

# MyJQL

FDU 2024 – Introduction to Database Systems – Project II  
2024.5

# Introduction

- **MyJQL**是一个类似**redis**的键值数据库，键值类型均为变长字符串；
- **key**上建立**B**-树索引，以提高访问效率；
- 具有**get**，**set**，**del**三种操作。

key	value
aa	adfasdfasdfasdfads
aab	sdfadfasdfasdfasdfasdf
abc	adfasdfsadfasdfafasdfasdf
bbbb	dfsadfasdfadf
ccccc	dfasdf
def	sdfasdfasdfasdfasdf

# An Intuitive Practice

- 数据保存于文本文件，启动时全部读入，退出时全部写出。

key	value
aa	adfasdfasdfasdfads
aab	sdfadfasdfasdfasdfasdf
abc	adfasdfsadfasdfafasdfasdf
bbbb	dfsadfasdfadf
ccccc	dfasdf
def	sdfasdfasdfasdfasdf



xxx.txt

-----  
aa [sep] adfasdfasdfasdfads [sep]  
aab [sep] sdfadfasdfasdfasdf [sep]  
abc [sep] adfasdfsadfasdfafasdfasdf [sep]  
bbbb [sep] dfsadfasdfadf [sep]  
ccccc [sep] dfasdf [sep]  
def [sep] sdfasdfasdfasdfasdf [sep]

- 数据量比较大时，这样做显然不佳。

# A Different Practice

---

- 数据量比较大时，这样做.....



# Logical View (1)

- 分两张表，键值分别为指向另一张表中的字符串的指针。

key	value
aa	adfasdfasdfasdfads
aab	sdfadfasdfasdfasdfasdf
abc	adfasdfsadfasdfafasdfasdf
bbbb	dfsadfasdfadf
ccccc	dfasdf
def	sdfasdfasdfasdfasdf

