los discos duros son metálicos. Los discos están *permanentemente* montados sobre un eje, lo que permite que tengan una velocidad de giro entre 20 y 50 veces mayor que los discos flexibles — Entre 4,200 y 15,000 revoluciones por minuto (RPM), esto es, con una demora rotacional de entre 2 y 7.14 milisegundos.

Además, a excepción de algunos modelos tempranos, los discos duros constituyen un paquete cerrado y sellado que incluye las cabezas de lectura y escritura, y toda la electrónica de control. Esto permite que los discos duros tengan densidades de almacenamiento y velocidades de transmisión muy superiores a la de los discos flexibles: Los primeros discos duros que se comercializaron para computadoras personales eran de 10MB (aproximadamente 70 discos flexibles de su época), y actualmente hay ya discos de 4TB. La velocidad máxima de transferencia sostenida hoy en día es superior a los 100MB por segundo, 100 veces más rápido que la última generación de discos flexibles.

Para medir la eficiencia de un disco duro, además de la *demora rotacional* presentada unos párrafos atrás, el otro dato importante es el tiempo que toma la cabeza en moverse a través de la superficie del disco. Hoy en día, las velocidades más comunes son de 20ms para un *recorrido completo* (desde el primer hasta el último sector), y entre 0.2ms y 0.8ms para ir de un cilindro al inmediato siguiente. Como punto de comparación, el recorrido completo en una unidad de disco flexible toma aproximadamente 100ms, y el tiempo de un cilindro al siguiente va entre 3 y 8ms.

Notación C-H-S

En un principio y hasta la década de los noventa, el sistema operativo siempre hacía referencia a la ubicación de un bloque de información en el disco es conocido como la *notación C-H-S* — Indicando el cilindro, cabeza y sector (*Cylinder, Head, Sector*) para ubicar a cada bloque de datos. Esto permite mapear el espacio de almacenamiento de un disco a un espacio tridimensional, con cual resulta trivial ubicar a un conjunto de datos en una región contigua.

La *cabeza* indica a cuál de las superficies del disco se hace referencia; en un disco flexible hay sólo una o dos cabezas (cuando aparecieron las unidades de doble lado eran en un lujo, y al paso de los años se fueron convirtiendo en la norma), pero en un disco duro es común tener varios *platos* paralelos. Todas las cabezas van fijas a un mismo motor, por lo que no pueden moverse de forma independiente.

El *cilindro* indica la distancia del centro a la orilla del disco. Al cilindro también se le conoce como *pista* (*track*), una metáfora heredada de la época en que la música se distribuía principalmente en discos de vinil, y se podía ver a simple vista la frontera entre una pista y la siguiente.

Un *sector* es un segmento de arco de uno de los cilindros, y contiene siempre la misma cantidad de información (históricamente 512 bytes; en actualmente se están adoptando gradualmente sectores de 4096 bytes. Refiérase a la sección C.1.1 para una mayor discusión al respecto.)

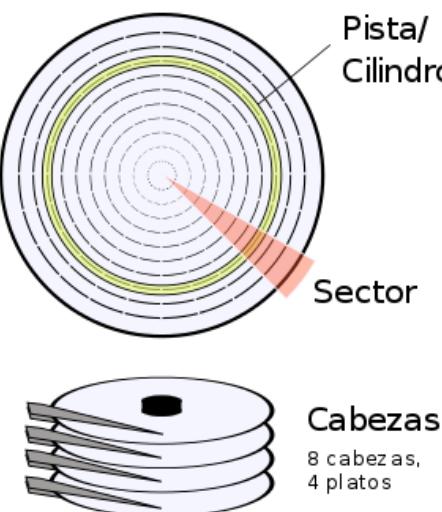


Figura C.1: Coordenadas de un disco duro, ilustrando su geometría basada en cabeza, cilindro y sector. (Imagen de la Wikipedia: *Cilindro Cabeza Sector*)

Un archivo almacenado secuencialmente ocupa *sectores adyacentes* a lo largo de una misma pista y con una misma cabeza.

Algoritmos de planificación de acceso a disco

Las transferencias desde y hacia los discos son uno de los procesos más lentos de los que gestiona el sistema operativo. Cuando éste tiene varias solicitudes de transferencia pendientes, resulta importante encontrar un mecanismo óptimo para realizar la transferencia, minimizando el tiempo de demora. A continuación se describirán a grandes rasgos tres de los algoritmos históricos de planificación de acceso a disco — Para abordar después el por qué estos hoy en día casi no son empleados.

Como con los demás escenarios en que se han abordado algoritmos, para analizar su rendimiento, el análisis se realizará sobre una *cadena de referencia*. Este ejemplo supone un disco hipotético de 200 cilindros, la cadena de solicitudes 83, 175, 40, 120, 15, 121, 41, 42, y teniendo la cabeza al inicio de la operación en el cilindro 60.

En la figura C.2 puede apreciarse de forma gráfica la respuesta que presentarían los distintos algoritmos ante la cadena de referencia dada.

FIFO Del mismo modo que cuando fueron presentados los algoritmos de asignación de procesador y de reemplazo de páginas, el primero y más sencillo de implementar es el *FIFO* — *Primero llegado, primero servido*.

Este algoritmo puede verse como muy *justo*, aunque sea muy poco eficiente: El movimiento total de cabezas para el caso planteado es de 622 cilindros, equivalente a poco más que recorrer de extremo a extremo el disco completo tres veces. Esto es, despreciando la demora rotacional la demora mecánica para que el brazo se detenga por completo antes de volver a moverse, esta lectura tomaría un mínimo de 60ms, siendo el recorrido completo del disco 20ms.

Puede identificarse como causante de buena parte de esta demora a la quinta posición de la cadena de referencia: Entre solicitudes para los cilindros contiguos 120 y 121, llegó una solicitud al 15.

Atender esta solicitud en FIFO significa un desplazamiento de $(120 - 15) + (121 - 15) = 211$ cilindros, para volver a quedar prácticamente en el mismo lugar de inicio. Una sola solicitud resulta responsable de la tercera parte del tiempo total.

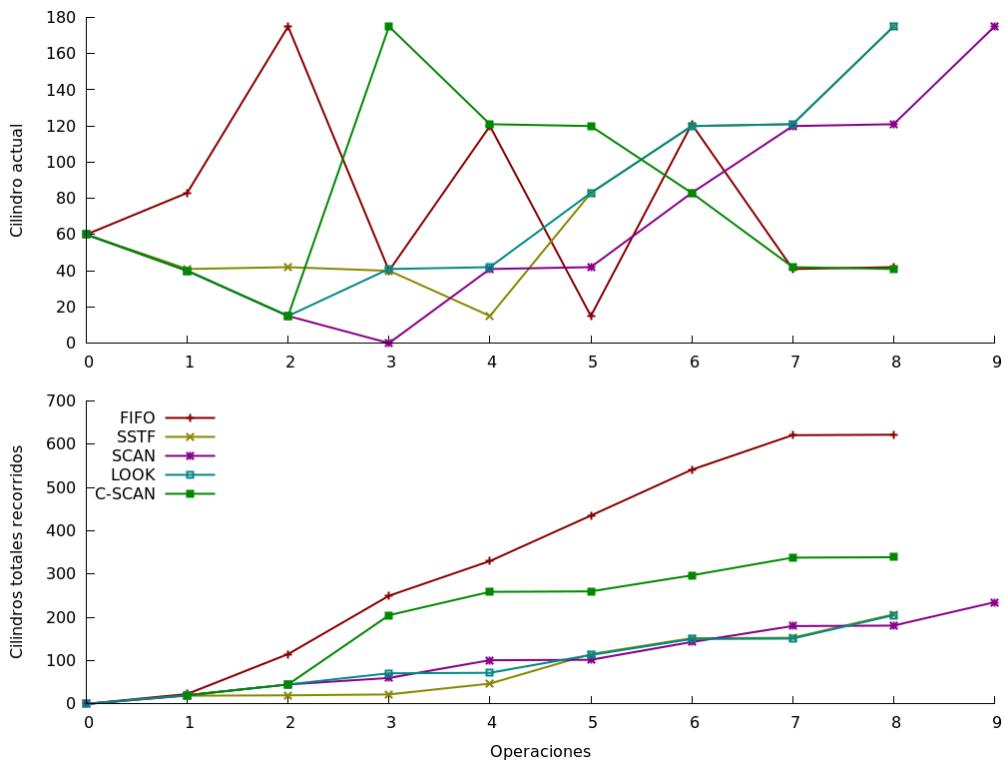


Figura C.2: Movimientos de las cabezas bajo los diferentes algoritmos planificadores de acceso a disco, indicando la distancia total recorrida por la cabeza bajo cada uno, iniciando con la cabeza en la posición 60. Para SCAN, LOOK y C-SCAN, se asume que la cabeza inicia avanzando en dirección decreciente.

