

# CDB2

Los arboles surgieron para almacenar archivos de indice de forma más eficiente y ordenada.

*Permite localizar de forma rapida la información del archivo.*

*Ante cambio, solo debe reorganizarse localmente.*

Existen distintos tipos de arboles:

- Binarios: Maximo 2 hijos por cada nodo.
- AVL: Binario de busqueda auto-balanceado.
- Multicamino: tiene N caminos, con  $N > 2$ .
- Balanceados: la altura de los subárboles de cualquier nodo no difiere en más de una unidad.

## Binario

Es un arbol de grado 2, se utiliza para representar expresion arimeticas y de decisión. Tambien puede ser utilizados como arboles binarios de busqueda.

Las estructuras de datos utilizadas para almecenarlos deben ser registros de long fija, la informacion del archivo no está ordenada, cada nodo es un registro. Existe un desperdicio de espacio ya que hay muchos campos vacios.

## Balanceado

Un arbol donde la altura de la trayectoria más corta hacia una hoja no difiere de la altura hacia dicha hoja. Estos se desordenan facilmente asi que se buscó una solucion y esta fue los llamados AVL.

**AVL:** No esta balanceado completamente.

Sin embargo todas estas opciones sigue teniendo 2 problemas ppales la busqueda requiere demasiados desplazamientos y mantener un indice de orden es muy costoso. Los arboles **balanceados** ofrecen una solucion admisible al 2do problema.

## Binarios paginados

Los desplazamientos en memoria secundaria tienen un costo de tiempo relativamente alto pero una vez en posición los conjuntos de bits de interés están contiguos, entonces el desplazamiento es lento pero la transferencia es rápida.

Por eso se llegó a la idea de dividir a los árboles en páginas almacenadas cada una en un bloque de direcciones contiguas en disco, esto reduce el número de desplazamientos considerablemente.

Tiene un pequeño inconveniente que es que requiere transferir muchos datos que a veces no son usados, aún así sigue siendo más rápido.

Esta idea era buena pero nunca se encontró una forma de aplicarla eficientemente.

## Multicamino

Un árbol con  $N$  hijos,  $N > 2$ . Estos resuelven el problema que tenían los árboles binarios que era la cantidad de niveles muy elevada, sin embargo se debe mantener balanceado para asegurar eficiencia. Para ello estos árboles se contruyen desde la base hasta la raíz.

Para su construcción surgieron varias operaciones:

- Inserción: Busca hasta que llega al nivel hoja, se inserta en un nodo terminal. Hay 2 casos, uno cuando hay espacios que solo se realizan desplazamientos internos y otra donde el lugar está lleno y se produce la promoción
- Eliminación: Busca por cada índice hasta encontrar el nodo a eliminar, y de nuevo hay 2 casos, uno en donde sea una hoja y dicho nodo no quede vacío luego, y el otro en donde se produce underflow y debes o bien reordenar con el hermano o concatenarlo, perdiendo uno de los nodos.

## Árboles B

Son árboles multicamino, se construyen de forma ascendente. La inserción de un nuevo registro puede producir overflow, en dichos casos puedes tener que redistribuir o contener con el hermano de dicho nodo. Cuando el hermano tiene elementos suficientes se puede realizar la redistribución y cuando se debe concatenar donde es posible perder un nivel de árbol.

Para hacer más eficiente en términos de utilización de espacios se hizo una mejora y se la llamo B+, donde cada nodo está lleno por lo menos  $2/3$  partes, mismo manejo interno que un árbol B pero su manejo de nodos se realiza en

buffers lo que permite minimizar los acceso a disco. En la ram se almacena la raiz para ahorra en acceso a mem sec.

## Conclusion

Los arboles B y B+ permiten el acceso aleatorio por clave de forma eficiente ademas los B+ exigen que sus nodos se mantengan más llenos que en el otro. Los B+ son los unico que permiten un recorrido ordenado secuencial rapido. Este ultimo consta de 2 partes: *Indices* en la raiz y nodos interiores donde las claves son duplicadas en la *secuencia* que es donde se encuetran todas la claves colocadas en las hojas de arbol.

