

Taller File System

Tomás Felipe Melli

Noviembre 2024

Índice

| | | |
|----------|---|-----------|
| 1 | Introducción | 2 |
| 1.1 | Conocimientos previos | 2 |
| 1.1.1 | Definición de File System | 2 |
| 1.1.2 | Definición de Archivo | 2 |
| 1.1.3 | Definición de Inodo | 2 |
| 2 | Funciones | 5 |
| 2.1 | blockgroup_for_inode (unsigned int inode) | 5 |
| 2.2 | blockgroup_inode_index(unsigned int inode) | 5 |
| 2.3 | load_inode(unsigned int inode) | 5 |
| 2.4 | get_block_address(struct Ext2FSInode * inode, unsigned int block_number) | 6 |
| 2.5 | get_file_inode_from_dir_inode(struct Ext2FSInode * from, const char * filename) | 7 |
| 2.6 | inode_for_path(const char * path) | 9 |
| 3 | Ejercicio tipo parcial : recorrer en EXT2 y en FAT | 10 |

1 Introducción

En este taller se pondrán en práctica los conocimientos de File System. En particular, se trabajará con un FS basado en **inodos**. La idea es poder, implementar la lógica de las funciones principales para trabajar con estas estructuras :

- a) `blockgroup_for_inode` : para poder encontrar el **grupo** en que se encuentra el **inodo** pasado por parámetro.
- b) `blockgroup_inode_index` : para encontrar el índice de cierto inodo dentro de la **tabla de inodos**.
- c) `load_inode` : para cargar el **bloque** que contiene al inodo que buscamos en memoria y devolverlo. Le pasamos el nro de inodo.
- d) `get_block_address` : con esta función podremos obtener la **LBA** o **logical block address** del bloque que contiene al inodo que pasamos por parámetro (también con su nro de bloque).
- e) `get_file_inode_from_dir_inode` : es una función que a partir del **filename** queremos traernos su inodo. Tenemos la entrada de cierto directorio para buscar allí.

1.1 Conocimientos previos

1.1.1 Definición de File System

Un **File System** o sistema de archivos es un **ordenamiento lógico dentro de un disco**. Es decir, se busca establecer un orden para la información del usuario y del SO dentro de un dispositivo de almacenamiento. El responsable de este orden es un **módulo del kernel**. Este se encargará de :

- **Organización interna** : estructura interna del archivo.
- **Organización externa** : cómo se ordenan los archivos ? Hay limitaciones (mismo nombre,...)? Aparece el concepto de **directorío** donde se conforma la estructura jerárquica de árbol

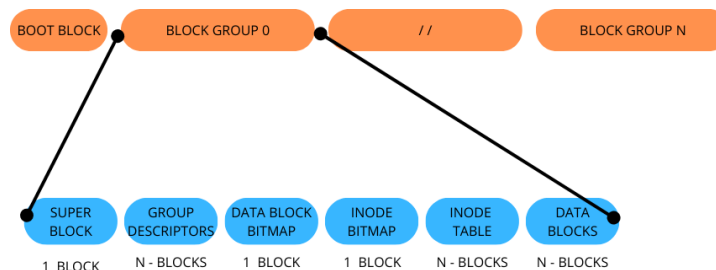
1.1.2 Definición de Archivo

Hablamos mucho de ordenar los archivos, de darle cierta estructura arbórea... pero qué es un archivo ? La respuesta es que es **una secuencia sin estructura de bytes(a nivel SO)**. El SO los identifica con un nombre y la extensión le tira medio centro sobre lo que **podría tener, esto no es determinante en cuanto a lo que tendrá dentro**.

** Existen tipos : d (directory), c (char device), b (block device) , - (regular), ... tirá `ls -l` en algún lado y chequealo. Y cómo hace este módulo para modelar un sistema lleno de secuencias de bytes sin estructura ?