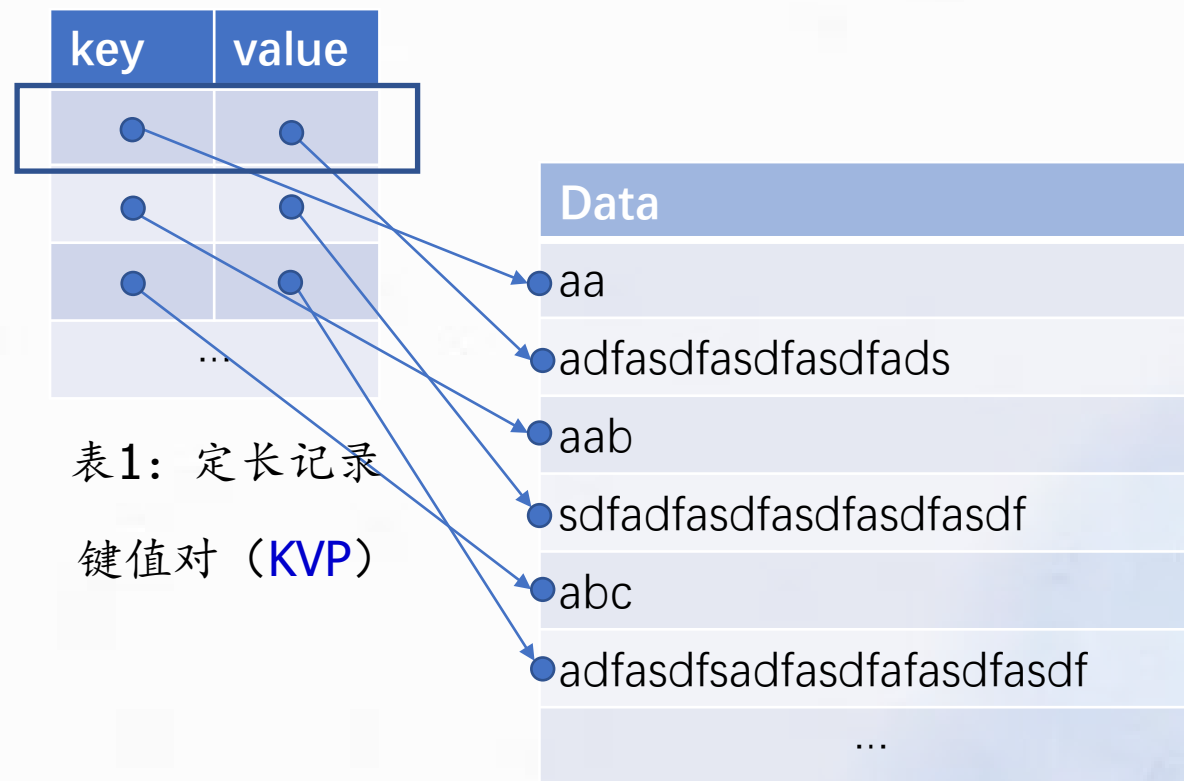


表2': 变长记录

# Logical View (2)

表2'可以进一步变化

- 长字符串切分成字符串块（**Chunk**），以链表形式保存。
- 字符串块最大不能超过一页的大小。

Data
aa
adfasdfasdfasdfs
aab
sdfadfasdfasdfasdfs
abc
...

表2': 变长记录

Data	Next
aa	
adfasdfasd	●
fasdfads	
aab	
sdfadfasdf	●
asdfasdfs	●
df	
abc	
...	

表2: 变长记录

# Logical View (3)

• 即：

key	value
aa	adfasdfasdfasdfads
aab	sdfadfasdfasdfasdf
abc	adfasdfsadfasdfafasdfasdf
bbbb	dfsadfasdfadf
ccccc	dfasdf
def	sdfasdfasdfasdfasdf

