



Capítulo 12: Indexación y asociación

- Conceptos básicos
- Índices ordenados
- Archivos de índice de árbol B+
- Archivos de índice de árbol B
- Asociación estática
- Asociación dinámica
- Comparación entre indexación ordenada y asociación
- Definición de índice en SQL
- Acceso multiclave
- Mapas de bits



Conceptos básicos

- Mecanismos de indexación empleados para acelerar el acceso a los datos deseados.
 - Por ejemplo, el catálogo de autores en una biblioteca
- **Clave de búsqueda** – atributo, del conjunto de atributos, empleado para buscar registros en un archivo.
- Un **archivo de índices** consta de registros (denominados **entradas de índice**) de la forma

clave de búsqueda	puntero
-------------------	---------

- Los archivos de índices generalmente son más pequeños que el archivo original
- Dos clases básicas de índices:
 - **Índices ordenados:** las claves de búsqueda se almacenan de forma ordenada
 - **Índices asociativos:** las claves de búsqueda están distribuidas uniformemente en “cajones”, empleando una “función de asociación”.



Métricas de evaluación de índices

- Tipos de acceso soportados eficientemente. Por ejemplo,
 - registros con un valor concreto en el atributo
 - o registros con un valor de atributo que se encuentra en un determinado rango de valores.
- Tiempo de acceso
- Tiempo de inserción
- Tiempo de borrado
- Costes de espacio



Índices ordenados

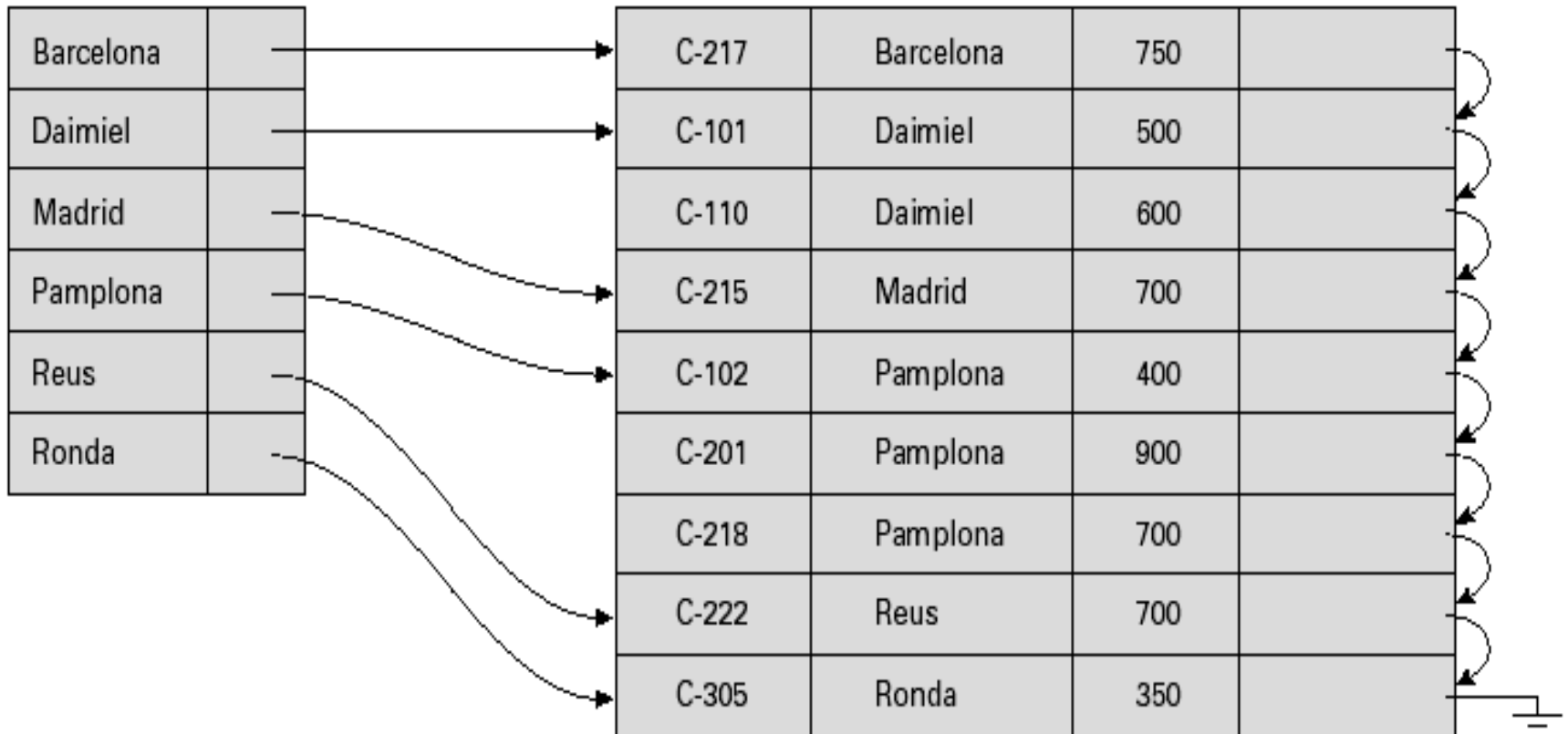
Técnicas de indexado evaluadas en base a:

- En un **índice ordenado**, las entradas de índices se almacenan ordenadas sobre el valor de la clave de búsqueda. Por ejemplo, el catálogo de autores en una biblioteca.
- **Índice primario**: en un archivo ordenado secuencialmente, el índice cuya clave de búsqueda determina el orden secuencial del archivo.
 - También denominado **índice con agrupación (clustering)**
 - La clave de búsqueda de un índice primario es generalmente, pero no necesariamente, la clave primaria.
- **Índice secundario**: un índice cuya clave de búsqueda determina un orden diferente del orden secuencial del archivo. También llamado **índice sin agrupación (non clustering)**.
- **Archivo secuencial indexado**: archivo secuencial ordenado con un índice primario.



Archivos de índice denso

- **Índice denso** — Registro del índice que aparece por cada valor de la clave de búsqueda del archivo.