1.1.3 Definición de Inodo

Un **inodo** es una estructura de datos que se utiliza en el FS que estudiaremos en este taller que se llama EXT2 (Extended File System 2) que se utiliza en Linux (y en otros) para **almacenar la información de un archivo**. Todo es un archivo en EXT2. La estructura del File System es :



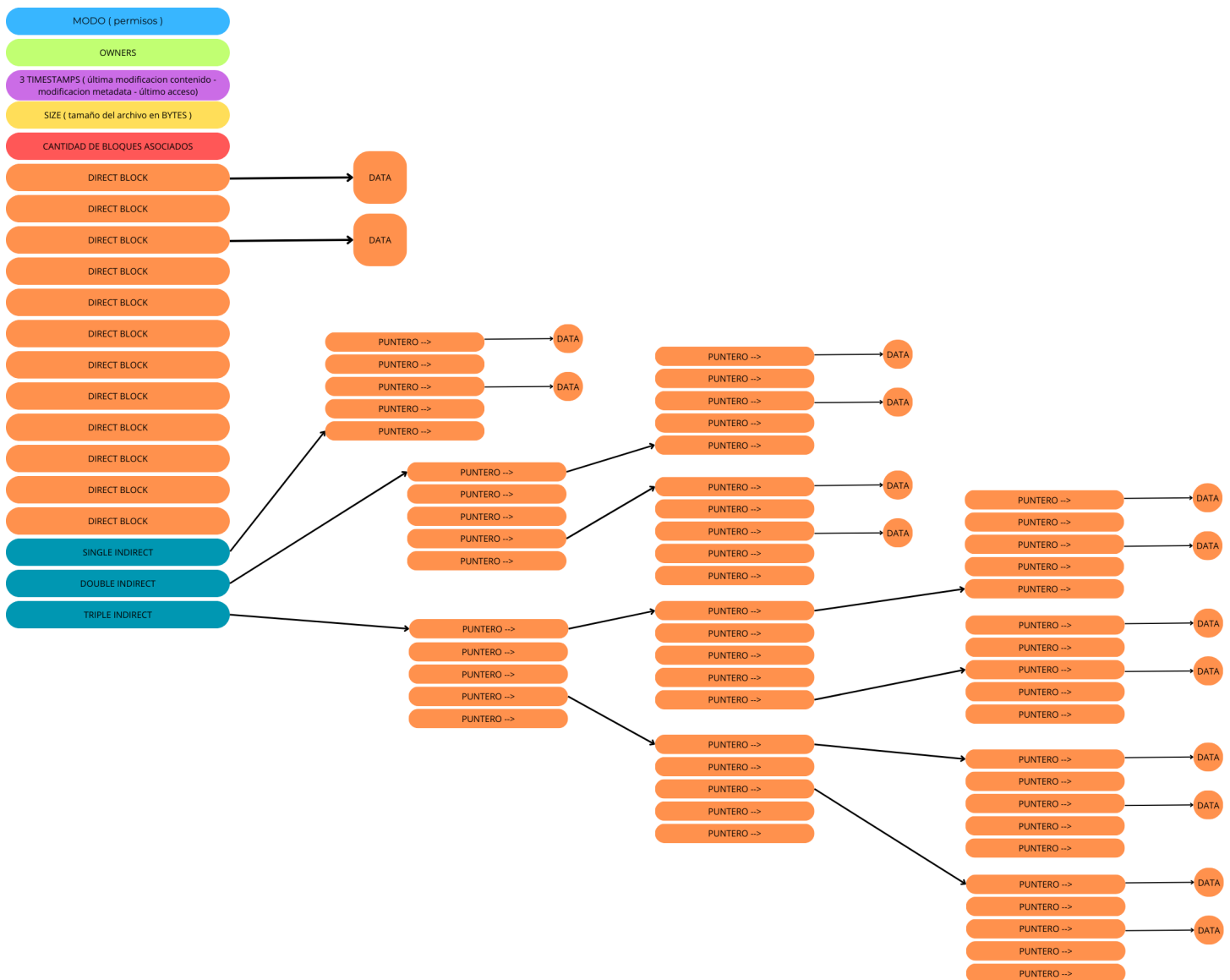
Tendremos :