SSTF Ahora bien, si el factor que impone la principal demora es el movimiento de la cabeza, el segundo algoritmo busca reducir al mínimo el movimiento de la cabeza: *SSTF* (*Shortest Seek Time First, Tiempo de búsqueda más corto a continuación*) es el equivalente en este ámbito del *Proceso más corto a continuación*, presentado en la sección 5.2.4 — con la ventaja de no estar prediciendo comportamiento futuro, sino partir de una lista de solicitudes pendientes. Empleando SSTF, el tiempo de desplazamiento para este caso se reduce a tan sólo 207 cilindros, muy cerca del mínimo absoluto posible.

Una desventaja de SSTF es que puede llevar a la inanición: Si hay una gran densidad de solicitudes para cilindros en determinada zona del disco, una solicitud para un cilindro alejado puede quedar a la espera indefinidamente.

Ejemplificando esto con una serie de solicitudes distinta a la cadena referencia: Si el sistema tuviera que atender solicitudes por los cilindros 15, 175, 13, 20, 14, 32, 40, 5, 6, 7, SSTF penalizaría a la segunda solicitud (175) hasta terminar con los cilindros bajos. Si durante el tiempo que tome responder a estas solicitudes llegan otras adicionales, el proceso que está esperando el contenido del cilindro 175 puede quedar en espera indefinida.

Familia de algoritmos de elevador (SCAN, LOOK, C-SCAN) En este tercer lugar se abordará ya no un sólo algoritmo, sino que una *familia*, dado que parten de la misma idea, pero con modificaciones menores llevan a que el patrón de atención resultante sea muy distinto.

El planteamiento base para el algoritmo básico de elevador (SCAN) busca evitar la inanición, minimizando al mismo tiempo el movimiento de las cabezas. Su lógica indica que la cabeza debe recorrer el disco de extremo a extremo, como si fuera un

elevador en un edificio alto, atendiendo a todas las solicitudes que haya pendientes en su camino. Si bien los recorridos para ciertos patrones pueden resultar en mayores desplazamientos a los que daría SSTF, la garantía de que ningún proceso esperará indefinidamente lo hace muy atractivo.

Atender la cadena de referencia bajo SCAN, asumiendo un estado inicial *descendente* (esto es, la cabeza está en el cilindro 60 y va bajando) da un recorrido total de 235 cilindros; empleando LOOK, se reduce a 205 cilindros, y evita el movimiento innecesario hasta el límite del disco.

Una primer (y casi obvia) modificación a este algoritmo sería, cada vez que la cabeza se detenga para satisfacer una solicitud, verificar si hay alguna otra solicitud pendiente en la *dirección actual*, y de no ser así, emprender el camino de regreso sin llegar a la orilla del disco. Esta modificación es frecuentemente descrita como *LOOK*. Sin embargo, el patrón de atención a solicitudes de SCAN y LOOK dejan qué desear: Al llegar a un extremo del recorrido, es bastante probable que no haya ninguna solicitud pendiente en la primer mitad del recorrido de vuelta (dado que acaban de ser atendidas). El tiempo que demora atender a una solicitud se compone de la suma del desplazamiento de la cabeza y la demora rotacional (que depende de cuál sector del cilindro fue solicitado). Para mantener una tasa de transferencia más predecible, el algoritmo C-SCAN (SCAN Circular) realiza las operaciones en el disco únicamente en un sentido — Si el algoritmo lee en orden *descendente*, al llegar a la solicitud del cilindro más bajo, saltará de vuelta hasta el más alto para volver a iniciar desde ahí. Esto tiene como resultado, claro, que el recorrido total aumente (aumentando hasta los 339 para la cadena de referencia presentada).

Limitaciones de los algoritmos presentados

Ahora bien, ¿por qué se mencionó que estos algoritmos hoy en día ya casi no se usan?

Hay varias razones. En primer término, todos estos algoritmos están orientados a reducir el traslado de *la cabeza*, pero ignoran la *demora rotacional*. Como se explicó, en los discos duros actuales, la demora rotacional va entre $\frac{1}{10}$ y $\frac{1}{3}$ del tiempo total de recorrido de la cabeza. Y si bien el sistema podría considerar esta demora como un factor adicional al planificar el siguiente movimiento de forma que se redujera el tiempo de espera, los algoritmos descritos obviamente requieren ser replanteados por completo.

Por otro lado, el sistema operativo muchas veces requiere dar distintas prioridades a los diferentes tipos de solicitud. Por ejemplo, sería esperable que diera preferencia a los accesos a memoria virtual por encima de las solicitudes de abrir un nuevo archivo. Estos algoritmos tampoco permiten expresar esta necesidad.

Pero el tercer punto es mucho más importante aún: Del mismo modo que los procesadores se van haciendo más rápidos y que la memoria es cada vez de mayor capacidad, los controladores de discos también son cada vez más *inteligentes*, y *esconden* cada vez más información del sistema operativo, por lo cual éste cada vez más carece de la información necesaria acerca del acomodo *real* de la información como para planificar correctamente sus accesos.

Uno de los cambios más importantes en este sentido fue la transición del empleo de la notación C-H-S al esquema de *direcccionamiento lógico de bloques* (*Logical Block Addressing*, LBA) a principios de los noventa. Hasta ese momento, el sistema operativo tenía información de la ubicación *física* de todos los bloques en el disco.

Una de las desventajas, sin embargo, de este esquema es que el mismo BIOS tenía que conocer la *geometría* de los discos — Y el BIOS presentaba límites duros en este sentido: Principalmente, no le era posible referenciar más allá de 64 cilindros. Al aparecer la interfaz de discos IDE (*Electrónica integrada al dispositivo*) e ir reemplazando a la ST-506, se introdujo LBA.

Este mecanismo convierte la dirección C-H-S a una dirección *lineal*, presentando el disco al sistema operativo ya no como un espacio *tridimensional*, sino que como un gran arreglo de bloques. En este primer momento, partiendo de que *CPP* denota el número de cabezas por cilindro y *SPP* el número de sectores por pista, la equivalencia de una dirección C-H-S a una LBA era:

$$LBA = ((C \times CPC) + H) \times SPP + S - 1$$

LBA significó mucho más que una nueva notación: marcó el inicio de la transferencia de inteligencia y control del CPU al controlador de disco. El impacto de esto se refleja directamente en dos factores:

Sectores variables por cilindro En casi todos los discos previos a LBA,¹ el número de sectores por pista se mantenía constante, se tratara de las pistas más internas o más externas. Esto significa que, a igual calidad de la cobertura magnética del medio, los sectores ubicados en la parte exterior del disco desperdiciaban mucho espacio (ya que el *área por bit* era mucho mayor).

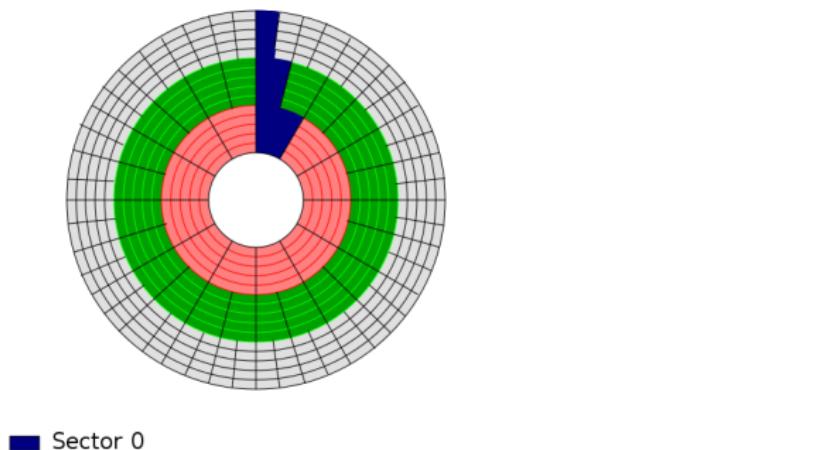


Figura C.3: Disco formateado bajo *densidad de bits por zona*, con más sectores por pista en las pistas exteriores. (Imagen de la Wikipedia: *Zone Bit Recording*)

Bajo LBA, los discos duros comenzaron a emplear un esquema de *densidad de bits por zona* (*zone bit recording*), con la que en los cilindros más externos se aumenta.

Reubicación de sectores Conforme avanza el uso de un disco, es posible que algunos sectores vayan resultando *difíciles* de leer por daños microscópicos a la superficie. El controlador es capaz de detectar estos problemas, y de hecho, casi siempre puede rescatar la información de dichos sectores de forma imperceptible al usuario.