Los B+ ocupan más que el B ya que tiene q almacenar los separadores.

Tabla de comparacion

	B	B+
Objetos	nodos	nodos terminales
Tiempo de busq	equivalente	equivalente
Procesamiento sec ord	lente	rapido(ptos)
Ins/Elim	equivalente	equivalente

## Dispersion

Para los casos en que se necesita un mecanismo de acceso a registros con una sola lectura se usa **hash**. Esta tecnica genera una direccion base unica para un clave dada, convirtiendo la clave del registro en un numero aleatorio que luego se transforma en una direccion de mem donde debe ser almacenado el reg. Si esa direccion ya se está ocupada se realizan un tratado especial.

Tiene los sig beneficios

- No requiere almacenamiento adicional
- Facilita la insercion y eliminacion de reg rapida
- Encuentra reg con pocos accesos a disco

Costos

- No es posible el uso de reg de long variable

- No existe orden físico de datos
- No permite claves duplicadas

Hay 2 tipos de dispersión. Con espacio de dir estático o con uno dinámico

Parámetros que afectan a su eficiencia

1. Función de dispersión
2. Tamaño de los compartimientos
3. Densidad de empaquetamiento
4. Tratamiento de saturación

Caja negra: a partir de una clave se obtiene la dirección donde debe estar el registro, como no hay relación entre el índice del reg y la dir puede haber 2 reg que se quieran almacenar en la misma dirección esto se llama colisión.

Funciones de dispersión:

Centros cuadrados: la clave se multiplica por sí misma y se toman los dígitos centrales.

División: la clave se divide por un nro aprox igual a la cantidad de dir, luego se obtiene el residuo

Desplazamiento: se realiza el desp del número, se promociona y se suman las partes resultantes.

Plegado: se realiza un plegado del nro, se suman partes resultantes.

Transformado de base: la base del número se modifica tomando su módulo

1. **Cuadrados Centrales:** Supongamos que la clave es 1234. Al cuadrarla obtenemos 1522756. Tomando los dígitos centrales, por ejemplo, los cuatro dígitos centrales, obtenemos 2275 como valor de dispersión.
2. **División:** Si la clave es 1234 y usamos un número de división de 100 (que es aproximadamente igual a la cantidad de direcciones), el residuo es 34, que sería el valor de dispersión.
3. **Desplazamiento:** Si la clave es 123456 y realizamos un desplazamiento a la izquierda de 2 dígitos, obtenemos 3456. Luego, sumamos las partes resultantes (34 y 56) para obtener 90 como valor de dispersión.
4. **Plegado:** Si la clave es 123456789 y decidimos plegarla en grupos de 3 dígitos, obtenemos los grupos 123, 456 y 789. Al sumar estos grupos, obtenemos 1368 como valor de dispersión.

5. **Transformación de base:** Si la clave es 1234 y decidimos transformarla a base 8, obtenemos 2322. Luego, podemos tomar el módulo de este número con respecto a un número de direcciones para obtener el valor de dispersión.

Division es la mejor de todas para

Las colisiones son inherentes a este tipo de sistemas por ende se busca minimizar los máximos posibles. Para ello se eligen los métodos que las distribuyan de forma más aleatoria, se utilizan grandes espacios de dir para almacenar pocos reg. Un aspecto que mejora esto considerablemente es que cada dir puede almacenar más de un reg, a estos espacios de dir se les llama cubetas

El tamaño de dichas cubetas si es muy grande hay menor riesgo de saturación pero mayor fragmentación y la búsqueda es más lenta dentro del mismo, además de que si es grande la recuperación de un reg es lenta. El tamaño va depender del sistema teniendo en cuenta las posibilidades de trans en operaciones E/S a buffer del SO.

Lo que se tiene en cuenta a la hora de eficiencia es la **densidad de empaquetamiento DE** y con ella es posible hacer una estimación de la saturación.

En gral si hay N direcciones el nro esperado de I registros asignados es igual  $N * P(i)$ , donde  $P(i) = (K/N)^i * e^{K-N} / i!$

Otro gran cambio que disminuye la saturación es que los esp de dir contengan más de un reg.