- **Superblock** : contiene los metadatos críticos del FS:

```
1      struct __attribute__((__packed__)) Ext2FSSuperblock
2      {
3          unsigned int inodes_count;
4          unsigned int blocks_count;
5          unsigned int reserved_blocks_count;
6          unsigned int free_blocks_count;
7          unsigned int free_inodes_count;
8          unsigned int first_data_block;
9          unsigned int log_block_size;
10         unsigned int log_fragment_size;
11         unsigned int blocks_per_group;
12         unsigned int fragments_per_group;
13         unsigned int inodes_per_group;
14         unsigned int mount_time;
15         unsigned int write_time;
16         unsigned short mount_count;
17         unsigned short max_mount_count;
18         unsigned short magic_bytes;
19         unsigned short state;
20         unsigned short errors;
21         unsigned short minor_revision_level;
22         unsigned int lastcheck;
23         unsigned int checkinterval;
24         unsigned int creator_os;
25         unsigned int revision_level;
26         unsigned short default_reserved_userid;
27         unsigned short default_reserved_groupid;
28         // -- EXT2_DYNAMIC_REV Specific --
29         unsigned int first_inode;
30         unsigned short inode_size;
31         unsigned short block_group_number;
32         unsigned int feature_compatibility;
33         unsigned int feature_incompatibility;
34         unsigned int readonly_feature_compatibility;
35         char uuid[16];
36         char volume_name[16];
37         char last_mounted[64];
38         unsigned int algo_bitmap;
39         // Performance hints
40         unsigned char prealloc_blocks;
41         unsigned char prealloc_dir_blocks;
42         char alignment[2];
43         // Journaling support
44         char journal_uuid[16];
45         unsigned int journal_inode;
46         unsigned int journal_device;
47         unsigned int last_orphan;
48         // Directory indexing support
49         unsigned int hash_seed[4];
50         unsigned char default_hash_version;
51         char padding[3];
52         // Other options
53         unsigned int default_mount_options;
54         unsigned int first_meta_bg;
55         char unused[760];
56     };
```

- **Descriptor de grupo** : ayuda al SO a manejar permisos y acceso a recursos.
- **Tabla de inodos** : contiene los inodos del grupo
- **Bloques de datos**
- **Bitmaps** de inodos y de bloques de datos

Y el Inodo luce algo como esto :



Y su representación en C es como sigue :

```

1 struct __attribute__((__packed__)) Ext2FSInode {
2     unsigned short mode;
3     unsigned short uid;
4     unsigned int size;
5     unsigned int atime;
6     unsigned int ctime;
7     unsigned int mtime;
8     unsigned int dtime;
9     unsigned short gid;
10    unsigned short links_count;
11    unsigned int blocks;
12    unsigned int flags;
13    unsigned int os_dependant_1;
14    unsigned int block[15];
15    unsigned int generation;
16    unsigned int file_acl;
17    unsigned int directory_acl;
18    unsigned int faddr;
19    unsigned int os_dependant_2[3];
20 };

```

La función principal de esta estructura es representar un **archivo**. Pero qué pasa con los directorios ? También se representan con estos, la estructura es igual. Lo que cambia es que cada bloque apuntado, contendrá **DirEntries** de la siguiente forma :

```

1 struct __attribute__((__packed__)) Ext2FSDirEntry {
2     unsigned int inode;
3     unsigned short record_length;
4     unsigned char name_length;
5     unsigned char file_type;
6     char name[];
7 };

```

2 Funciones

2.1 blockgroup_for_inode (unsigned int inode)

Hablamos un poco del **superbloque**, este tiene mucha info, nos interesará saber **cuántos inodos hay por grupo** ya que, como sabemos el número del inodo podemos saber en qué grupo estará de la siguiente forma :

$$\frac{\text{inode} - 1}{_superblock \rightarrow \text{inodes_per_group}}$$

2.2 blockgroup_inode_index(unsigned int inode)