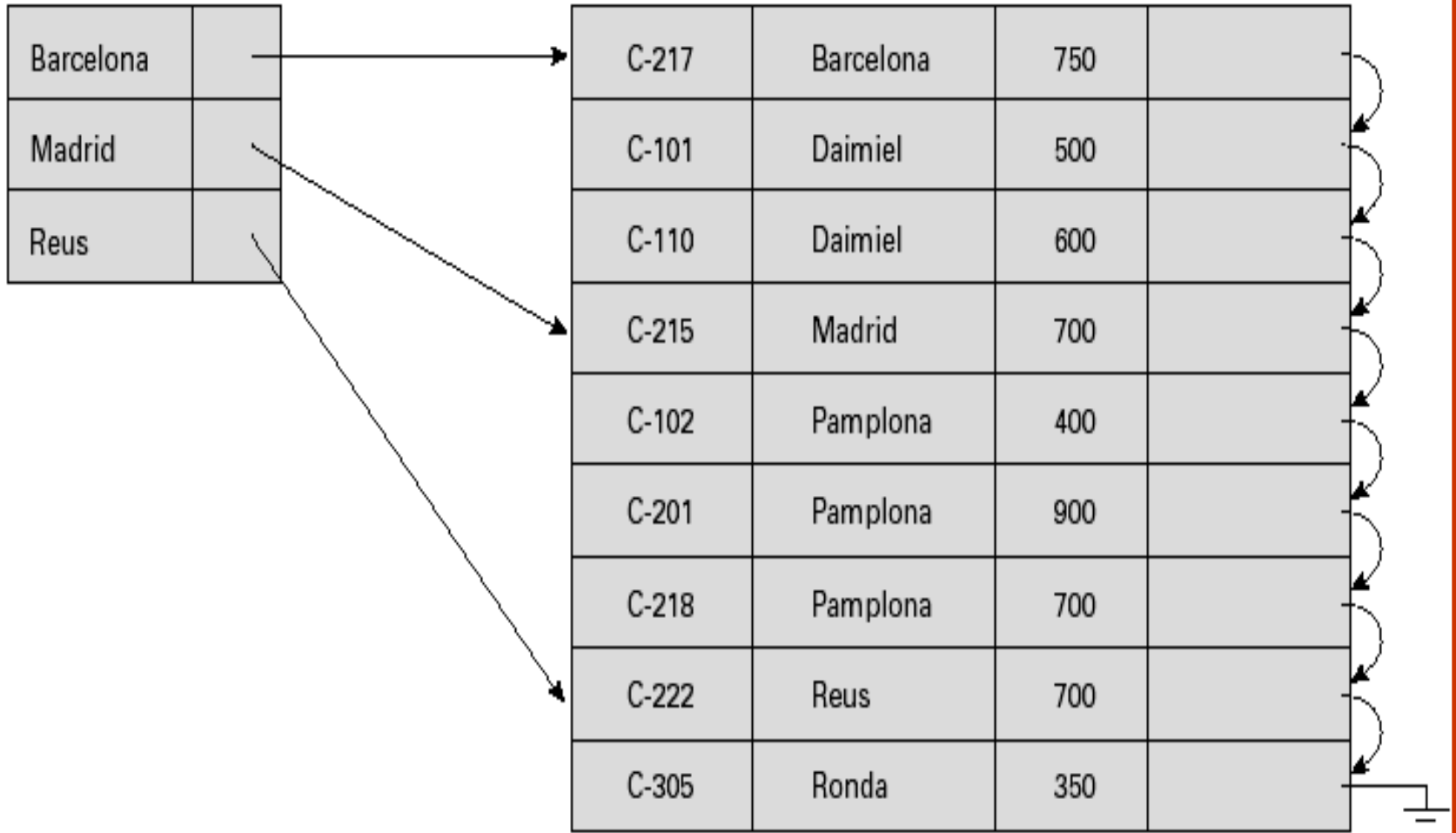


Archivos de índice disperso

- **Índice disperso:** contiene registros del índice, sólo para algunos valores de la clave de búsqueda.
 - Aplicable cuando los registros están ordenados secuencialmente sobre la clave de búsqueda
- Para localizar un registro con valor K de la clave de búsqueda:
 - Encontrar el registro del índice con mayor valor de clave de búsqueda $< K$
 - Búsqueda secuencial del archivo, empezando por el registro al que apunta el registro del índice
- Menos espacio y menores costes de mantenimiento para las inserciones y los borrados.
- Generalmente más lento, para localizar registros, que el índice denso.
- Buen equilibrio: índice disperso con una entrada del índice por cada bloque en el archivo, correspondiente al menor valor de la clave de búsqueda en el bloque.



Ejemplo de archivos de índice disperso

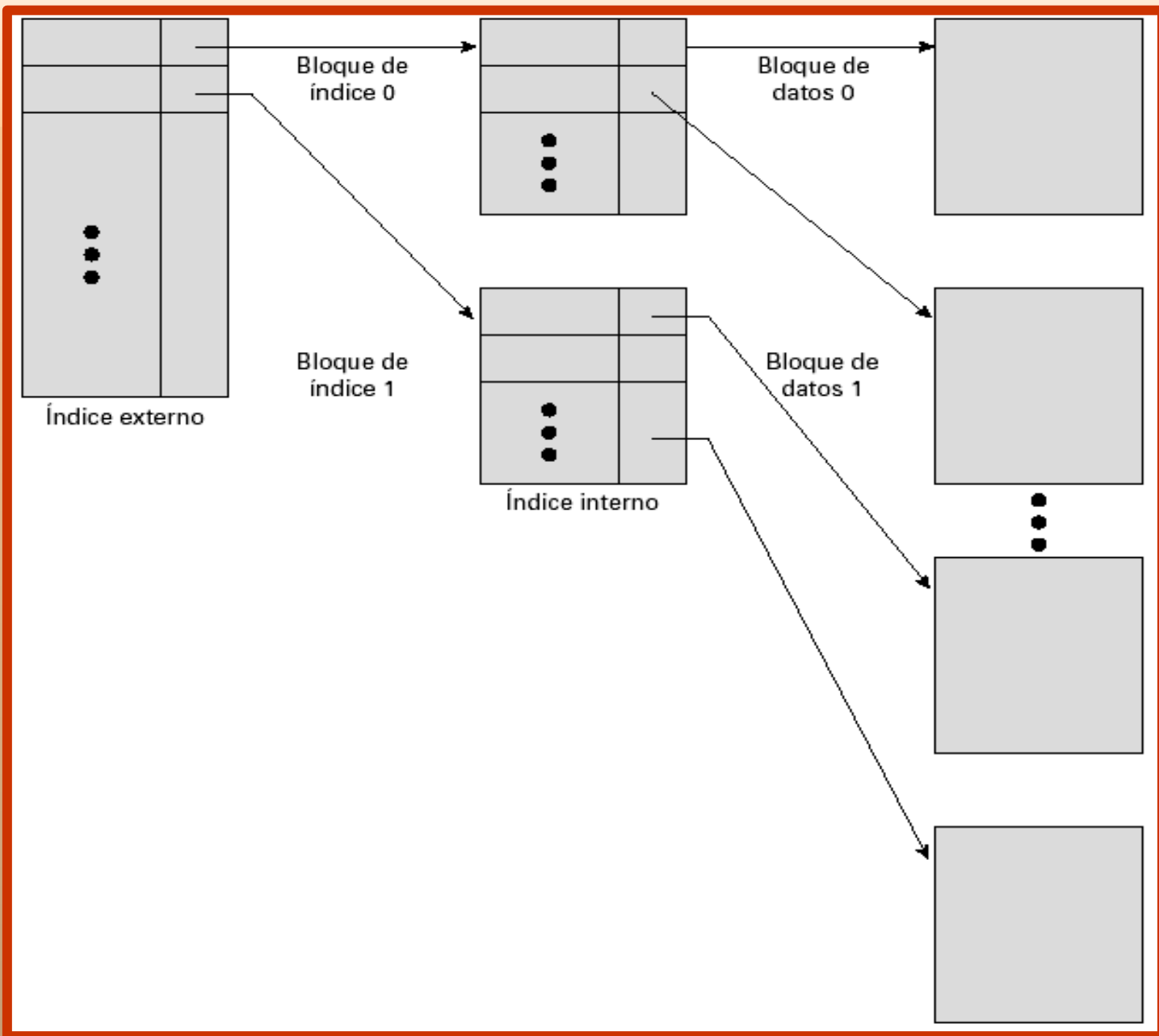




Índice multinivel

- Si el índice primario no encaja en memoria, el acceso se hace costoso.
- Para reducir el número de accesos a disco para los registros del índice, tratar de mantener el índice primario sobre disco como un archivo secuencial y construir un índice disperso en él.
 - índice externo – un índice disperso de índice primario
 - índice interno – el archivo del índice primario
- Si incluso el índice externo es demasiado grande para encajar en memoria principal, aún se puede crear otro nivel de índice, etcétera.
- Al insertar o borrar del archivo, se deben actualizar los índices a todos los niveles.

Índice multinivel (Cont.)





Actualización del índice: Borrado

- Si el registro borrado era el único registro del archivo con su valor de clave de búsqueda concreto, la clave de búsqueda se borra también del índice.
- Borrado del índice de un solo nivel:
 - Índices densos – el borrado de la clave de búsqueda es similar al borrado del registro del archivo.
 - Índices dispersos – si existe una entrada para la clave de búsqueda en el índice se borra, reemplazando la entrada en el índice con el siguiente valor de la clave de búsqueda en el archivo (ordenado por clave de búsqueda). Si el valor de la siguiente clave de búsqueda tiene una entrada del índice, se borra la entrada en vez de reemplazarla.