Los discos duros ST-506 típicamente iban acompañados por una *lista de defectos*, una lista de coordenadas C-H-S que desde su fabricación habían presentado errores. El usuario debía ingresar estos defectos al formatear el disco *a bajo nivel*.

Hoy en día, el controlador del disco detecta estos fallos y se los *salta*, presentando un mapa LBA lineal y completo. Los discos duros típicamente vienen con cierto número de *sectores de reserva* para que, conforme se van detectando potenciales daños, estos puedan reemplazarse de forma transparente.

¹Las unidades de disco *Commodore 1541* y {*Macintosh Superdrive*}, que empleaban velocidad variable por cilindro para aprovechar mejor el medio magnético, constituyen notorias excepciones; en ambos casos, sin embargo, terminaron desapareciendo por cuestiones de costos y de complejidad al sistema.

A estos factores se suma que a los controladores de disco se les agregó también una memoria caché dedicada para las operaciones de lectura y escritura. El controlador del disco es hoy en día capaz de implementar estos mismos algoritmos de forma completamente autónoma del sistema operativo.

Y si bien las diferentes unidades de disco duro habían mantenido sectores de 512 bytes desde los primeros discos duros, a partir de la aprobación del *Formato Avanzado* en 2010 que incrementa los sectores a 4096 bytes, presenta otra abstracción más: Un disco con sectores de 4096 bytes que es empleado por el sistema operativo como si fuera de 512² tiene que efectuar, dentro de la lógica de su controlador, una emulación — Y una modificación de un sólo sector se vuelve un *ciclo lectura-modificación-escritura* (RMW), que redunda en una espera de por lo menos una revolución adicional (8ms con un disco de 7200RPM) del disco antes de que la operación pueda completarse.

Resulta claro que, dados estos cambios en la manera en que debe referirse a los bloques del disco, el sistema operativo no cuenta ya con la información necesaria para emplear los algoritmos de planificación de acceso a disco.

C.1.2 Almacenamiento en estado sólido

Desde hace cerca de una década va creciendo consistentemente el uso de medios de almacenamiento de *estado sólido* — Esto es, medios sin partes móviles. Las características de estos medios de almacenamiento son muy distintas de las de los discos.

Si bien las estructuras lógicas que emplean hoy en día prácticamente todos los sistemas de archivos en uso mayoritario están pensadas siguiendo la lógica de los medios magnéticos rotativos, como se verá en esta sección, el empleo de estructuras más acordes a las características del medio físico. Este es indudablemente un área bajo intensa investigación y desarrollo, y que seguramente ofrecerá importantes novedades en los próximos años.

Lo primero que llama la atención de estos medios de almacenamiento es que, a pesar de ser fundamentalmente distintos a los discos magnéticos, se presentan ante el sistema operativo como si fueran lo mismo: En lo que podría entenderse como un esfuerzo para ser utilizados pronto y sin tener que esperar a que los desarrolladores de sistemas operativos adecuaran los controladores, se conectan a través de la misma interfaz y empleando la misma semántica que un disco rotativo.³ Esto no sólo evita que se aprovechen sus características únicas, adoptando restricciones y criterios de diseño que ahora resultan indudablemente artificiales, sino que incluso se exponen a mayor stress por no emplearse de la forma que les resultaría natural.

Antes de ver por qué, conviene hacer un breve repaso de los tipos de discos de estado sólido que hay. Al hablar de la tecnología sobre la cual se implementa este tipo de almacenamiento, los principales medios son:

NVRAM Unidades RAM No Volátil. Almacenan la información en chips de RAM estándar, con un respaldo de batería para mantener la información cuando se desconecta la corriente externa. Las primeras unidades de estado sólido eran de este estilo; hoy en día son poco comunes en el mercado, pero siguen existiendo.

Su principal ventaja es la velocidad y durabilidad: El tiempo de acceso o escritura de datos es el mismo que el que podría esperarse de la memoria principal del sistema, y al no haber demoras mecánicas, este tiempo es el mismo independientemente de la dirección que se solicite.

Su principal desventaja es el precio: En líneas generales, la memoria RAM es, por

²Al día de hoy, los principales sistemas operativos pueden ya hacer referencia al nuevo tamaño de bloque, pero la cantidad de equipos que corren sistemas *heredados* o de controladores que no permiten este nuevo modo de acceso limitan una adopción al 100%

³Las unidades de estado sólido cuentan con una {capa de traducción} que emula el comportamiento de un disco duro, y presenta la misma interfaz tanto de bus como semántica.

volumen de almacenamiento, cientos de veces más cara que el medio magnético. Y si bien el medio no se degrada con el uso, la batería sí, lo que podría poner en peligro a la supervivencia de la información.

Estas unidades típicamente se instalan internamente como una tarjeta de expansión.

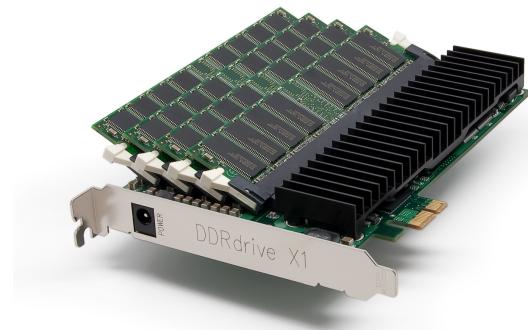


Figura C.4: Unidad de estado sólido basado en RAM: DDRdrive X1 (Imagen de la Wikipedia: *Solid state drive*)

Memoria flash Derivada de los EEPROM (*Electrically Erasable Programmable Read-Only Memory, Memoria de Sólo Lectura Programable y Borrable Eléctricamente*). Los EEPROM tienen la característica de que, además de lectura y escritura, hay un tercer tipo de operación que deben implementar: El *borrado*. Un EEPROM ya utilizado debe borrarse antes de volverse a escribir a él. La principal característica que distingue a las memorias *flash* de los EEPROMs tradicionales es que el espacio de almacenamiento está dividido en muchas *celdas*, y el controlador puede leer, borrar o escribir a cada uno de ellos por separado.⁴

El uso de dispositivos *flash* para almacenamiento de información inició hacia 1995 como respuesta a las necesidades de las industrias aeroespacial y militar, dada la frecuencia de los daños a la información que presentaban los medios magnéticos por la vibración. Hoy en día hay dispositivos *flash* de muy bajo costo y capacidad, aunque presentan una gran variabilidad tanto en su tiempo de acceso como en su durabilidad. En este sentido, existen dos tipos principales de dispositivos *flash*:

Almacenamiento primario (SSD) Las llamadas formalmente *unidad de estado sólido (Solid State Drive)*⁵ son unidades Flash de alta velocidad y capacidad, y típicamente presentan una interfaz similar a la que tienen los discos duros; hoy en día, la más común es SATA.

Su velocidad de lectura es muy superior y su velocidad de escritura (incluyendo el borrado) es comparable a la de los discos magnéticos. Su precio por el mismo volumen de almacenamiento es entre 5 y 10 veces el de los discos magnéticos. Estas unidades se emplean tanto como unidades independientes en servidores, equipos de alto desempeño e incluso algunas subportátiles (*netbooks*) o como un componente de la tarjeta madre en dispositivos móviles como teléfonos y tabletas.

Transporte de archivos Esta tecnología también está presente en las diversas unidades extraíbles o móviles, como las unidades USB, SD, Memory Stick, Compact Flash, etc. La principal diferencia entre estas son los diferentes conectores que

⁴Estos dispositivos se conocen como *flash* en referencia a los chips EPROM (antes de que fuera posible borrar *eléctricamente*): Estos chips tenían una ventana en la parte superior, y debían operar siempre cubiertos con una etiqueta. Para borrar sus contenidos, se retiraba la etiqueta y se les administraba una descarga lumínica — Un *flash*.

⁵Un error muy común es confundir la *D* con *Disk*, que denotaría que llevan un *disco*, un *medio rotativo*



Figura C.5: Unidad de estado sólido basado en Flash con interfaz SATA (Imagen de la Wikipedia: *Solid state drive*)

emplean; todas estas tecnologías presentan dispositivos que varían fuertemente en capacidad, velocidad y durabilidad.



Figura C.6: Unidad de estado sólido basado en Flash con interfaz USB (Imagen de la Wikipedia: *Solid state drive*)

Independientemente del tipo, las unidades de estado sólido presentan ventajas ante los discos rotativos, como un muy bajo consumo eléctrico, operación completamente silenciosa, y resistencia a la vibración o a los golpes. Además, el medio es *verdaderamente* de acceso aleatorio: Al no ser ya un disco, desaparecen tanto la demora de movimiento de cabezas como la rotacional.