En esta función queremos saber en qué posición dentro de la **inode table** dentro del grupo estará el inodo.

$$(\text{inode} - 1) \bmod _superblock \rightarrow \text{inodes_per_group}$$

2.3 load_inode(unsigned int inode)

La idea es cargarlo en memoria al inodo. Los pasos son :

- Saber en qué número de grupo está con `blockgroup_for_inode`.
- Conocer el **descriptor del grupo** con la función **block_group(unsigned int grupo)** que básicamente indexa en una tabla que tiene los descriptores de los grupos. Este objeto tiene la **dirección de la inode table**.
- Una vez capturada la dire de la tabla, tenemos que saber dónde indexar según el inodo que tenemos. Llamamos a `blockgroup_inode_index`.
- **Aclaración :** un inodo normalmente ocupa mucho menos espacio que lo que ocupa un bloque entero en memoria. Es por ello que tenemos que considerar esto a la hora de hacer la carga del inodo en memoria. Dicho esto, capturamos el **tamaño de bloque**, este dato lo tiene el **superblock** bajo el atributo `_superblock \rightarrow log_block_size`, pero ojo, porque este valor es \log_2 del tamaño. Entonces tenemos que llevarlo a KB (simplemente se multiplica por 2^{10} bytes. Entonces, una vez calculado el tamaño en KB, queremos saber **cuántos inodos habrá por bloque**

$$\text{cant_inodos_por_bloque} = \frac{\text{tam_bloque}}{_superblock \rightarrow \text{inode_size}}$$

- Pensemos en algo antes, la tabla de inodos está en memoria también, y no necesariamente ocupa 1 bloque, entonces si nuestro inodo está dentro de esa tabla, en particular podría no estar en el bloque 1 que la contiene. Para asegurarse qué bloque que contiene la tabla de inodos contiene al inodo que buscamos hacemos la siguiente cuenta :

$$\text{bloque_inodo} = \frac{\text{índice}}{\text{cant_inodos_por_bloque}}$$

- Una vez que sabemos el bloque de la tabla de inodos que contiene el inodo que buscamos, necesitamos el **offset** adecuado :

$$\text{offset_inodo} = (\text{índice} \% \text{cant_inodos_por_bloque}) \times _superblock \rightarrow \text{inode_size}$$

- Nos traemos el bloque con **read_block(unsigned int block_address, unsigned char * buffer)** y en el lugar adecuado leemos y copiamos la info que nos interesa.