Actualización del índice: Inserción

- Inserción de índices de un solo nivel:
 - Realizar una búsqueda empleando el valor de la clave de búsqueda que aparece en el registro a insertar.
 - Índices densos – si el valor de la clave de búsqueda no aparece en el índice, insertarlo.
 - Índices dispersos – si el índice almacena una entrada por cada bloque del archivo, no es necesario hacer ningún cambio al índice, a menos que se cree un nuevo bloque. En este caso, se inserta en el índice el primer valor de la clave de búsqueda en el nuevo bloque.
- Los algoritmos de inserciones multinivel (así como en el borrado) son simples extensiones de los algoritmos de un solo nivel

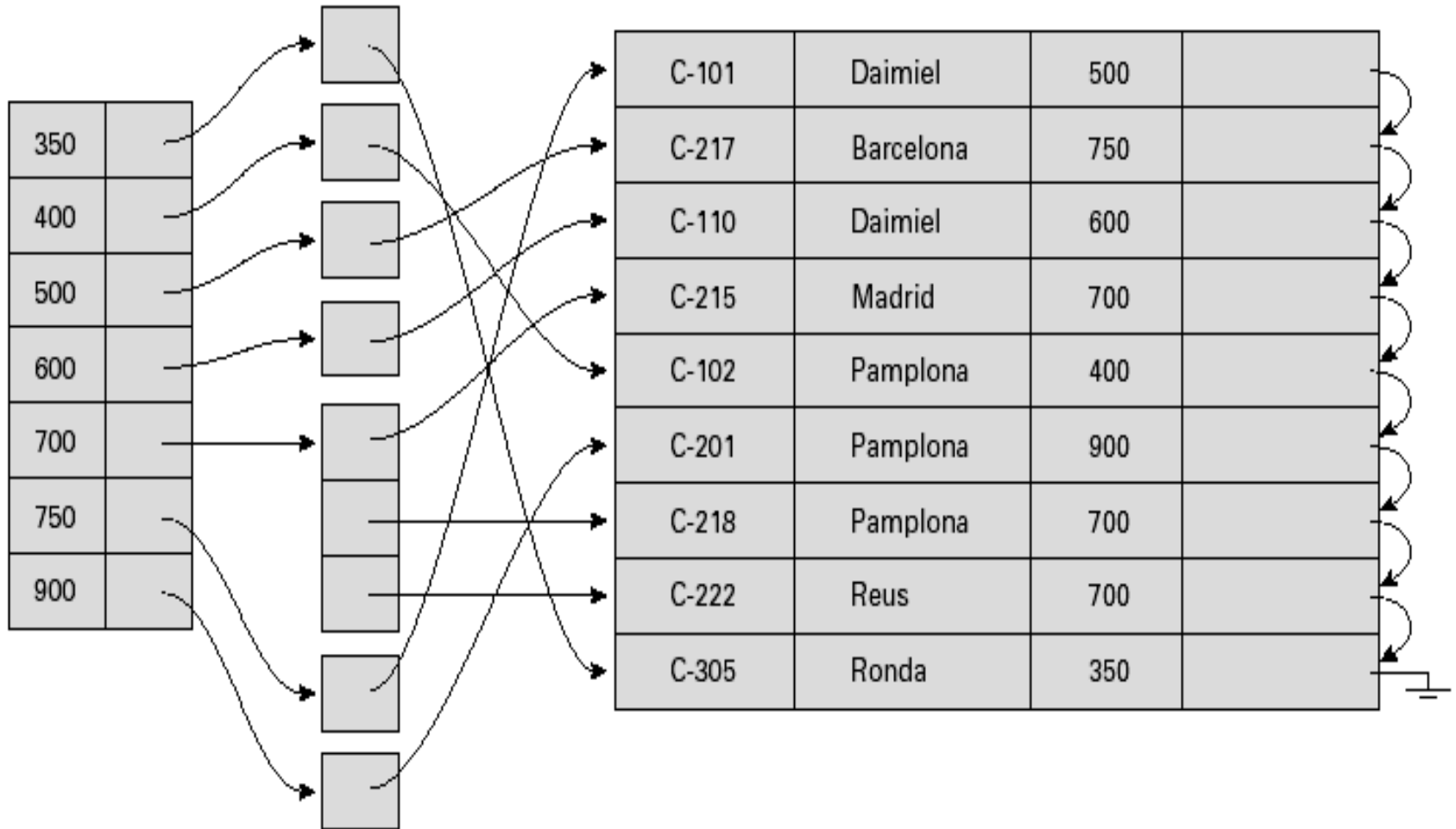


Índices secundarios

- Frecuentemente, se quieren encontrar todos los registros cuyos valores en un cierto campo (que no es la clave de búsqueda del índice primario) cumplen alguna condición.
 - Ejemplo 1: En la base de datos *cuenta*, almacenada secuencialmente por número de cuenta, se pueden encontrar fácilmente todas las cuentas de una determinada oficina
 - Ejemplo 2: como antes, pero donde se quieren encontrar todas las cuentas con un determinado saldo o rango de saldos
- Se puede tener un índice secundario con un registro del índice por cada valor de la clave de búsqueda; el registro del índice apunta a un cajón que contiene punteros a todos los registros actuales, con ese valor particular de clave de búsqueda.



Índice secundario sobre el campo *saldo de cuenta*





Índices primario y secundario

- Los índices secundarios han de ser densos.
- Los índices ofrecen importantes ventajas en la búsqueda de registros.
- Cuando se modifica un archivo, se debe actualizar cada índice del archivo; Actualizar los índices implica sobrecargas en la modificación de la base de datos.
- La búsqueda secuencial empleando índices primarios es eficiente, pero utilizando un índice secundario es costosa
 - cada acceso de registro puede coger un nuevo bloque del disco



Archivos de índice de árbol B+

Los índices de árbol B⁺ son una alternativa a los archivos secuenciales indexados.

- Inconvenientes de los archivos secuenciales indexados: el rendimiento baja cuando el archivo crece, dado que se crean muchos bloques de desbordamiento. Es necesario reorganizar periódicamente todo el archivo.
- Ventajas de los archivos de índice de árbol B⁺: se reorganiza automáticamente por sí mismo con pequeños cambios locales, a pesar de las inserciones y los borrados. No es necesario reorganizar todo el archivo para mantener el rendimiento.
- Inconvenientes de los árboles B⁺: inserciones extras, sobrecarga de borrados y costes de espacio.
- Las ventajas de los árboles B⁺ superan a los inconvenientes, por lo que se emplean ampliamente.



Archivos de índice de árbol B+ (Cont.)

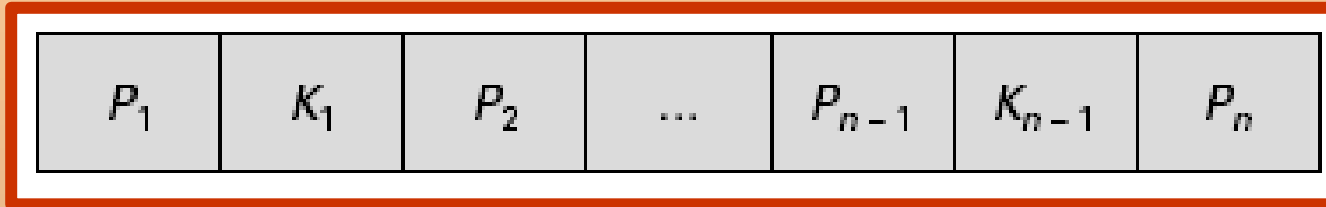
Un árbol B⁺ es un árbol con raíz que satisface las siguientes propiedades:

- Todos los caminos, desde la raíz a las hojas, tienen la misma longitud
- Cada nodo que no es ni raíz ni hoja, tiene entre $\lceil n/2 \rceil$ y n hijos.
- Un nodo hoja tiene entre $\lceil (n-1)/2 \rceil$ y $n-1$ valores
- Casos especiales:
 - Si la raíz no es una hoja, tiene al menos 2 hijos.
 - Si la raíz es una hoja (es decir, no hay otros nodos en el árbol), puede tener entre 0 y $(n-1)$ valores.



Estructura de nodos del árbol B+

■ Nodo típico



- K_i son los valores de la clave de búsqueda
- P_i son los punteros a los hijos (para nodos no hoja) o a los registros o cajones de registros (para nodos hoja).