Desgaste del medio

La memoria Flash presenta patrones de desgaste muy distintos de los que presentan otros medios. La memoria Flash tiene capacidad de aguantar un cierto número de operaciones de borrado por página⁶ antes de comenzar a degradarse y fallar. Las estructuras tradicionales de sistemas de archivos basados en disco *concentran* una gran cantidad de modificaciones frecuentes a lo largo de la operación normal del sistema en ciertas regiones clave: Las tablas de asignación y directorios registran muchos más cambios que la región de datos.

Casi todos los controladores de discos Flash cuentan con mecanismos de *nivelamiento de escrituras* (*write leveling*). Este mecanismo busca reducir el desgaste focalizado modificando

⁶Dependiendo de la calidad, va entre las 3,000 y 100,000

el mapeo de los sectores que ve el sistema operativo respecto a los que son grabados *en verdad* en el medio: En vez de actualizar un bloque (por ejemplo, un directorio) *en su lugar*, el controlador le asigna un nuevo bloque de forma transparente, y marca el bloque original como libre.

Los mecanismos más simples de nivelamiento de escrituras lo hacen únicamente intercambiando los bloques libres con los recién reescritos; mecanismos más avanzados buscan nivelar el nivel de reescritura en toda la unidad reubicando periódicamente también a los bloques que no son modificados, para no favorecerlos injustamente y hacer un mejor balanceo de uso.

Emulación de discos

Hoy en día, casi la totalidad de medios de estado sólido se presentan ante el sistema con una interfaz que emula la de los discos, la *FTL (Flash Translation Layer, Capa de Traducción de Flash)*. La ventaja de esta emulación es que no hizo falta desarrollar controladores adicionales para comenzar a emplear estos medios. La desventaja, sin embargo, es que al ocultarse el funcionamiento *real* de las unidades de estado sólido, el sistema operativo no puede aprovechar las ventajas estructurales — Y más importante aún, no puede evitar las debilidades inherentes al medio.

Uno de los ejemplos más claros de esta falta de control real del medio la ilustra el [artículo de Valerie Aurora \(2009\)](#), que menciona que tanto la poca información públicamente disponible acerca del funcionamiento de los controladores como los patrones de velocidad y desgaste de los mismos apuntan a que la estructura subyacente de casi todos los medios de estado sólido es la de un *sistema de archivos estructurado en bitácora*. Aurora indica que hay varias operaciones que no pueden ser traducidas eficientemente a través de esta capa de emulación, y que seguramente permitirían un mucho mejor aprovechamiento del medio. Como se mencionó en la sección 8.3.5 (*Sistemas de archivo estructurados en bitácora*), si bien varios de estos sistemas de archivos han presentado implementaciones completamente utilizables, la falta de interés ha llevado a que muchos de estos proyectos sean abandonados.

En su [artículo de 2012](#), Neil Brown apunta a que Linux tiene una interfaz apta para hablar directamente con dispositivos de estado sólido, llamada `mtd` — *memory technology devices, dispositivos de tecnología de memoria*.

Si bien los discos duros se han empleado por ya 50 años y los sistemas de archivos están claramente desarrollados para aprovechar sus detalles físicos y lógicos, el uso de los dispositivos de estado sólido apenas está despegando en la última década. Y si bien esta primer aproximación que permite emplear esta tecnología transparentemente es *suficientemente buena* para muchos de los usos básicos, sin duda hay espacio para mejorar. Este es un tema que seguramente brinda amplio espacio para investigación y desarrollo para los próximos años.

C.2 RAID: Más allá de los límites físicos

En la sección 8.1.1 se presentó muy escuetamente al concepto de *volumen*, mencionando que un volumen *típicamente* coincide con una partición, aunque no siempre es el caso — Sin profundizar más al respecto. En esta sección se presentará uno de los mecanismos que permite combinar diferentes *dispositivos físicos* en un sólo volumen, llevando —bajo sus diferentes modalidades— a mayor confiabilidad, rendimiento y espacio disponible.

El esquema más difundido para este fin es conocido como *RAID, Arreglo Redundante de Discos Baratos (Redundant Array of Inexpensive Disks)*⁷, propuesto en 1988 por David

⁷Ocasionalmente se presenta a RAID como acrónimo de *Arreglo Redundante de Discos Independientes (Redundant Array of Independent Disks)*

Patterson, Garth Gibson y Randy Katz ante el diferencial que se presentaba (y se sigue presentando) entre el avance en velocidad y confiabilidad de las diversas áreas del cómputo en relación al almacenamiento magnético.

Bajo los esquemas RAID queda sobreentendido que los diferentes discos que forman parte de un volumen son del mismo tamaño. Si se remplaza un disco de un arreglo por uno más grande, la capacidad *en exceso* que tenga éste sobre los demás discos será desperdiciada.

Por muchos años, para emplear un *arreglos RAID* era necesario contar con controladores dedicados, que presentaban al conjunto como un dispositivo único al sistema operativo. Hoy en día, prácticamente todos los sistemas operativos incluyen la capacidad de integrar varias unidades independientes en un arreglo por software; esto conlleva un impacto en rendimiento, aunque muy pequeño. Hay también varias tecnologías presentes en distintos sistemas operativos modernos que heredan las ideas presentadas por RAID, pero integrándolos con funciones formalmente implementadas por capas superiores.

RAID no es un sólo esquema, sino que especifica un *conjunto de niveles*, cada uno de ellos diseñado para mejorar distintos aspectos del almacenamiento en discos. Se exponen a continuación las características de los principales niveles en uso hoy en día.

C.2.1 RAID nivel 0: División en *franjas*

El primer nivel de RAID brinda una ganancia tanto en espacio total, dado que presenta a un volumen grande en vez de varios discos más pequeños (simplificando la tarea del administrador) como de velocidad, dado que las lecturas y escrituras al volumen ya no estarán sujetas al movimiento de una sola cabeza, sino que habrá una cabeza independiente por cada uno de los discos que conformen al volumen.

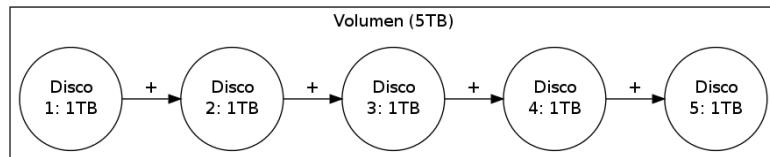


Figura C.7: Cinco discos organizados en RAID 0

Los discos que participan en un volumen RAID 0 no están sencillamente *concatenados*, sino que los datos son *divididos en franjas* (en inglés, el proceso se conoce como *striping*, de la palabra *stripe*, franja; algunas traducciones al español se refieren a este proceso como *bandeado*). Esto hace que la carga se reparta de forma uniforme entre todos los discos, y asegura que todas las transferencias mayores al tamaño de una franja provengan de más de un disco independiente.

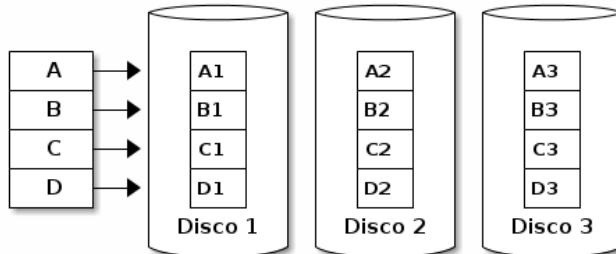


Figura C.8: División de datos en *franjas*

La confiabilidad del volumen, sin embargo, disminuye respecto a si cada uno de los discos se manejara por separado: Basta con que uno de los discos presente daños para que la información contenida en el volumen se pierda.

Un arreglo RAID nivel 0 puede construirse con un mínimo de dos discos.

C.2.2 RAID nivel 1: Espejo

Este nivel está principalmente orientado a aumentar la confiabilidad de la información: Los datos son grabados de forma *simultánea e idéntica* en todos los discos que formen parte del volumen. El costo de mantener los datos en espejo, claro está, es el del espacio empleado: En su configuración habitual, de dos discos por volumen, el 50% del espacio de almacenamiento se pierde por fungir como respaldo del otro 50%.

La velocidad de acceso a los datos bajo RAID 1 es mayor a la que se lograría con un disco tradicional: Basta con obtener los datos de uno de los discos; el controlador RAID (sea el sistema operativo o una implementación en hardware) puede incluso programar las solicitudes de lectura para que se vayan repartiendo entre ambas unidades. La velocidad de escritura se ve levemente reducida, dado que hay que esperar a que ambos discos escriban la información.