Solución:

```

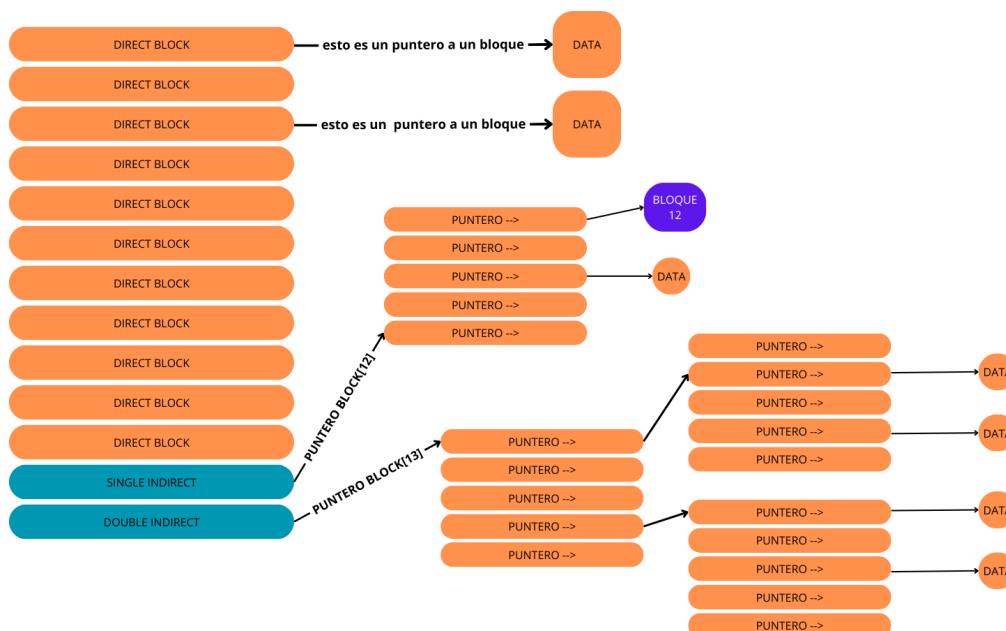
1 struct Ext2FSInode * Ext2FS::load_inode(unsigned int inode_number)
2 {
3     unsigned int nro_grupo = blockgroup_for_inode(inode_number);
4     Ext2FSBlockGroupDescriptor* descriptor = block_group(nro_grupo);
5     unsigned int dire_tabla = descriptor->inode_table;
6     unsigned int indice = blockgroup_inode_index(inode_number);
7     unsigned int tam_bloque = 1024 << _superblock->log_block_size;
8     unsigned int cant_inodos_por_bloque = tam_bloque / _superblock->inode_size;
9     unsigned int bloque_inodo = indice / cant_inodos_por_bloque ;
10    unsigned int offset_inodo = (indice % cant_inodos_por_bloque) * _superblock->inode_size;
11
12    unsigned char* res = (unsigned char*) malloc(sizeof(Ext2FSInode));
13    unsigned char* lectura = (unsigned char*) malloc(sizeof(tam_bloque));
14
15    read_block(dire_tabla + bloque_inodo, lectura);
16
17    for (int i = 0; i < sizeof(Ext2FSInode); i++)
18    {
19        res[i] = lectura[offset_inodo + i];
20    }
21
22    delete[] lectura;
23    return (Ext2FSInode*)res;
24 }

```

2.4 get_block_address(struct Ext2FSInode * inode, unsigned int block_number)

Recordemos la estructura del inodo. Teníamos un array llamado **block[15]**. Estos son direcciones de bloques en memoria. Se llaman LBA como ya hablamos. **no son punteros** son enteros sin signo que se interpretan como direcciones lógicas en el disco. Nos interesará lo siguiente : **obtenerla**. Para ello, como se intuye, vamos a tener que recorrer el inodo en busca de ella. Hay ciertas situaciones a las que tendremos que estar atentos :

- Si el número de bloque se encuentra en los **bloques directos** entonces simplemente devolvemos inode → block[block_number].
- Si el número de bloque excede el ítem anterior, tenemos que buscar en el **bloque de indirección simple**. En este bloque vamos a encontrar muchas LBAs, cuántas ? cuál es el número de bloque más grande que puede ser **contenida** por este bloque de indirección simple ? Estas son las preguntas adecuadas para poder encontrar la LBA del bloque que nos piden.
- Si el número de bloque excede el indirecto simple, hay que moverse al **indirecto doble**, esta estructura se vuelve más pesadita de leer y todo. Hay que buscar en una LBA que tiene LBAs de bloques que tienen otras LBAs que esas últimas pueden ser las que apuntan al nro de bloque que nos interesa.



Solución:

```
1 unsigned int Ext2FS::get_block_address(struct Ext2FSInode * inode, unsigned int block_number)
2 {
3     if (block_number < 12)
4     {
5         return inode->block[block_number];
6     }
7     else
8     {
9         unsigned int res;
10        unsigned int limite_indirecto_simple = _superblock->log_block_size / sizeof(unsigned int);
11
12        if (block_number < limite_indirecto_simple)
13        {
14            unsigned char* lectura = (unsigned char*) malloc(1024 << _superblock->log_block_size);
15            unsigned int dire_bloque_indirectos_simple = inode->block[12];
16            read_block(dire_bloque_indirectos_simple, lectura);
17            unsigned int dire_bloque = block_number - 12;
18            res = ((unsigned int*)lectura)[dire_bloque];
19            delete[] lectura;
20        }
21        else
22        {
23            unsigned int dire_bloque_indirectos_doble = inode->block[13];
24            unsigned char* lectura = (unsigned char*) malloc(1024 << _superblock->log_block_size);
25            read_block(dire_bloque_indirectos_doble, lectura);
26            unsigned int primero = block_number - limite_indirecto_simple;
27            unsigned int dire_bloque_apuntado_primer_puntero = ((unsigned int*) lectura)[primero];
28            read_block(dire_bloque_apuntado_primer_puntero, lectura);
29            unsigned int offset = primero % _superblock->inode_size;
30            res = ((unsigned int*)lectura)[offset];
31            delete[] lectura;
32        }
33        return res;
34    }
35 }
```