■ En un nodo las claves de búsqueda están ordenadas

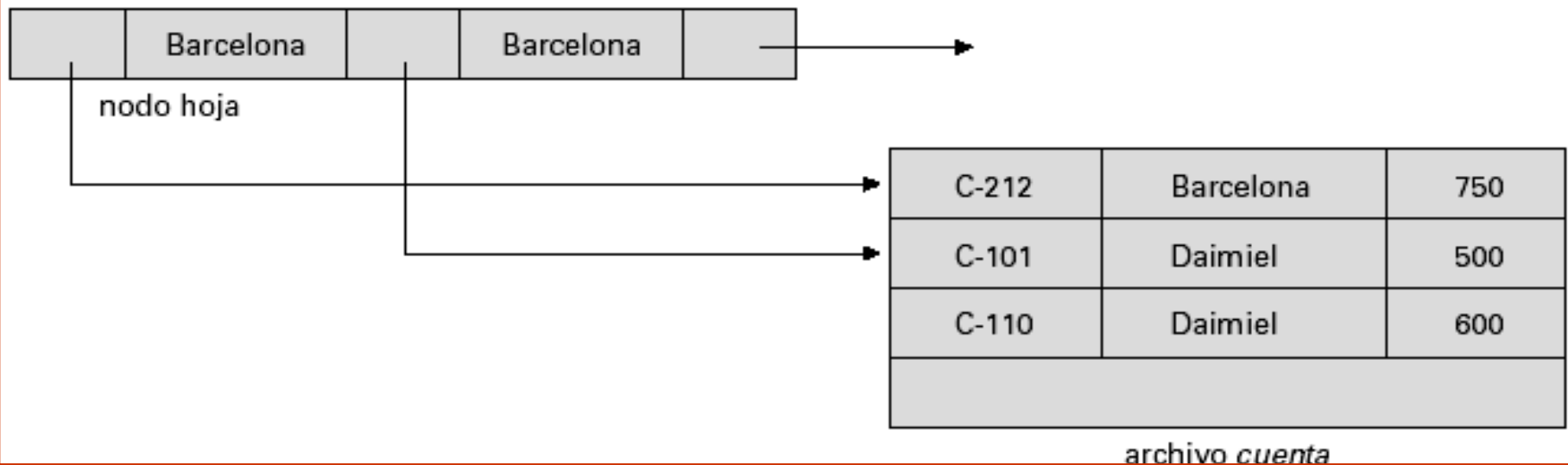
$$K_1 < K_2 < K_3 < \dots < K_{n-1}$$



Nodos hoja en árboles B⁺

Propiedades de un nodo hoja:

- Para $i = 1, 2, \dots, n-1$, el puntero P_i apunta a un registro del archivo con valor de clave de búsqueda K_i , o a un cajón de punteros a los registros del archivo, cada registro con valor de clave de búsqueda K_i . Sólo es necesaria una estructura de cajones si la clave de búsqueda no es una clave primaria.
- Si L_i, L_j son nodos hoja e $i < j$, los valores de clave de búsqueda de L_i son menores que los de L_j
- P_n apunta al siguiente nodo hoja, ordenado por clave de búsqueda





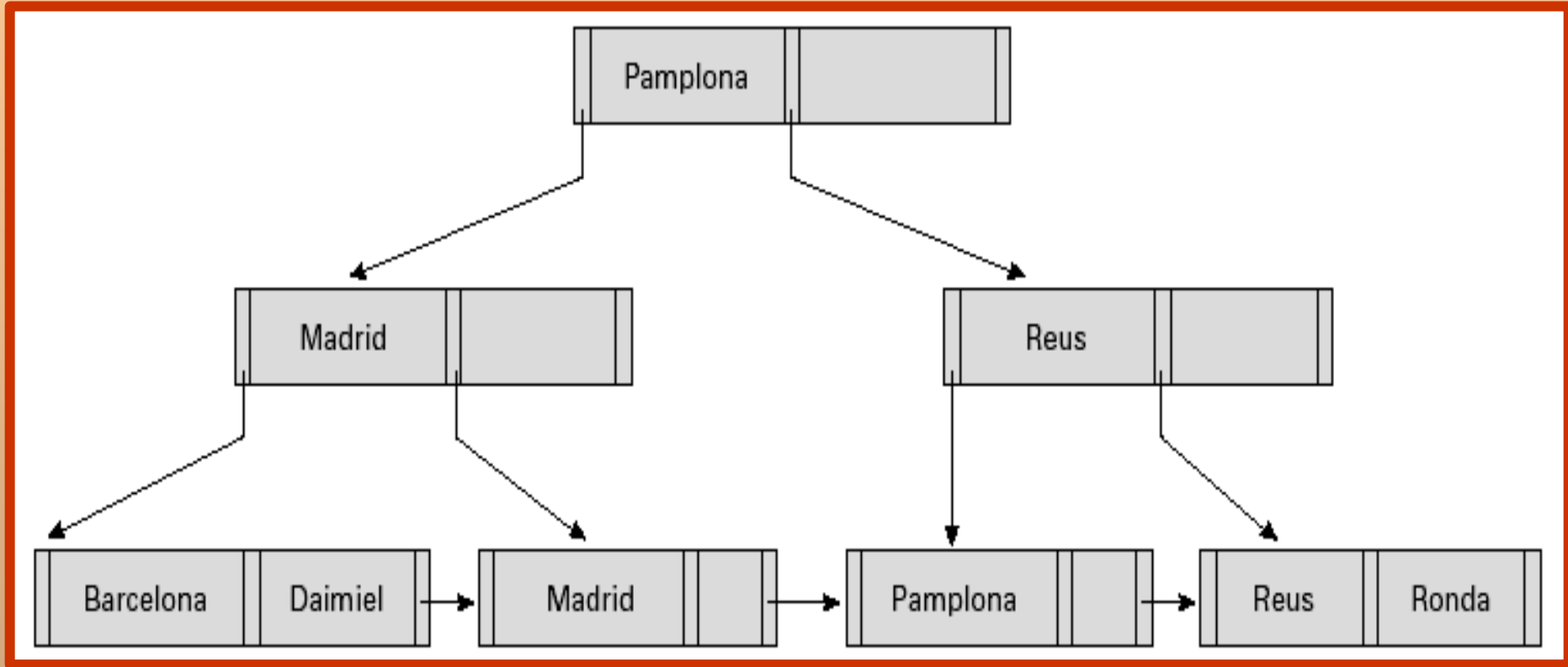
Nodos que no son hoja en árboles B⁺

- Los nodos que no son hoja forman un índice disperso multinivel sobre los nodos hoja. Para un nodo no hoja con m punteros:
 - Todas las claves de búsqueda en el subárbol al que apunta P_1 son menores que K_1
 - Para $2 \leq i \leq n-1$, todas las claves de búsqueda en el subárbol al que apunta P_i tienen valores mayores o iguales que K_{i-1} y menores que K_{m-1}





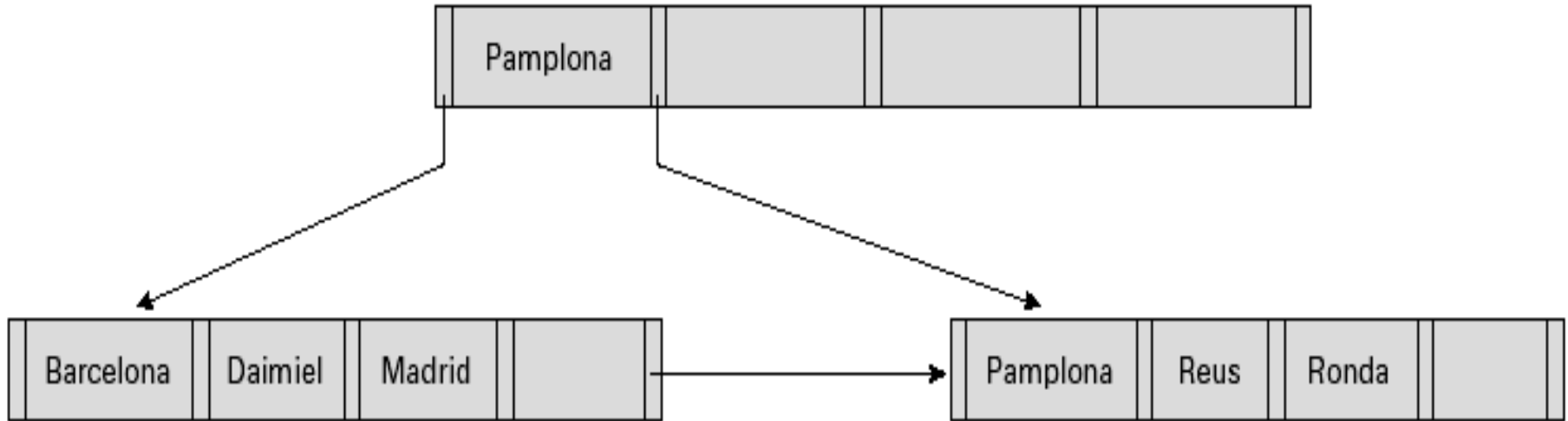
Ejemplo de un árbol B⁺



Árbol B⁺ del archivo *cuenta* ($n = 3$)



Ejemplo de árbol B⁺



Árbol B⁺ del archivo *cuenta* ($n = 5$)

- Los nodos hoja deben tener entre 2 y 4 valores ($\lceil (n-1)/2 \rceil$ y $n-1$, con $n = 5$).
- Otros nodos que no son hoja, excepto la raíz, deben tener entre 3 y 5 hijos ($\lceil n/2 \rceil$ y n con $n = 5$).
- La raíz debe tener al menos 2 hijos.