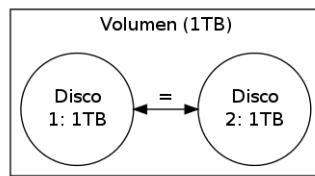


Figura C.9: Dos discos en espejo con RAID 1

Un arreglo RAID nivel 1 se construye típicamente con dos discos.

C.2.3 Los niveles 2, 3 y 4 de RAID

Los siguientes tres niveles de RAID combinan propiedades de los primeros junto con un algoritmo de verificación de integridad y corrección de errores. Estos han caído casi por completo en el desuso dado que los otros niveles, y muy en particular el nivel 5, ofrecen las mismas características, pero con mayor confiabilidad.

C.2.4 RAID nivel 5: Paridad dividida por bloques

El nivel 5 de RAID proporciona un muy buen equilibrio respecto a las características que se han mencionando: brinda el espacio total de almacenamiento de todos los discos que formen parte del volumen *menos uno*. Para cada una de las *franjas*, RAID5 calcula un bloque de *paridad*.

Para obtener una mayor tolerancia a fallos, este bloque de paridad no siempre va al mismo disco, sino que se va repartiendo entre todos los discos del volumen, *desplazándose* a cada franja, de modo que *cualquiera de los discos puede fallar*, y el arreglo continuará operando sin pérdida de información. Esta debe notificarse al administrador del sistema, quien reemplazará al disco dañado lo antes posible (dado que, de no hacerlo, la falla en un segundo disco resultaría en la pérdida de toda la información).

En equipos RAID profesionales es común contar con discos de reserva *en caliente (hot spares)*: Discos que se mantienen apagados pero listos para trabajar. Si el controlador detecta un disco dañado, sin esperar a la intervención del administrador, desactiva al disco afectado y activa al *hot spare*, reconstruyendo de inmediato la información a partir de los datos en los discos *sanos*.

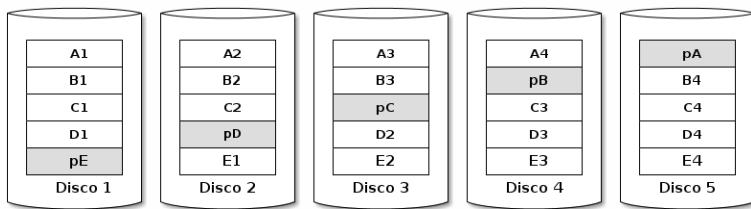


Figura C.10: División de datos en *franjas*, con paridad, para RAID 5

Dependiendo de la configuración, la velocidad de acceso de este nivel puede ser es ligeramente menor que la obtenida de los discos sin RAID, o ligeramente menor a la que se logra con RAID nivel 0. Dado que la electrónica en los discos actuales notificará explícitamente al sistema operativo en caso de fallo de lectura, cuando el sistema requiere leer datos, estos pueden ser solicitados únicamente a $n - 1$ discos (e ignorar al de paridad); si el arreglo RAID está configurado para verificar la paridad en lecturas, todas las lecturas tendrán que obtener la franja correspondiente de todos los discos del arreglo para poder calcularla.

RAID 5 opera con un algoritmo de verificación y recuperación sorprendentemente eficiente y simple: El de una *suma XOR*, ilustrado en la figura C.11. La operación booleana XOR (de *Exclusive OR*) suma los bits individuales, columna por columna. Si es un número par, almacena un 0, si es impar, almacena un 1. Esta operación es muy eficiente computacionalmente.

	Franja A	Franja B
Disco 1	0001 1000	0001 0001
Disco 2	0101 1101	0110 1001
Disco 3	0100 1011	1011 0100
Disco 4	1001 0011	1010 1010
Disco 5	1001 1101	0110 0110

Figura C.11: Para cada franja, el disco de paridad guarda la suma XOR de los bits de las franjas correspondientes de los otros discos; no importa cuál disco falle, sus datos pueden recuperarse haciendo un XOR de los datos de los demás.

Las escrituras son invariablemente más lentas respecto tanto ante la ausencia de RAID como en niveles 0 y 1, dado que siempre tendrá que recalcularse la paridad; en el caso de una escritura mínima (menor a una franja) tendrá que leerse la franja entera de todos los discos participantes en el arreglo, recalcularse la paridad, y grabarse en el disco correspondiente.

Cuando uno de los discos falla, el arreglo comienza a trabajar en el *modo interino de recuperación de datos* (*Interim data recovery mode*, también conocido como *modo degradado*), en el que todas las lecturas involucran a todos los discos, ya que tienen que estar recalculando y *rellenando* la información que provendría del disco dañado.

Para implementar RAID nivel 5 son necesarios por lo menos 3 discos, aunque es común verlos más *anchos*, pues de este modo se desperdicia menos espacio en paridad. Si bien teóricamente un arreglo nivel 5 puede ser arbitrariamente ancho, en la práctica es muy raro ver arreglos con más de 5 discos: Tener un arreglo más ancho aumentaría la probabilidad de falla. Si un arreglo que está ya operando en el modo interino de recuperación de datos se encuentra con una falla en cualquiera de sus discos, tendrá que reportar un fallo irrecuperable.

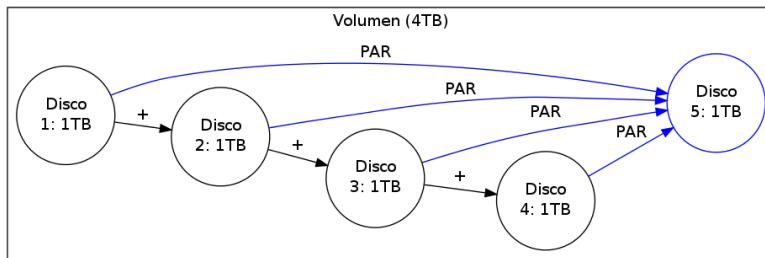


Figura C.12: Cinco discos organizados en RAID 5

C.2.5 RAID nivel 6: Paridad por redundancia P+Q

Se trata nuevamente de un nivel de RAID muy poco utilizado. Se basa en el mismo principio que el de RAID 5 pero, empleando dos distintos algoritmos para calcular la paridad, permite la pérdida de hasta dos de los discos del arreglo. La complejidad computacional es sensiblemente mayor a la de RAID 5, no sólo porque se trata de un segundo cálculo de paridad, sino porque este cálculo debe hacerse empleando un algoritmo distinto y más robusto — Si bien para obtener la paridad P basta con hacer una operación XOR sobre todos los segmentos de una franja, la segunda paridad Q típicamente emplea al *algoritmo Reed-Solomon*, *paridad diagonal* o *paridad dual ortogonal*. Esto conlleva a una mayor carga al sistema, en caso de que sea RAID por software, o a que el controlador sea de mayor costo por implementar mayor complejidad, en caso de ser hardware dedicado.

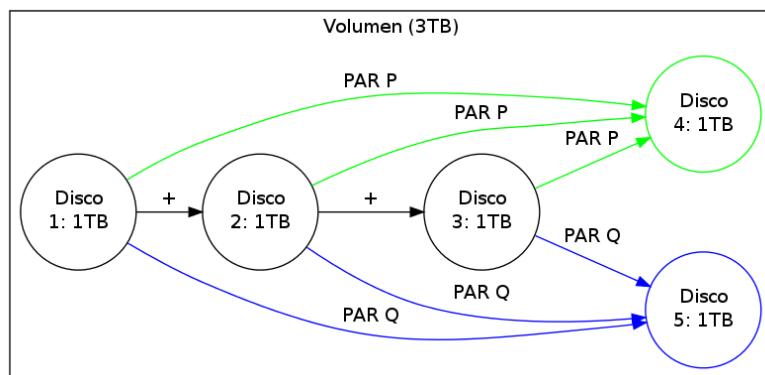


Figura C.13: Cinco discos organizados en RAID 6

El nivel 6 de RAID puede implementarse con 4 o más unidades, y si bien el espacio dedicado a la redundancia se incrementa a dos discos, la redundancia adicional que ofrece este esquema permite crear volúmenes con un mayor número de discos.

C.2.6 Niveles combinados de RAID

Viendo desde el punto de vista de la abstracción presentada, RAID toma una serie de dispositivos de bloques y los *combina* en otro dispositivo de bloques. Esto significa que puede tomarse una serie de volúmenes RAID y combinarlos en uno solo, aprovechando las características de los diferentes niveles.

Si bien pueden combinarse arreglos de todo tipo, hay combinaciones más frecuentes que otras. Con mucho, la más popular es la de los niveles 1 + 0 — Esta combinación, frecuentemente llamada sencillamente *RAID 10*, ofrece un máximo de redundancia y rendimiento, sin sacrificar demasiado espacio.