2.5 get_file_inode_from_dir_inode(struct Ext2FSInode * from, const char * filename)

Ya hablamos sobre los **directorios** y cómo se representan. Lo que nos piden es básicamente iterar sobre toda la info que tiene un inodo, llamado **from**, y dentro de él vamos a encontrar **DirEntries**. Como ya hicimos toda la movida de sacar la dire en el ejercicio anterior, no nos preocuparemos por la indirección. Sólo, traer un bloque lleno de DirEntries, chequear si el nombre coincide con el **filename** que nos pasan, sino continúe.

****Aclaraciones :** cuando hablamos de FS y cómo el módulo maneja las limitaciones, mencionamos los nombres, lo que no dijimos fue : hay límite de tamaño del nombre ? Eso influye en el tamaño de un DirEntry a otro a la hora de recorrerlos. Puede suceder que comience en un bloque y termine en el siguiente ? Sí y, por ese motivo hay que tener 2 bloques de lectura, porque no sabemos a priori qué datos caen el uno y qué datos caen en otro, por tanto si hacemos DirEntry->record_length tal vez ese dato está en el siguiente bloque y se rompe todo. La estrategia es :

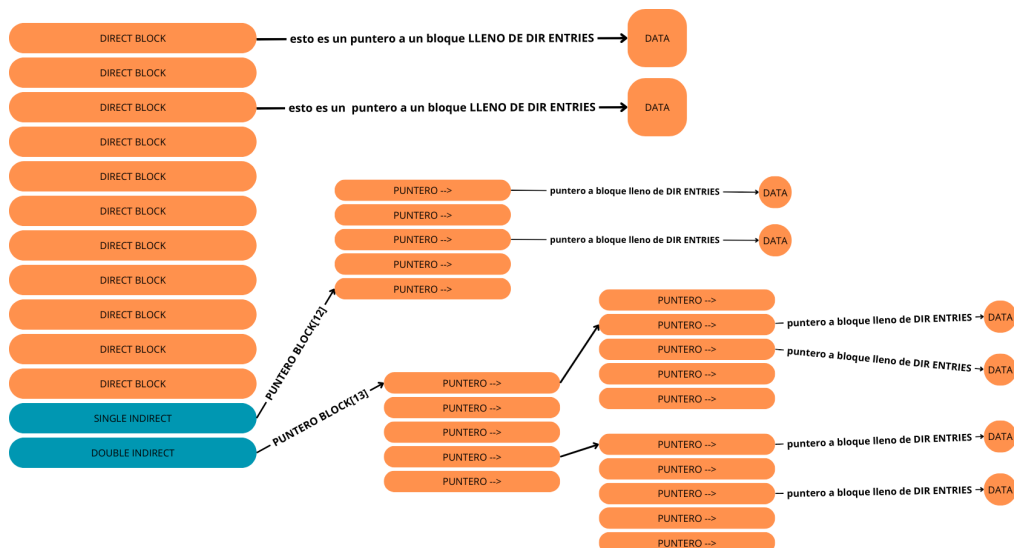
- Primero necesitamos saber de qué tamaño son los bloques.
- Como dijimos que tenemos que tener 2 bloques de datos por si una dir entry arranca en un bloque y termina en el consecutivo, pedimos memoria, dos blocksize y mantenemos referencia al punto en donde comienza el segundo bloque.
- Arrancamos leyendo por el bloque 0 con el offset en 0.
- Necesitamos la direcciones (LBAs) de los bloques a leer. Usamos la función anterior para ello.
- Hacemos las lecturas
- Ingresamos al bucle : tenemos que ver si ya terminamos de leer todas las entradas del un bloque. Para ello, miramos si el offset se pasó del tamaño de bloque.
- Si se pasó, hay que traer el siguiente bloque, acomodar el offset y la dinámica es la misma, capturamos las LBAs y leemos.
- Si no se pasó, tenemos que capturar cierta dir entry y hacer la lógica de la comparación. Chequeamos si el filename coincide, en caso de coincidir, cargamos el inodo. Caso contrario, continuamos buscando.