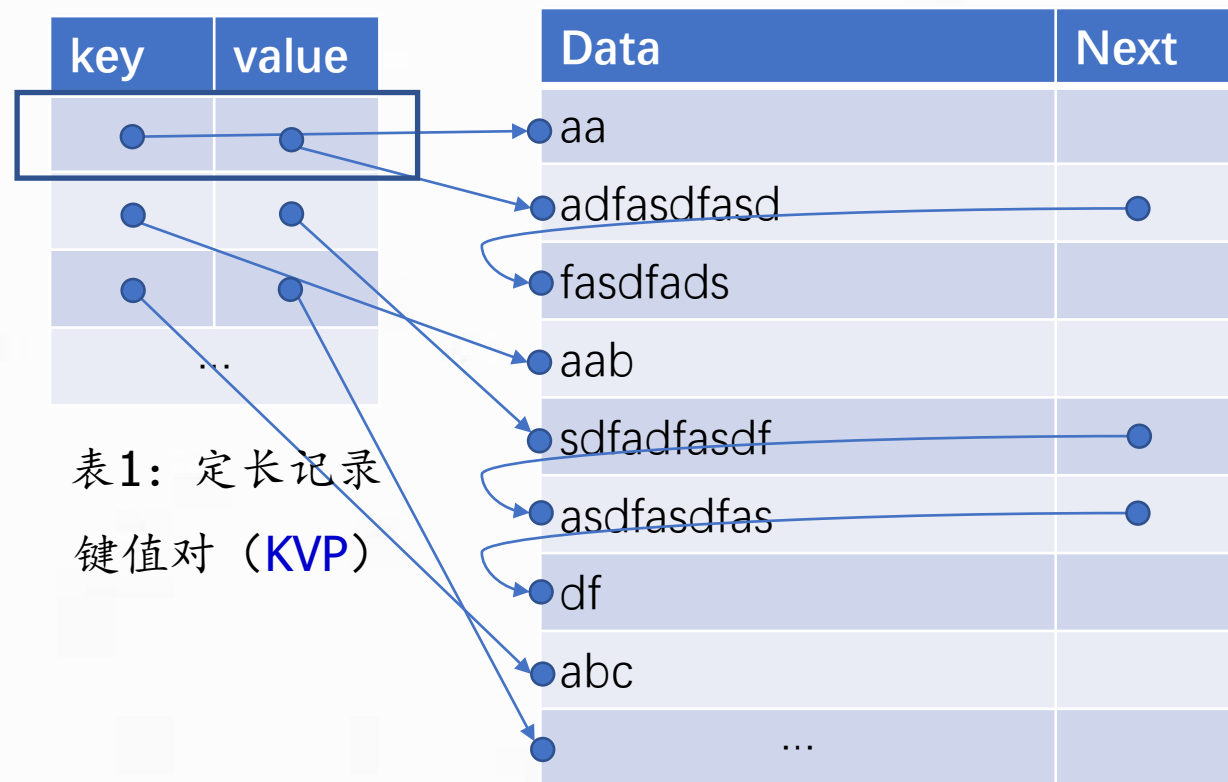
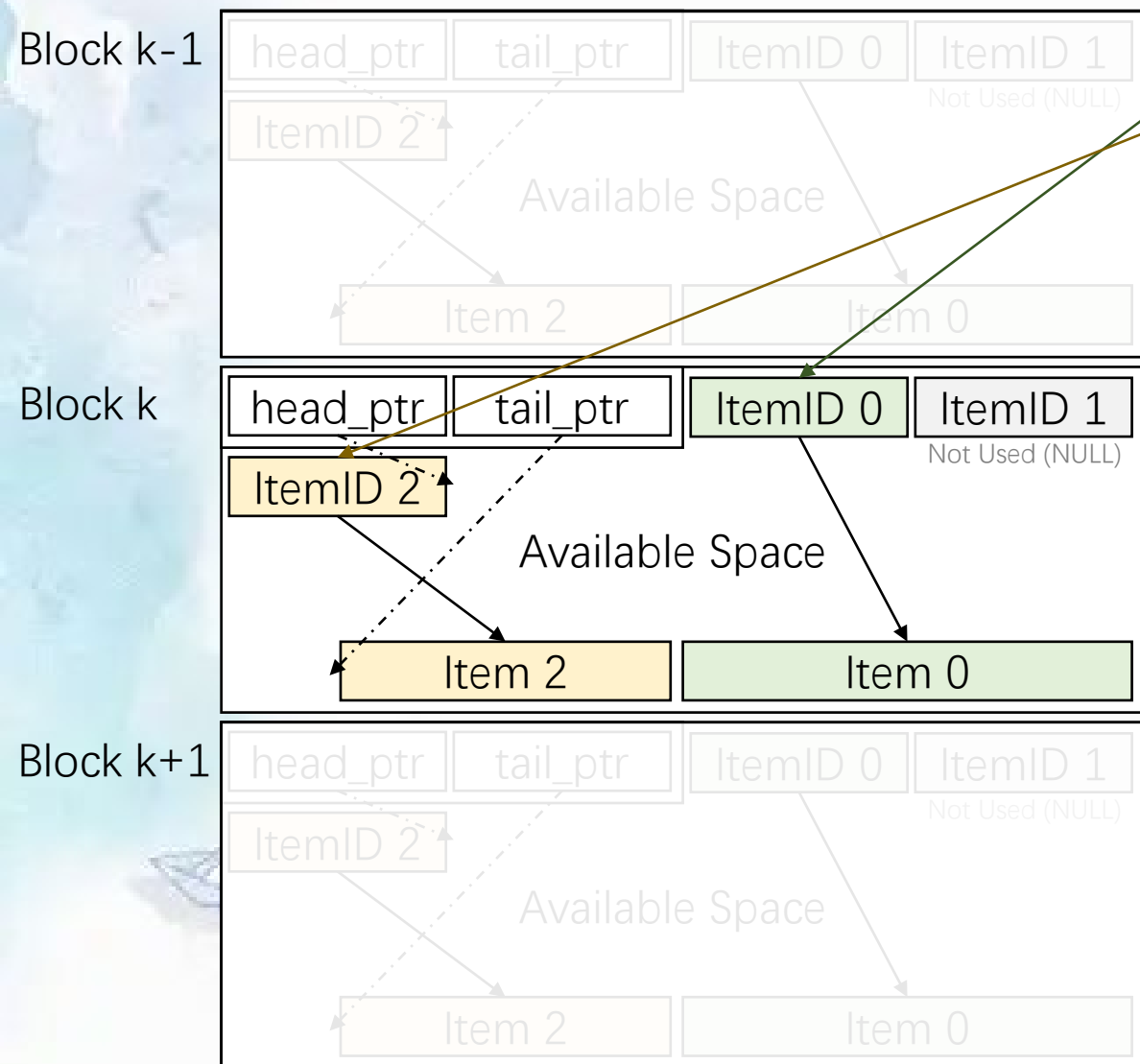