Observaciones sobre los árboles B⁺

- Dado que las conexiones entre nodos se hacen mediante punteros, los bloques “lógicamente” cercanos no necesitan estar “físicamente” próximos.
- Los niveles que no son hoja del árbol B⁺ forman una jerarquía de índices dispersos.
- El árbol B⁺ contiene un número relativamente pequeño de niveles (logarítmicos en el tamaño del archivo principal), por lo que las búsquedas se pueden realizar de forma eficiente.
- Las inserciones y los borrados sobre el archivo principal se pueden gestionar eficientemente, del mismo modo que el índice se puede reconstruir de forma logarítmica (como se verá).



Asociación estática

- Un **cajón** es una unidad de almacenamiento que contiene uno o más registros (generalmente un cajón es un bloque de disco).
- En una **organización de archivos asociativa** se obtiene el cajón de un registro directamente desde su valor de clave de búsqueda, empleando una **función de asociación**.
- La función de asociación h es una función desde el conjunto de todos los valores de claves de búsqueda K , hasta el conjunto de todas las direcciones de cajones B .
- La función de asociación se emplea para localizar registros para accesos, inserciones y borrados.
- Los registros con diferentes valores de claves de búsqueda pueden asociarse al mismo cajón; de este modo, para localizar un registro, se ha de recorrer secuencialmente el cajón entero.



Ejemplo de organización de archivos asociativa (Cont.)

La organización de archivos asociativa del archivo *cuenta*, utilizando *nombre-sucursal* como clave (véase la figura de la siguiente transparencia).

- Hay 10 cajones,
- La representación binaria del carácter i ésimo se asume que sea el entero i .
- La función de asociación devuelve la suma de las representaciones binarias de los caracteres módulo 10
 - Por ejemplo, $h(\text{Navacerrada}) = 5$ $h(\text{Ronda}) = 3$ $h(\text{Galapagar}) = 3$



Ejemplo de organización de archivos asociativa

Cajón 0

--	--	--

Cajón 1

--	--	--

Cajón 2

--	--	--

Cajón 3

C-217	Barcelona	750
C-101	Daimiel	500

Cajón 4

C-222	Reus	700

Cajón 5

C-102	Pamplona	400
C-201	Pamplona	900
C-218	Pamplona	700

Cajón 6

--	--	--

Cajón 7

C-215	Madrid	700

Cajón 8

C-305	Ronda	350
C-110	Daimiel	600

Cajón 9

--	--	--

La organización de archivos asociativa del archivo *cuenta* utilizando *nombre-sucursal* como clave (para obtener más detalles, véase la transparencia anterior).



Funciones de asociación

- La peor función de asociación asocia todos los valores de las claves de búsqueda al mismo cajón; esto hace que el tiempo de acceso sea proporcional al número de valores de claves de búsqueda en el archivo.
- Una función de asociación ideal es **uniforme**, es decir, cada cajón se asigna al mismo número de valores de claves de búsqueda, desde el conjunto de *todos* los valores posibles.
- La función de asociación ideal es **random**; así cada cajón tendrá asignado el mismo número de registros, independientemente de la *distribución* real de los valores de las claves de búsqueda en el archivo.
- Las funciones de asociación típicas realizan cálculos sobre la representación binaria interna de la clave de búsqueda.
 - Por ejemplo, para una clave de búsqueda de secuencia de caracteres, se podrían añadir las representaciones binarias de todos los caracteres en la secuencia y se podría devolver la suma del número de cajones.



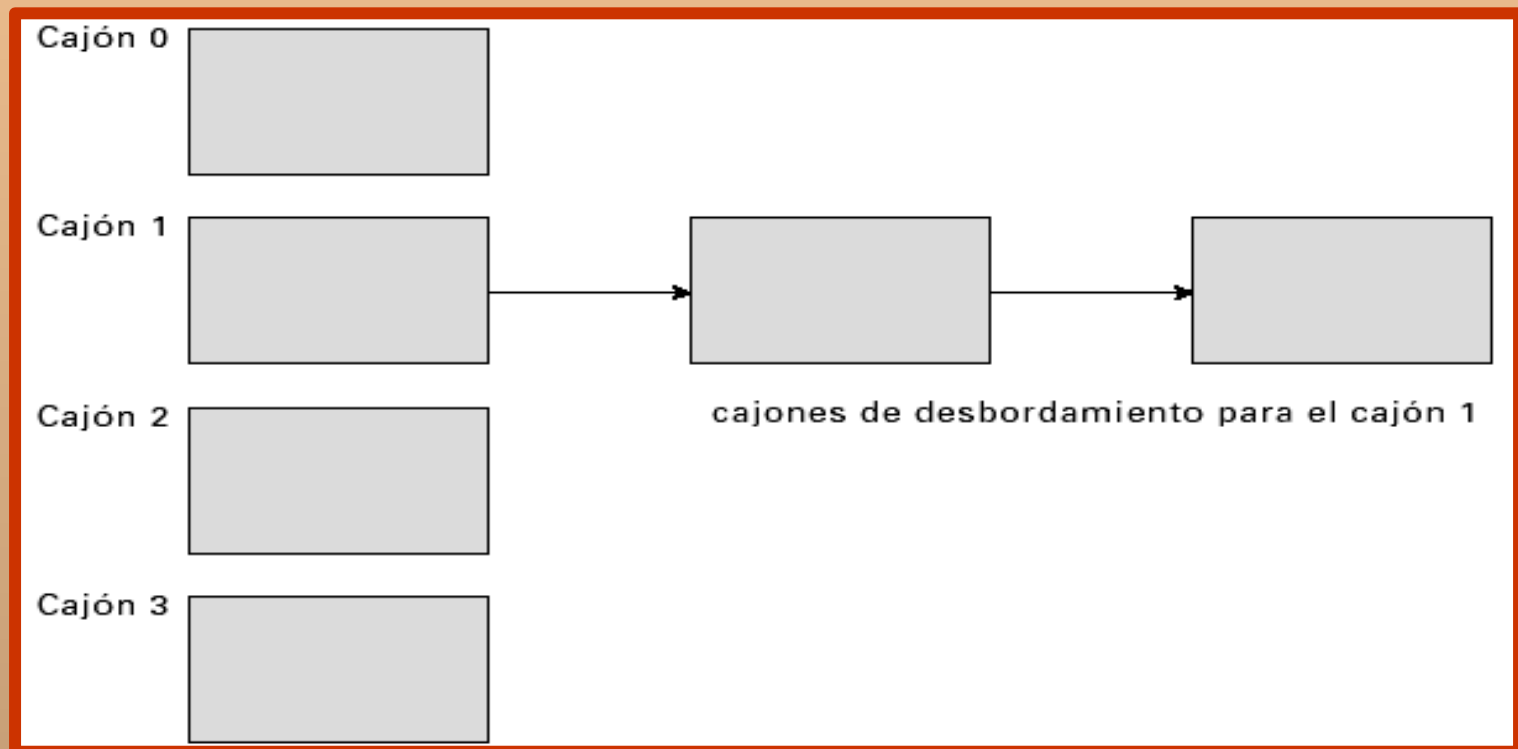
Gestión de desbordamiento de cajones

- El desbordamiento de cajones puede producirse por
 - Insuficientes cajones
 - Desviación en la distribución de los registros. Esto puede tener lugar por dos razones:
 - t múltiples registros tienen el mismo valor de clave de búsqueda
 - la función de asociación elegida produce una distribución no uniforme de los valores de las claves
- Aunque se puede reducir la probabilidad de desbordamiento de cajones, no se puede eliminar y se gestiona empleando *cajones de desbordamiento*.



Gestión de desbordamiento de cajones (Cont.)