Solución:

```

1 struct Ext2FSInode * Ext2FS::get_file_inode_from_dir_inode(struct Ext2FSInode * from, const char *
   filename)
2 {
3     if(from == NULL)
4         from = load_inode(EXT2_RDIR_INODE_NUMBER);
5     //std::cerr << *from << std::endl;
6     assert(INODE_ISDIR(from));
7
8     unsigned int block_size = 1024 << _superblock->log_block_size;
9
10    unsigned char * block_buf = (unsigned char *) malloc(2 * block_size);
11    unsigned char * block_sig = (unsigned char *) block_buf + block_size;
12
13    unsigned int block_number = 0;
14    unsigned int offset = 0;
15
16    unsigned int block_address = get_block_address(from, block_number);
17    unsigned int block_address2 = get_block_address(from, block_number+1);
18
19    read_block(block_address, block_buf);
20    read_block(block_address2, block_sig);
21
22    while (true) {
23        if(offset > block_size) {
24            block_number++;
25
26            offset -= block_size;
27
28            block_address = get_block_address(from, block_number);
29            block_address2 = get_block_address(from, block_number+1);
30
31            read_block(block_address, block_buf);
32            read_block(block_address2, block_sig);
33        }
34
35        Ext2FSDirEntry *entrada_actual = (Ext2FSDirEntry*) (block_buf + offset);
36        int size_name = entrada_actual->name_length;
37
38        char* filename_inode = entrada_actual->name;
39
40        if(size_name == strlen(filename) && !strcmp(filename, filename_inode, size))
41        {
42            unsigned int inodo_actual = entrada_actual->inode;
43            free(block_buf);
44            return load_inode(inodo_actual);
45        }
46
47        offset += entrada_actual->record_length;
48    }
49    return NULL;
50 }

```



2.6 inode_for_path(const char * path)

La idea es, yo te paso `/home/user/documentos/archivo.txt` con la función `strtok(mypath, "/")`; separamos todos los directorios y **retornamos el inodo**. Mientras tengamos directorios para recorrer, ejecutamos la función anterior `get_file_inode_from_dir_inode(inodo_anterior, filename` que nos interesa encontrar en este), que nos permitirá actualizar el inodo para buscar el filename del siguiente y así hasta dar con el inodo del archivo.txt.

```
1 // Estructura Ext2FSInode * Ext2FS::inode_for_path
2 struct Ext2FSInode * Ext2FS::inode_for_path(const char * path)
3 {
4     char * mypath = new char[strlen(path)+1];
5     mypath[strlen(path)] = '\0';
6     strncpy(mypath, path, strlen(path));
7
8     char * pathtok = strtok(mypath, PATH_DELIM);
9     struct Ext2FSInode * inode = NULL;
10
11     while(pathtok != NULL)
12     {
13         struct Ext2FSInode * prev_inode = inode;
14         // std::cerr << "pathtok: " << pathtok << std::endl;
15         inode = get_file_inode_from_dir_inode(prev_inode, pathtok);
16         pathtok = strtok(NULL, PATH_DELIM);
17
18         delete prev_inode;
19     }
20
21     delete[] mypath;
22     return inode;
23 }
```