Con RAID nivel 10 se crean volúmenes que suman por franjas unidades en espejo (un volumen RAID 0 compuesto de varios volúmenes RAID 1). En caso de fallar cualquiera

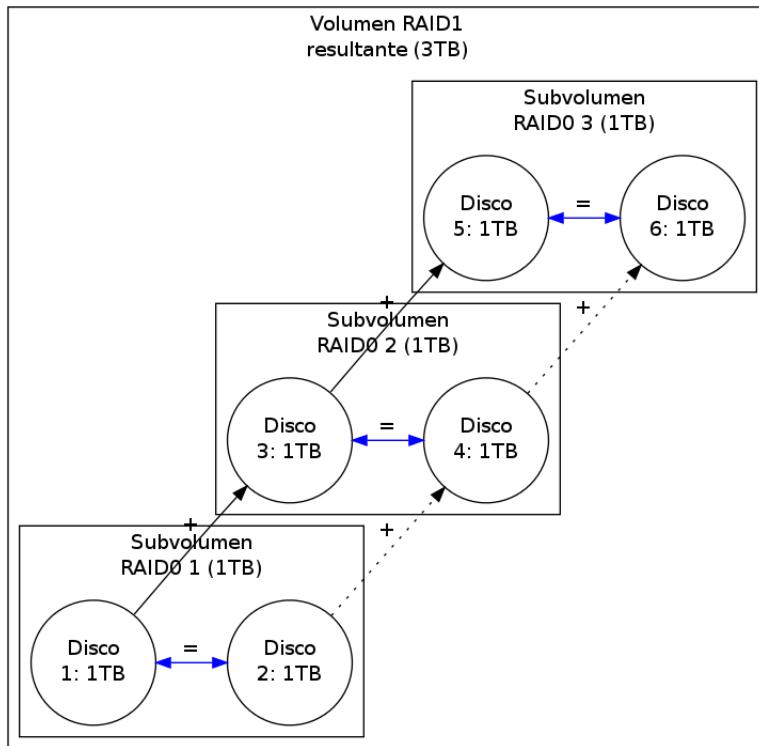


Figura C.14: Seis discos organizados en RAID 1+0

de las unidades del arreglo, ésta puede ser reemplazada fácilmente, y su reemplazo no significará un trabajo tan intensivo para el arreglo entero (sólo para su disco espejo).

Bajo este esquema, en el peor de los casos, un volumen con n discos físicos está conformado por $\frac{n}{2}$ volúmenes nivel 1, y por tanto puede soportar la pérdida de hasta $\frac{n}{2}$ discos — Siempre que estos no formen parte de un mismo volumen nivel 1.

Esta combinación ilustra cómo el orden de los factores *sí altera* al producto: Si en vez de la concatenación de varias unidades espejeadas (un volumen nivel 0 compuesto de varios volúmenes nivel 1) se armara el arreglo en orden inverso (esto es, como el espejo de varias unidades concatenadas por franjas), ante un primer análisis parecería se obtienen los mismos beneficios — Pero analizando lo que ocurre en caso de falla, resulta claro que el nivel de redundancia resulta mucho menor.

En este caso, el arreglo soportará también el fallo de hasta $\frac{n}{2}$ de sus discos, pero únicamente si ocurren *en el mismo volumen RAID 1* del espejo.

Dado que RAID opera meramente agregando *dispositivos de bloques* en un nuevo dispositivo del mismo tipo, no tiene conocimiento de la información subyacente. Por tanto, si se perdieran al mismo tiempo el disco 1 (del subvolumen 1) y el disco 5 (del subvolumen 2), resultaría en pérdida de datos.⁸

C.3 Manejo avanzado de volúmenes

Los esquemas RAID vienen, sin embargo, de fines de la década de 1980, y si bien han cambiado el panorama del almacenamiento, en los más de 20 años desde su aparición, han sido ya superados. El énfasis en el estudio de RAID (y no tanto en los desarrollos posteriores) se justifica dada la limpieza conceptual que presentan, y dado que esquemas

⁸O por lo menos, en una tarea de reconstrucción manual, dada que la información completa existe. Sin embargo, ambos volúmenes RAID 1 estarían dañados e incompletos.

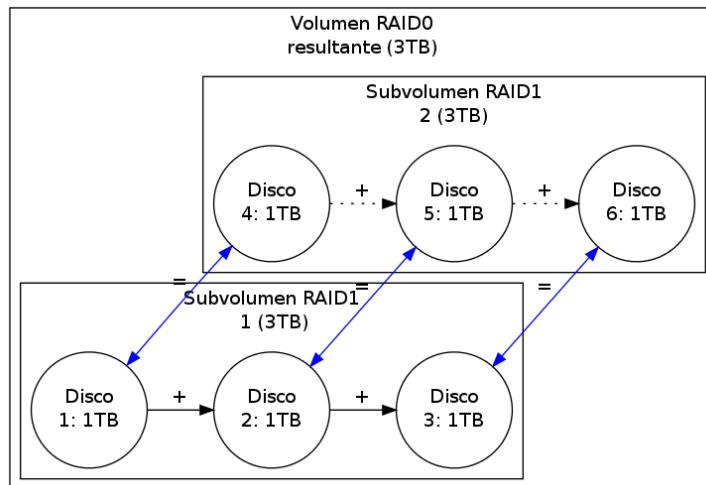


Figura C.15: Seis discos organizados en RAID 0+1

posteriores incluso hacen referencia explícita al nivel de RAID que *estarían reemplazando* en su documentación.

A continuación se presentan brevemente dos esquemas avanzados de gestión de volúmenes, principalmente ilustrando la dirección en que parece ir avanzando la industria en este campo. Dado que no presentan nuevos conceptos sino que sólo ilustran cómo se integran los que se han expuesto en las últimas páginas, la exposición se limitará a presentar ejemplos de aplicación, sin entrar más que a un nivel descriptivo de su funcionamiento.

C.3.1 LVM: el Gestor de Volúmenes Lógicos

Una evolución natural de los conceptos de RAID es el *LVM2* (segunda generación del *Logical Volume Manager*, o *Gestor de Volúmenes Lógicos*) de Linux. La lógica de operación de LVM está basada en los siguientes conceptos:

Volumen físico Cada uno de los discos o unidades disponibles

Grupo de volúmenes Conjunto de volúmenes físicos que serán administrados como una sola entidad

Volumen lógico Espacio dentro del grupo de volúmenes que se presenta como un dispositivo, y que puede alojar sistemas de archivos.

El esquema es limpio y elegante: LVM es una interfaz que permite, como dos pasos independientes, agregar diferentes *volúmenes físicos* a un *grupo de volúmenes*, para posteriormente –y siguiendo las necesidades del administrador del sistema, ya independientes del tamaño de las unidades físicamente existentes– crear las *unidades lógicas*, donde se alojarán los sistemas de archivos propiamente.

Este esquema permite naturalmente una funcionalidad comparable con RAID 0: Puede crearse un grupo de volúmenes con todos los discos que disponibles, y dentro de este crear un volumen lógico único. Dependiendo de la configuración, este volumen lógico puede crecer abarcando todos los discos en cuestión, sea como simple concatenación o dividiéndose en franjas.

Permite también la creación de unidades *espejo*, con una operación a grandes rasgos equivalente a la de RAID1. Incluso, dentro de un mismo grupo de volúmenes, pueden existir tanto volúmenes lógicos espejeados como otros que no lo estén, a diferencia de la estricta rigidez de RAID.

Para los niveles 4, 5 y 6 de RAID, la correspondencia es más directa aún: Al crear un volumen, se le puede solicitar a LVM al crear un volumen lógico que cree un volumen con

ese nivel de RAID — Obviamente, siempre que cuente con suficientes volúmenes físicos.

El esquema de LVM no brinda, pues, funcionalidad estrictamente distinta a la que presenta RAID — Pero da al administrador del sistema flexibilidad: ampliar o reducir el espacio dedicado a cada uno de los volúmenes, incluso en un sistema en producción y con datos.

LVM ofrece varias funcionalidades adicionales, como las *fotografías (snapshots)* o varios esquemas de reemplazo de disco; si bien hay mucho más que podría decirse de LVM, no se profundiza más en esta herramienta dado que excede del objetivo del presente material.

C.3.2 ZFS

Si bien LVM realiza una importante tarea de simplificación en la administración del sistema, su operación sigue siendo orientada a *bloques*: Los volúmenes lógicos deben aún ser formateados bajo el sistema de archivos que el administrador del sistema considere acorde para la tarea requerida.

ZFS⁹ fue desarrollado por Sun Microsystems desde el año 2001, forma parte del sistema operativo *Solaris* desde el 2005, y hoy en día puede emplearse desde los principales sistemas operativos libres.¹⁰ Y si bien ZFS resulta suficientemente atractivo tan sólo por haber sido diseñado para que el usuario nunca más se tope con un límite impuesto por el sistema operativo, el principal cambio que presenta al usuario es una forma completamente distinta de referirse al almacenamiento.