- **Cadena de desbordamiento** – los cajones de desbordamiento de un determinado cajón se encadenan juntos en una lista enlazada.
- El esquema anterior se denomina **asociación cerrada**.
 - Una alternativa, denominada **asociación abierta**, que no emplea cajones de desbordamiento, no es adecuada para aplicaciones de bases de datos.





Deficiencias de la asociación estática

- En la asociación estática la función h asocia valores de claves de búsqueda a un determinado conjunto B , de direcciones de cajones.
 - Las bases de datos crecen con el tiempo. Si el número inicial de cajones es demasiado pequeño, disminuirá el rendimiento debido a los muchos desbordamientos.
 - Si se anticipa el tamaño del archivo en algún momento del futuro, y en consecuencia el número de cajones asignados, un aumento significativo de espacio se desperdiciará inicialmente.
 - Si disminuye la base de datos, nuevamente se desperdiciará espacio.
 - Una opción es la reorganización periódica del archivo con una nueva función de asociación, pero es muy costoso.
- Estos problemas se pueden evitar empleando técnicas que permitan que el número de cajones se modifique dinámicamente.

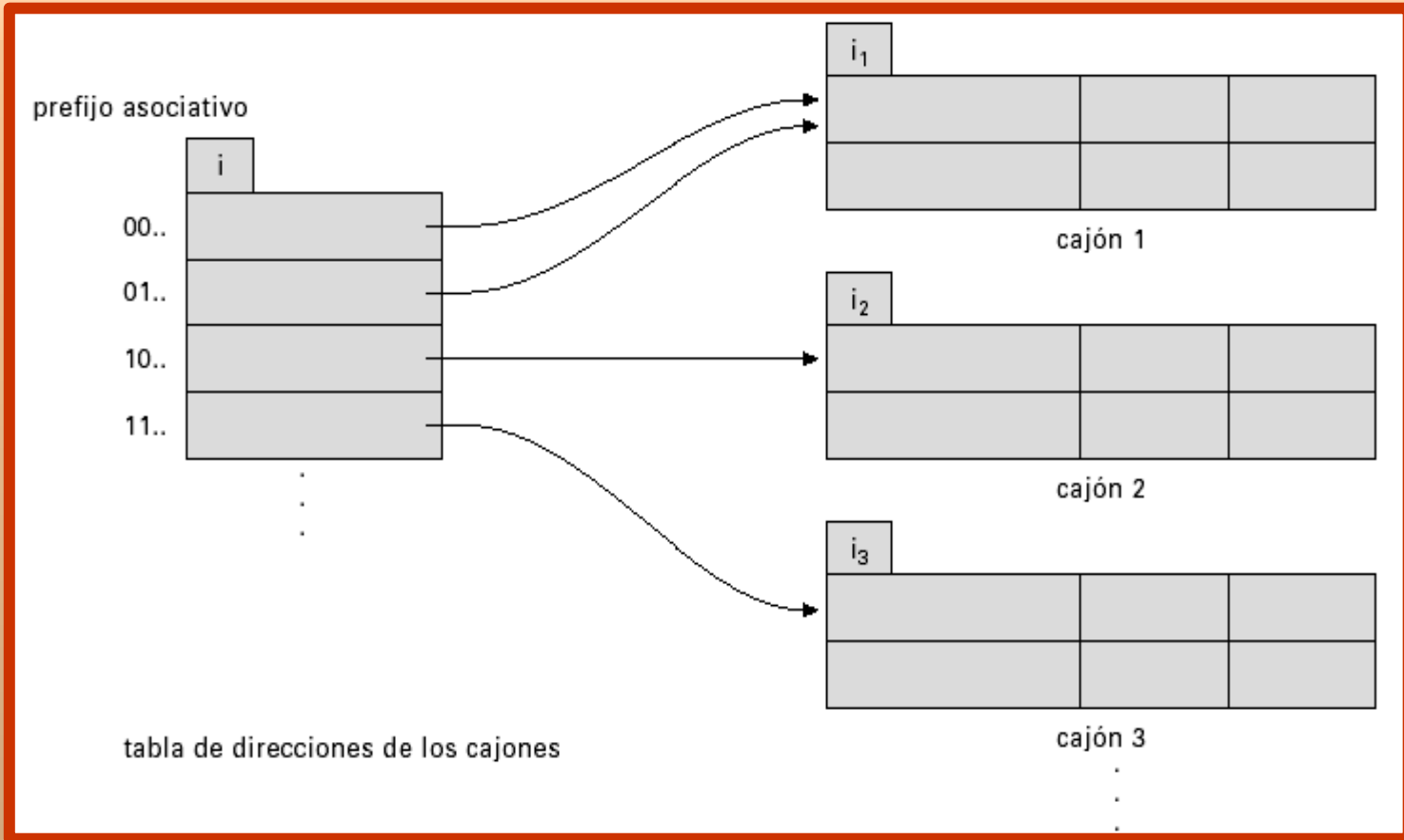


Asociación dinámica

- Buena para las bases de datos que aumentan y disminuyen de tamaño
- Permite modificar dinámicamente la función de asociación
- **Asociación extensible** – una forma de asociación dinámica
 - La función de asociación genera valores en un amplio rango — generalmente b -bit enteros, con $b = 32$.
 - En cualquier momento se emplea sólo un prefijo de la función de asociación, para indexar en una tabla de direcciones de cajones.
 - Sea la longitud del prefijo i bits, donde $0 \leq i \leq 32$.
 - El tamaño de la tabla de direcciones de los cajones es $= 2^i$. Inicialmente $i = 0$
 - El valor de i crece y disminuye según lo hace el tamaño de la base de datos.
 - Múltiples entradas en la tabla de direcciones de los cajones pueden apuntar a un cajón.
 - Así, el número real de cajones es $< 2^i$
 - El número de cajones también cambia dinámicamente debido a las agrupaciones y divisiones de los cajones.



Estructura de asociación extensible general



En esta estructura, $i_2 = i_3 = i$, mientras que $i_1 = i - 1$
(para obtener más detalles, véase la siguiente
transparencia)



Uso de la estructura de asociación extensible: Ejemplo

<i>nombre-sucursal</i>	<i>h(nombre-sucursal)</i>
Barcelona	0010 1101 1111 1011 0010 1100 0011 0000
Daimiel	1010 0011 1010 0000 1100 0110 1001 1111
Madrid	1100 0111 1110 1101 1011 1111 0011 1010
Pamplona	1111 0001 0010 0100 1001 0011 0110 1101
Reus	0011 0101 1010 0110 1100 1001 1110 1011
Ronda	1101 1000 0011 1111 1001 1100 0000 0001

prefijo asociativo

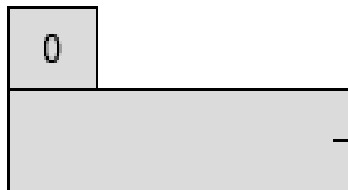
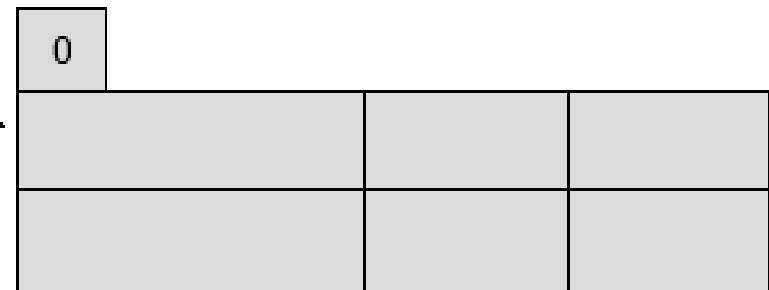


tabla de direcciones de los cajones



cajón 1

Estructura asociativa inicial, tamaño del cajón = 2



Ejemplo (Cont.)

- Estructura de asociación después de la inserción de un registro “Galapagar” y dos “Centro”

prefijo asociativo

1

tabla de direcciones de los cajones

1		
C-217	Barcelona	750

1		
C-101	Daimiel	500
C-110	Daimiel	600



Ejemplo (Cont.)