表2: 变长记录

# Physical View



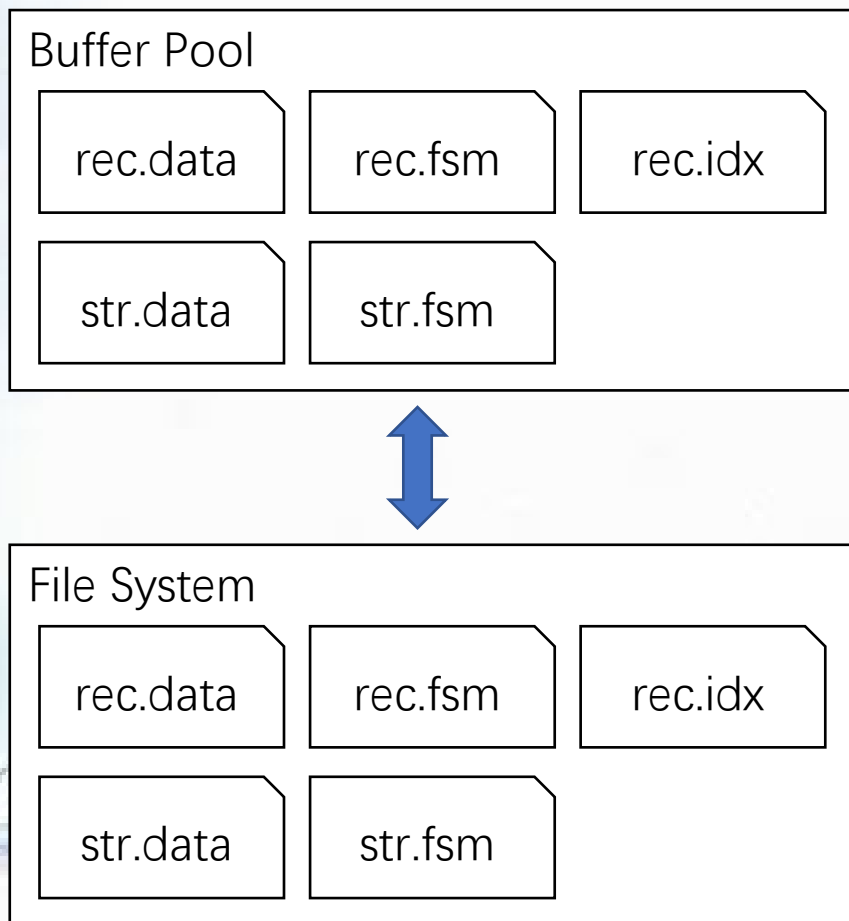
○ RID: (k, 0)

○ RID: (k, 2)

- 一个数据文件 (\*.data, 即存储数据的文件, 而非索引等其他数据结构的文件) 由若干个块 (Block) 构成。
- 每个块由头 (Header, 包含 head\_ptr 和 tail\_ptr 等字段), 若干 ItemID 和 Item 构成。
- ItemID (定长) 从头部开始分配, 作为标识, 指向尾部的 Item (变长), Item 保存具体数据 (KVP 或字符串块)。
- head\_ptr 至 tail\_ptr 指向的空间为可用空间。



# File & IO



- 所有文件IO操作通过缓冲池完成。
- 每个文件有独立的缓冲池。
- \*.data: 保存具体数据。
- \*.fsm: 维护空闲空间信息。
- \*.idx: 保存B-树索引。

# Overview

- 文件列表：
  - file\_io.h
  - buffer\_pool.h
  - block.h
  - hash\_map.h
  - table.h
  - str.h
  - b\_tree.h
  - myjql.h