En primer término, al igual que LVM presenta una primer integración entre conceptos, permitiendo unir de diferentes maneras varios dispositivos físicos en un dispositivo lógico, ZFS incluye en la misma lógica administrativa al sistema de archivos: En la configuración estándar, basta conectar una unidad al sistema para que ésta aparezca como espacio adicional disponible para los usuarios. El espacio combinado de todas las unidades conforma un *fondo de almacenamiento (storage pool)*.

La lógica de ZFS parte de que operará una *colección* de sistemas de archivos en una organización jerárquica. Pero a diferencia del esquema tradicional Unix en que cada sistema de archivos es preparado desde un principio para su función, en ZFS se pueden aplicar límites a jerarquías completas. Bajo un esquema ZFS, la creación y el montaje de un sistema de archivos es una operación sencilla — Al grado que se presenta como recomendación que, para cada usuario en el sistema, se genere un sistema de archivos nuevo e independiente.

Una de las principales diferencias con los sistemas de archivos tradicionales es el manejo del espacio vacío: El espacio disponible total del fondo de almacenamiento se reporta como disponible para todos los sistemas de archivos que formen parte de éste. Sin embargo, se pueden indicar *reservas* (mantener un mínimo del espacio especificado disponible para determinado subconjunto de sistemas de archivos dentro de la colección) y *límites* (evitar que el uso de una colección exceda el almacenamiento indicado) para las necesidades de las diferentes regiones del sistema.

C.4 Otros recursos

- *OpenSolaris ZFS Deduplication: Everything You Need to Know*

<http://constantin.glez.de/blog/2010/03/opensolaris-zfs-deduplication-everything-you-need-to-know/>
Constantin Gonzalez (2010)

⁹El nombre ZFS proviene de *Zettabyte File System*. Los diseñadores de ZFS indican, sin embargo, que esto no es para indicar que ZFS sea capaz de direccionar hasta zettabytes de información, sino que será el *último sistema de archivos* que cualquier administrador requerirá.

¹⁰ZFS no puede ser incorporado íntegramente al núcleo de Linux por *incompatibilidad de licencias*: Si bien ambos son software libre, los modelos de licenciamiento GPL (de Linux) y CDDL (de ZFS) son incompatibles.

- *The Z File System (ZFS)*
http://www.freebsd.org/doc/en_US.ISO8859-1/books/handbook/filesystems-zfs.html
FreeBSD Handbook
- *A hash-based DoS attack on Btrfs*
<http://lwn.net/Articles/529077/>
Jonathan Corbet (2012); Linux Weekly News
- *4K Sector Disk Drives: Transitioning to the Future with Advanced Format Technologies*
http://storage.toshiba.com/docs/services-support-documents/toshiba_4kwhitepaper.pdf
Michael E. Fitzpatrick (2011); Toshiba
- *The Design and Implementation of a Log-Structured File System*
<http://www.cs.berkeley.edu/~brewer/cs262/LFS.pdf>
Mendel Rosenblum, J. K. Ousterhout (1992); ACM Transactions on Computer Systems
- *LogFS — Finally a scalable flash file system*
http://www.informatik.uni-osnabrueck.de/papers_pdf/2005_07.pdf
Jörn Engel, Robert Mertens (2005)
- *Log-structured file systems: There's one in every SSD*
<https://lwn.net/Articles/353411/>
Valerie Aurora (2009); Linux Weekly News
- *JFFS2, UBIFS, and the growth of flash storage*
<http://lwn.net/Articles/528617/>
Neil Brown (2012); Linux Weekly News
- *A case for Redundant Arrays of Inexpensive Disks*
<http://www.cs.cmu.edu/~garth/RAIDpaper/Patterson88.pdf>
Patterson, Gibson, Katz (1988); ACM SIGMOD
- *Unidades de estado sólido. El reto de la computación forense en el mundo de los semiconductores*
<http://insecurityit.blogspot.mx/2013/06/>
Cano Martinez (2013); IT-Insecurity
- *Non-volatile memory based on the ferroelectric photovoltaic effect*
<http://dx.doi.org/10.1038/ncomms2990>
Guo, You, Zhou et. al. (2013); Nature Communications
- *eMMC/SSD File System Tuning Methodology*
http://elinux.org/images/b/b6/EMMC-SSD_File_System_Tuning_Methodology_v1.0.pdf
Cogent Embedded (2013)

Índice de figuras

2.1	La <i>microcomputadora Altair 8800</i> , primer computadora personal con distribución masiva, a la venta a partir de 1975. (Imagen de la Wikipedia: <i>Altair 8800</i>)	19
2.2	La <i>Commodore Pet 2001</i> , en el mercado desde 1977, una de las primeras con intérprete de BASIC. (Imagen de la Wikipedia: <i>Commodore PET</i>)	19
2.3	La computadora IBM PC modelo 5150 (1981), iniciadora de la arquitectura predominantemente en uso hasta el día de hoy. (Imagen de la Wikipedia: <i>IBM Personal Computer</i>)	20
2.4	Apple Macintosh (1984), popularizó la interfaz usuario gráfica (GUI). (Imagen de la Wikipedia: <i>Macintosh</i>)	21
2.5	Commodore Amiga 500 (1987), la computadora más popular de la familia <i>Amiga</i> , con amplias capacidades multimedia y multitarea preventiva; una verdadera maravilla para su momento. (Imagen de la Wikipedia: <i>Amiga</i>)	21
2.6	Esquematización de los componentes en un sistema monolítico	23
2.7	Esquematización de los componentes en un sistema microkernel	23
2.8	Esquematización de los componentes en un sistema híbrido	24
3.1	Jerarquía de memoria entre diversos medios de almacenamiento.	26
3.2	Ejemplo de registros: Intel 8086/8088 (Imagen de la Wikipedia: <i>Intel 8086 y 8088</i>)	27
3.3	Diagrama de la comunicación entre componentes de un sistema de cómputo basado en <i>punte norte</i> y <i>punte sur</i>	31
3.4	Esquema simplificado del chipset Intel 875 (para el procesador Pentium 4) ilustrando la velocidad de cada uno de los canales	32
3.5	Transición del flujo entre espacio usuario y espacio núcleo en una llamada al sistema	34
3.6	Esquema de la ejecución de tres procesos en un sistema secuencial, multi-programado, multiprocesado, e híbrido	37
3.7	La <i>Ley de Moore</i> , en su artículo publicado en 1965, prediciendo la miniaturización por diez años	37
3.8	La <i>Ley de Moore</i> se sostiene al día de hoy: conteo de transistores por procesador de 1971 al 2012	38
3.9	Ley de Amdahl: ejecución de un programa con 500 unidades de tiempo total de trabajo con uno, dos y cuatro procesadores.	41
3.10	Ganancia máxima al parallelizar un programa, según la Ley de Amdahl	42
4.1	Diagrama de transición entre los estados de un proceso	46
4.2	Patrón de hilos jefe/trabajador	48
4.3	Patrón de hilos <i>Equipo de trabajo</i>	49
4.4	Patrón de hilos <i>Línea de ensamblado</i>	49