Estructura de asociación después de la inserción del registro “Madrid”.

prefijo asociativo

2

tabla de direcciones de los cajones

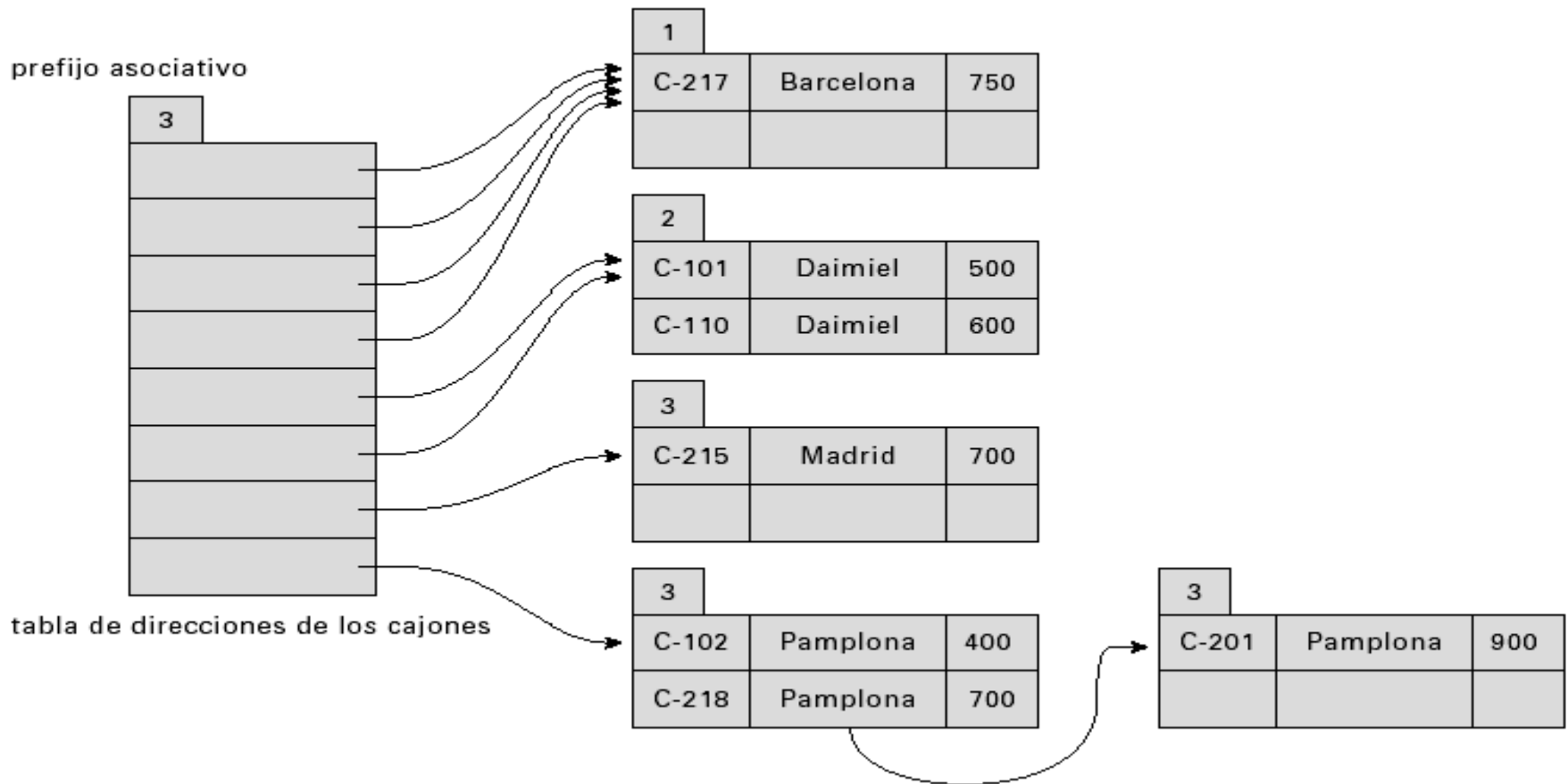
1		
C-217	Barcelona	750

2			
C-101	Daimiel	500	
C-110	Daimiel	600	

2		
C-215	Madrid	700



Ejemplo (Cont.)

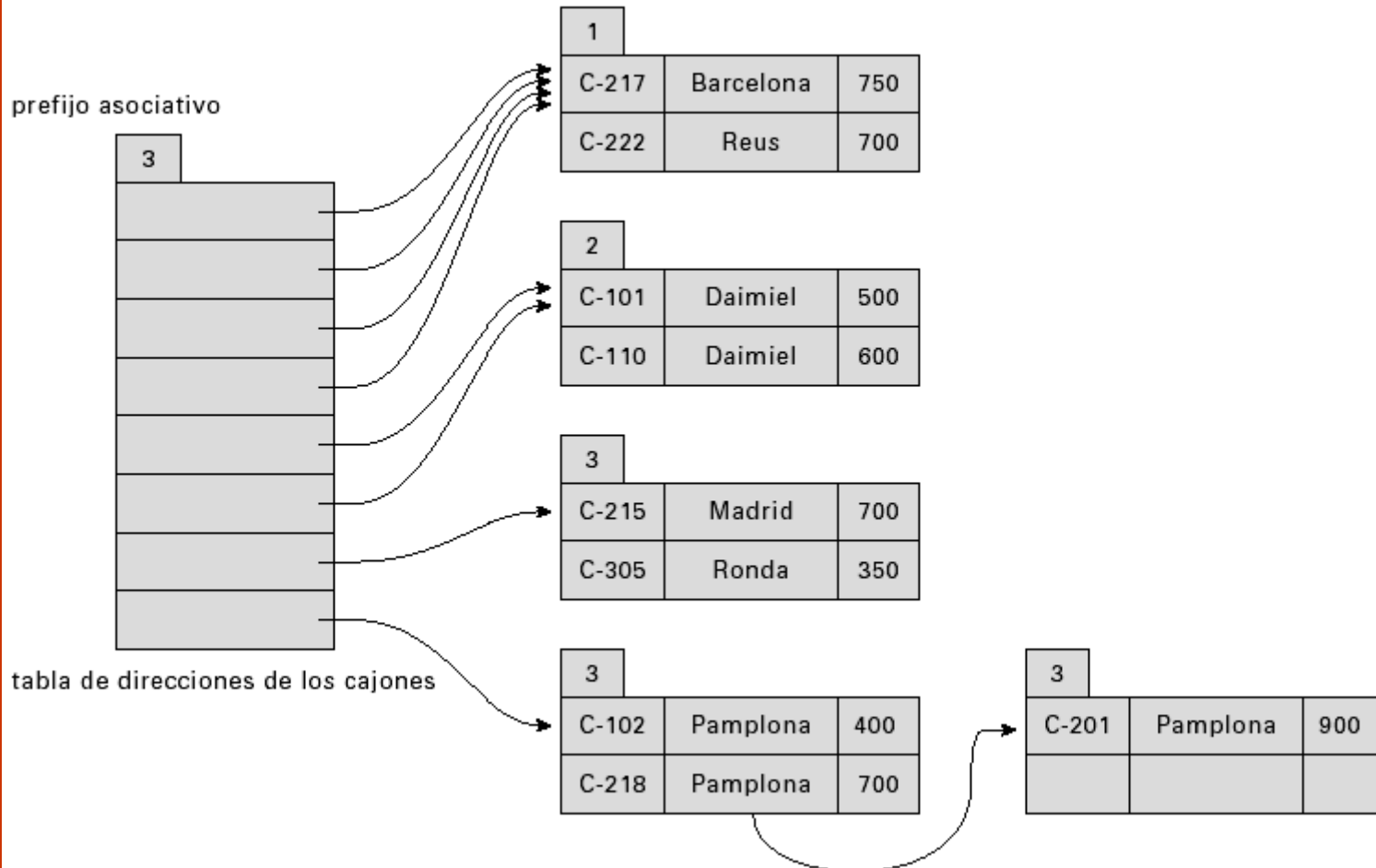


Estructura de asociación después de la inserción
de tres registros "Navacerrada".



Ejemplo (Cont.)