3 Ejercicio tipo parcial : recorrer en EXT2 y en FAT

```
1 int my_copy_ext2(char* src_path, char* dst_path) {
2     // Obtenemos los inodos para los archivos src y dst identificados por su path.
3     struct Ext2FSInode* src_inode = inode_for_path(src_path);
4     struct Ext2FSInode* dst_inode = inode_for_path(dst_path);
5
6     // Creamos un buffer donde vamos a cargar los bloques que vamos copiamos.
7     unsigned char* buffer = (unsigned char*) malloc(BLOCK_SIZE);
8
9     // Calculamos el total de bloques a copiar a partir del atributo size del
10    // archivo fuente (src). Notar que redondeamos hacia arriba con ceil().
11    int block_count = ceil(src_inode->size / BLOCK_SIZE);
12
13    // Copiamos bloque a bloque de src a dst.
14    for (int block_number = 0; block_number < block_count; block_number++) {
15        // Obtenemos la direcci n (LBA) del bloque a copiar de src.
16        unsigned int src_block_address = get_block_address(src_inode, block_number);
17
18        // Leemos el bloque de src y lo cargamos en el buffer.
19        read_block(src_block_address, buffer);
20
21        // Pedimos un bloque libre y lo asignamos al inodo de dst con el
22        // n mero de bloque correspondiente.
23        unsigned int dst_block_address = get_free_block_address(dst_inode);
24        add_block_address_to_inode(dst_inode, dst_block_address, block_number);
25
26        // Escribimos el buffer en el bloque asignado a dst.
27        write_block(dst_block_address, buffer);
28    }
29
30    free(buffer);
31    return 0;
32 }
33
34 int my_copy_fat(char* src_path, char* dst_path) {
35     // Obtenemos las direcciones (LBA) del primer bloque de cada archivo.
36     unsigned int src_block_address = get_init_block_for_path(src_path);
37     unsigned int dst_block_address = get_init_block_for_path(dst_path);
38
39     // Creamos un buffer donde vamos a cargar los bloques que vamos copiamos.
40     unsigned char* buffer = (unsigned char*) malloc(BLOCK_SIZE);
41
42     // A priori no sabemos cu ntos bloques vamos a copiar.
43     while (src_block_address != EOF) {
44         // Leemos el bloque de src y lo escribimos en dst.
45         read_block(src_block_address, buffer);
46         write_block(dst_block_address, buffer);
47
48         // Navegamos la tabla de FAT del archivo src. En src_block_address ya
49         // tenemos el pr ximo bloque de datos.
50         src_block_address = FAT[src_block_address];
51         if (src_block_address == EOF) {
52             // Si llegamos al EOF, marcamos en la tabla de dst el EOF y terminamos.
53             FAT[dst_block_address] = EOF;
54         } else {
55             // Caso contrario, pedimos un nuevo bloque libre para dst y lo
56             // configuramos como el siguiente bloque de datos de dst para la
57             // pr xima iteraci n del while.
58             next_dst_block_address = get_free_block_address();
59             FAT[dst_block_address] = next_dst_block_address;
60             dst_block_address = next_dst_block_address;
61         }
62     }
63
64    free(buffer);
65    return 0;
66 }
67
68 // Versi n alternativa con while(true).
69 int my_copy_fat_while_true(char* src, char* dst) {
70     uint src_block_addr = get_init_block_for_path(src);
71     uint dst_block_addr = get_init_block_for_path(dst);
72
73     unsigned char* buffer = (unsigned char*) malloc(BLOCK_SIZE);
```

```

74
75     while (true) {
76         read_block(src_block_addr, buffer);
77         write_block(dst_block_addr, buffer);
78
79         if (FAT[src_block_addr] == EOF) {
80             FAT[dst_block_addr] = EOF;
81             break;
82         } else {
83             src_block_addr = FAT[src_block_addr];
84             next_dst_block_addr = get_free_block_address();
85             FAT[dst_block_addr] = next_dst_block_addr;
86             dst_block_addr = next_dst_block_addr;
87         }
88     }
89
90     free(buffer);
91     return 0;
92 }

```