4.5	Sincronización: La exclusión de las <i>secciones críticas</i> entre a varios procesos se protegen por medio de regiones de <i>exclusión mutua</i>	58
4.6	Esquema clásico de un bloqueo mutuo simple: Los procesos <i>A</i> y <i>B</i> esperan mutuamente para el acceso a las unidades de cinta 1 y 2.	77
4.7	Espectro liberal—conservador de esquemas para evitar bloqueos	78
4.8	Evasión de bloqueos: Los procesos <i>A</i> (horizontal) y <i>B</i> (vertical) requieren del acceso exclusivo a un plotter y una impresora, exponiéndose a bloqueo mutuo.	80
4.9	Al emplear categorías de recursos, un ciclo <i>no necesariamente</i> indica un bloqueo	84
4.10	Situación en que se presenta espera circular, incluso empleando categorías de recursos	84
4.11	Detección de ciclos denotando bloqueos: Grafo de procesos y recursos en un momento dado	85
5.1	Planificador a largo plazo	89
5.2	Planificador a mediano plazo, o <i>agendador</i>	90
5.3	Planificador a corto plazo, o <i>despachador</i>	90
5.4	Diagrama de transición entre los estados de un proceso	91
5.5	Ejecución de cuatro procesos con <i>quantums</i> de 5 <i>ticks</i> y cambios de contexto de 2 <i>ticks</i>	93
5.6	Primero llegado, primero servido (FCFS)	95
5.7	Ronda (<i>Round Robin</i>)	96
5.8	Ronda (<i>Round Robin</i>), con $q = 4$	97
5.9	El proceso más corto a continuación (SPN)	98
5.10	Promedio exponencial (predicción de próxima solicitud de tiempo) de un proceso.	98
5.11	Ronda egoísta (SRR) con $a = 2$ y $b = 1$	99
5.12	Representación de un sistema con cinco colas de prioridad y siete procesos listos	100
5.13	Retroalimentación multinivel (FB) básica	101
5.14	Retroalimentación multinivel (FB) con q exponencial	102
5.15	Proporción de penalización registrada por cada proceso contra el porcentaje del tiempo que éste requiere (Finkel, p.33)	105
5.16	Tiempo <i>perdido</i> contra porcentaje de tiempo requerido por proceso (Finkel, p.34)	106
5.17	Mapeo de hilos <i>muchos a uno</i> (Imagen: Beth Plale; ver <i>otros recursos</i>)	106
5.18	Mapeo de hilos <i>uno a uno</i> (Imagen: Beth Plale; ver <i>otros recursos</i>)	107
5.19	Mapeo de hilos <i>muchos a muchos</i> (Imagen: Beth Plale; ver <i>otros recursos</i>)	107
5.20	Descomposición de una instrucción en sus cinco pasos clásicos para organizarse en un <i>pipeline</i>	110
5.21	Alternando ciclos de cómputo y espera por memoria, un procesador que implementa hilos hardware (<i>hyperthreaded</i>) se presenta como dos procesadores	111
6.1	Espacio de direcciones válidas para el proceso 3 definido por un registro base y un registro límite	117
6.2	Patrones de acceso a memoria, demostrando la localidad espacial / temporal (Silberschatz, p.350)	118
6.3	Regiones de la memoria para un proceso	119
6.23	Los picos y valles en la cantidad de fallos de página de un proceso definen a su <i>conjunto activo</i>	140
6.25	Estado de la memoria después del <i>strcpy()</i>	143
6.27	Marco de stack con un <i>canario</i> aleatorio protector de 12 bytes (qR' z2a&5f50s): Si este es sobreescrito por un buffer desbordado, se detendrá la ejecución del programa.	145

6.4	Proceso de compilación y carga de un programa, indicando el tipo de resolución de direcciones (Silberschatz, p.281)	149
6.5	Compactación de la memoria de procesos en ejecución	150
6.6	Ejemplo de segmentación	150
6.7	Página y desplazamiento, en un esquema de direccionamiento de 16 bits y páginas de 512 bytes	150
6.8	Esquema del proceso de paginación, ilustrando el rol de la MMU	151
6.9	Ejemplo (minúsculo) de paginación, con un espacio de direccionamiento de 32 bytes y páginas de 4 bytes	151
6.10	Esquema de paginación empleando un <i>buffer de traducción adelantada</i> (TLB) . .	152
6.11	Paginación en dos niveles: Una tabla externa de 10 bits, tablas intermedias de 10 bits, y marcos de 12 bits (esquema común para procesadores de 32 bits)	152
6.12	Uso de memoria compartida: Tres procesos comparten la memoria ocupada por el texto del programa (azul), difieren sólo en los datos.	152
6.13	Memoria de dos procesos inmediatamente después de la creación del proceso hijo por <code>fork()</code>	153
6.14	Cuando el proceso hijo modifica información en la primer página de su memoria, se crea como una página nueva.	153
6.15	Esquema general de la memoria, incorporando espacio en almacenamiento secundario, representando la memoria virtual	154
6.16	Pasos que atraviesa la respuesta a un fallo de página	155
6.17	Relación ideal entre el número de marcos y fallos de página	155
6.18	Comportamiento del algoritmo FIFO que exhibe la anomalía de Belady al pasar de 3 a 4 marcos. La cadena de referencia que genera este comportamiento es 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5 (Belady, 1969)	156
6.19	Algoritmo FIFO de reemplazo de páginas	156
6.20	Algoritmo óptimo de reemplazo de páginas (OPT)	157
6.21	Algoritmo reemplazo de páginas menos recientemente utilizadas (LRU) . . .	157
7.1	Capas de abstracción para implementar los sistemas de archivos	160
7.2	Archivo de acceso secuencial	166
7.3	Archivo de acceso aleatorio	166
7.4	Acceso relativo a índice: Un índice apuntando al punto justo de un archivo sin estructura	167
7.5	Directorio simple, limitado a un sólo nivel de profundidad	169
7.6	Directorio estucurado en árbol	170
7.7	Directorio como un <i>grafo dirigido acíclico</i> : El directorio <code>proyecto</code> está tanto en el directorio <code>/home/usr1</code> como en el directorio <code>/home/usr2</code>	170
7.8	Directorio como un <i>grafo dirigido</i> , mostrando los <i>enlaces ocultos</i> al directorio actual . y al directorio <i>padre</i>	172
7.9	Directorio basado en grafo dirigido que incluye ciclos	174
7.10	Árbol formado del montaje de <code>sda1</code> en la raíz, <code>sda2</code> como <code>/usr</code> , <code>sdb1</code> como <code>/home</code> , y el directorio virtual <code>proc</code>	175
7.11	Vista de un sistema de archivos Windows	176
8.1	Formato de la entrada del directorio bajo FAT (Mohammed, 2007)	185
8.2	Entradas representando archivo con nombre largo bajo VFAT (Imagen: Peter Clark, ver <i>otros recursos</i>)	186
8.3	Asignación contigua de archivos: Directorio con inicio y longitud	191
8.4	Asignación ligada de archivos: Directorio con apuntador sólo al primer <i>cluster</i> .	192
8.5	Asignación indexada de archivos: Directorio con apuntador al i-nodo (llevando a un i-nodo de tamaño extremadamente ineficiente)	193

8.6 Estructura típica de un i-nodo en Unix, mostrando además el número de accesos a disco necesarios para llegar a cada <i>cluster</i> (con sólo tres <i>cluster</i> por lista)	194
8.7 Ejemplo de entradas en la tabla de asignación de archivos (Imagen: Peter Clark, ver <i>otros recursos</i>)	196
8.8 Inconsistencias en un sistema de archivos tipo FAT	198
8.9 Sistema de archivos con bitácora	201
B.1 Arquitectura de la infraestructura de lenguajes comunes (CLI) de .NET (Imagen de la Wikipedia: <i>Common Language Infrastructure</i>)	216
C.1 Coordenadas de un disco duro, ilustrando su geometría basada en cabeza, cilindro y sector. (Imagen de la Wikipedia: <i>Cilindro Cabeza Sector</i>)	227
C.2 Movimientos de las cabezas bajo los diferentes algoritmos planificadores de acceso a disco, indicando la distancia total recorrida por la cabeza bajo cada uno, iniciando con la cabeza en la posición 60. Para SCAN, LOOK y C-SCAN, se asume que la cabeza inicia avanzando en dirección decreciente.	228
C.3 Disco formateado bajo <i>densidad de bits por zona</i> , con más sectores por pista en las pistas exteriores. (Imagen de la Wikipedia: <i>Zone Bit Recording</i>)	230
C.4 Unidad de estado sólido basado en RAM: DDRdrive X1 (Imagen de la Wikipedia: <i>Solid state drive</i>)	232
C.5 Unidad de estado sólido basado en Flash con interfaz SATA (Imagen de la Wikipedia: <i>Solid state drive</i>)	233
C.6 Unidad de estado sólido basado en Flash con interfaz USB (Imagen de la Wikipedia: <i>Solid state drive</i>)	233
C.7 Cinco discos organizados en RAID 0	235
C.8 División de datos en <i>franjas</i>	235
C.9 Dos discos en espejo con RAID 1	236
C.10 División de datos en <i>franjas</i> , con paridad, para RAID 5	237
C.11 Para cada franja, el disco de paridad guarda la suma XOR de los bits de las franjas correspondientes de los otros discos; no importa cuál disco falle, sus datos pueden recuperarse haciendo un XOR de los datos de los demás.	237
C.12 Cinco discos organizados en RAID 5	238
C.13 Cinco discos organizados en RAID 6	238
C.14 Seis discos organizados en RAID 1+0	239
C.15 Seis discos organizados en RAID 0+1	240



Edición: Marzo de 2014.

Este texto forma parte de la Iniciativa Latinoamericana de Libros de Texto abiertos (LATIn), proyecto financiado por la Unión Europea en el marco de su [Programa ALFA III EuropeAid](#).



Los textos de este libro se distribuyen bajo una Licencia Reconocimiento-CompartirIgual 3.0 Unported (CC BY-SA 3.0) http://creativecommons.org/licenses/by-sa/3.0/deed.es_ES