- Estructura de asociación después de la inserción de los registros “Reus” y “Ronda”





Asociación extensible vs otros esquemas

- Ventajas de la asociación extensible:
 - Las prestaciones de la asociación no disminuyen con el crecimiento del archivo
 - Costes de espacio mínimos
- Inconvenientes de la asociación extensible
 - Nivel extra de falta de dirección para encontrar el registro deseado
 - La tabla de direcciones de los cajones puede hacerse muy grande (más que la memoria)
 - ¡Necesita una estructura de árbol para localizar el registro deseado en la estructura!
 - El cambio del tamaño de la tabla de direcciones de cajones es una operación costosa



Comparación entre indexación ordenada y asociación

- Coste de una reorganización periódica
- Frecuencia relativa de inserciones y borrados
- ¿Es deseable optimizar el tiempo de acceso medio, a costa del tiempo de acceso del peor de los casos?
- Tipo esperado de consultas:
 - La asociación es generalmente mejor para recuperar registros que tienen un valor determinado de la clave.
 - Si las consultas de rangos son comunes, es preferible que los índices estén ordenados



Acceso multiclave

- Emplear múltiples índices para ciertos tipos de consultas.

- Ejemplo:

select *número-cuenta*

from *cuenta*

where *nombre-sucursal* = "Navacerrada" **and** *saldo* = 1.000

- Estrategias posibles para el procesamiento de consultas empleando índices sobre atributos simples:

1. Usar el índice sobre *nombre-sucursal* para encontrar cuentas con saldos de \$1.000; probar con *nombre-sucursal* = "Navacerrada".
2. Usar el índice sobre *saldo* para encontrar cuentas con saldos de \$1.000; probar con *nombre-sucursal* = "Navacerrada".
3. Emplear el índice de *nombre-sucursal* para encontrar punteros a todos los registros que pertenecen a la sucursal de Navacerrada. Análogamente, emplear el índice sobre *saldo*. Tomar la intersección de los dos conjuntos de punteros obtenidos.



Índices sobre múltiples atributos

Supóngase que tenemos un índice con la clave de búsqueda combinada (*nombre-sucursal, saldo*).

- Con la cláusula **where**
where *nombre-sucursal* = “Navacerrada” **and** *saldo* = 1.000
el índice sobre la clave de búsqueda combinada tomará sólo los registros que cumplan ambas condiciones.
Emplear índices independientes es menos eficiente — se pueden tomar muchos registros (o punteros) que sólo cumplen una de las condiciones.
- También se puede manejar eficientemente
where *nombre-sucursal* = “Navacerrada” **and** *saldo* < 1.000
- Pero no se puede gestionar eficientemente
where *nombre-sucursal* < “Navacerrada” **and** *saldo* = 1.000
Se pueden tomar muchos registros que cumplen la primera condición, pero no la segunda.

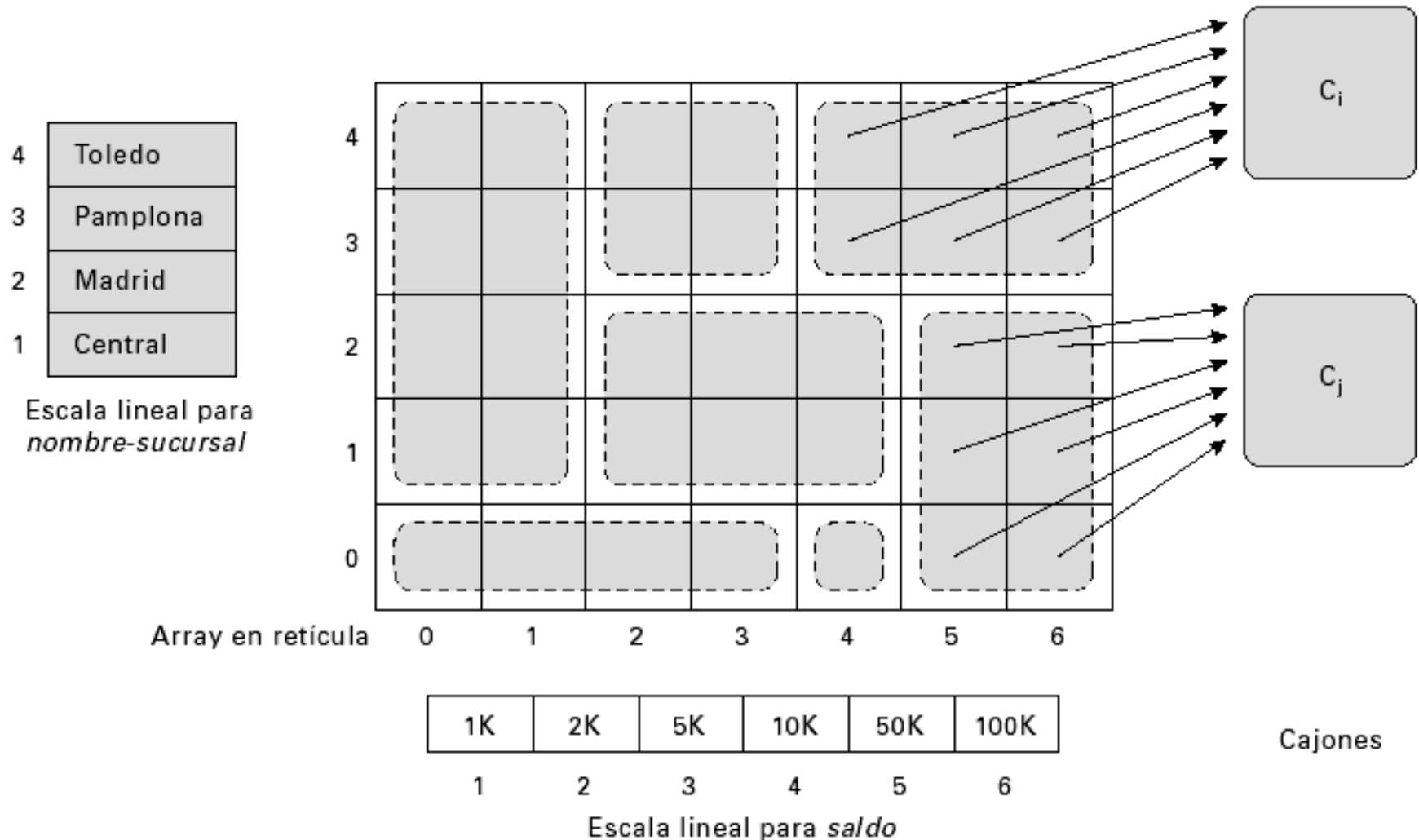


Archivos en retícula

- Estructura empleada para acelerar el procesamiento de consultas generales de claves de búsqueda múltiples, que implican una o más operadores de comparación.
- El **archivo en retícula** tiene un solo array en retícula y una escala lineal por cada atributo de la clave de búsqueda. El array en retícula tiene un número de dimensiones igual al de atributos de la clave de búsqueda.
- Múltiples celdas del array en retícula pueden apuntar al mismo cajón
- Para encontrar el cajón de un valor de clave de búsqueda, localizar la fila y la columna de su celda empleando las escalas lineales y seguir el puntero



Ejemplo de archivo en retícula para cuenta





Índices de mapas de bits

- Los índices de mapas de bits son un tipo especial de índice, diseñado para consultas eficientes sobre claves múltiples
- Los registros en una relación se asume que se numeran secuencialmente desde, por ejemplo, 0
 - Dado un número n debe ser fácil recuperar el registro n
 - Particularmente fácil si los registros son de tamaño fijo
- Aplicable sobre atributos que toman un número relativamente pequeño de valores distintos
 - Por ejemplo, sexo, país, provincia, ...
 - Por ejemplo, nivel de ingresos (ingresos descompuestos en un pequeño número de niveles tales como 0-9.999, 10.000-19.999, 20.000-50.000, 50.000- infinito)
- Un mapa de bits es simplemente un array de bits



Índices de mapas de bits (Cont.)

- En su forma más simple, un índice de mapa de bits sobre un atributo tiene una mapa de bits por cada valor del atributo
 - El mapa de bits tiene tantos bits como registros
 - En un mapa de bits para el valor v , el bit para un registro es 1 si el registro tiene el valor v para el atributo, de lo contrario es 0

número de registro	nombre	sexo	dirección	nivel-ingresos
0	Juan	m	Pamplona	L1
1	Diana	f	Barcelona	L2
2	María	f	Jaén	L1
3	Pedro	m	Barcelona	L4
4	Katzalin	f	Pamplona	L3

Mapas de bits para sexo		Mapas de bits para nivel-ingresos	
m	10010	L1	10100
f	01101	L2	10100
		L3	10100
		L4	10100
		L5	10100