# file\_io.h

- `file_io.c`已给，若非必要，无须修改。
- 重要宏：
  - `PAGE_SIZE`
    - 页大小，规定为128。
    - 为了方便调试，建议在引入B-树前使用64，引入B-树后使用128。
- 重要函数：
  - `FileIOResult read_page(Page *page, const FileInfo *file, off_t addr);`
    - 读入位于`addr`的`page`。
  - `FileIOResult write_page(const Page *page, FileInfo *file, off_t addr);`
    - 将`page`写出至`addr`，若`addr==`文件大小，则追加新的页。

[参考链接：MySQL的缓冲池管理](#)

# buffer\_pool.h

- **BufferPool**数据结构请自行设计，缓冲池内需要缓冲一定量的页，可以固定缓冲池大小，也可以动态申请。
- **void init\_buffer\_pool(const char \*filename, BufferPool \*pool);**
  - 打开**filename**，并关联**pool**为其缓冲池。
- **void close\_buffer\_pool(BufferPool \*pool);**
  - 关闭缓冲池，将缓冲的页写回文件。
  - 注意：无需考虑程序异常终止，测试保证程序一定会正常退出。
- **Page \*get\_page(BufferPool \*pool, off\_t addr);**
  - 获取地址为**addr**的页，并锁定（保证该页不会被意外换出）。
- **void release(BufferPool \*pool, off\_t addr);**
  - 释放地址为**addr**的页，该页之后可以被换出。
- 提示：文件读写调用**file\_io.h**中的读写函数。

# block.h

ItemID (32-bit unsigned int):

00000000 00000000 00000000 00000000

Unused bit

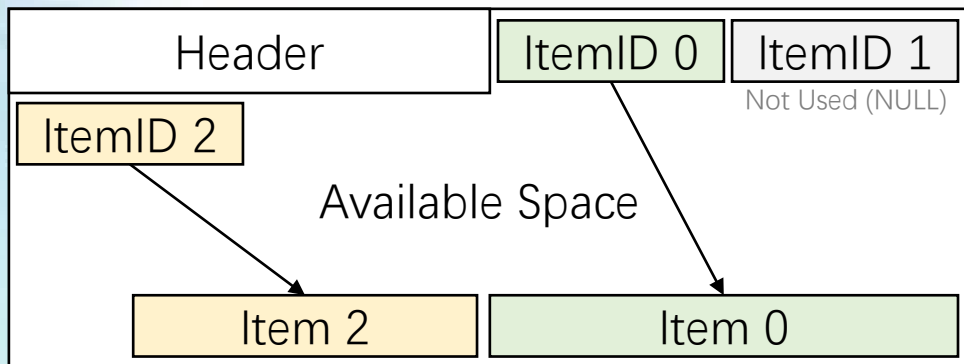
Offset: Item的起始地址

Size: Item的大小

Availability

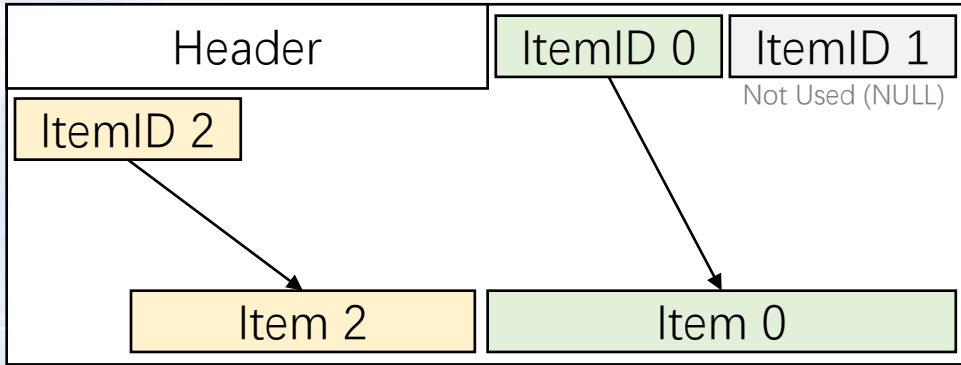
0: 已使用, 指向具体Item

1: 空闲可用, 空指针

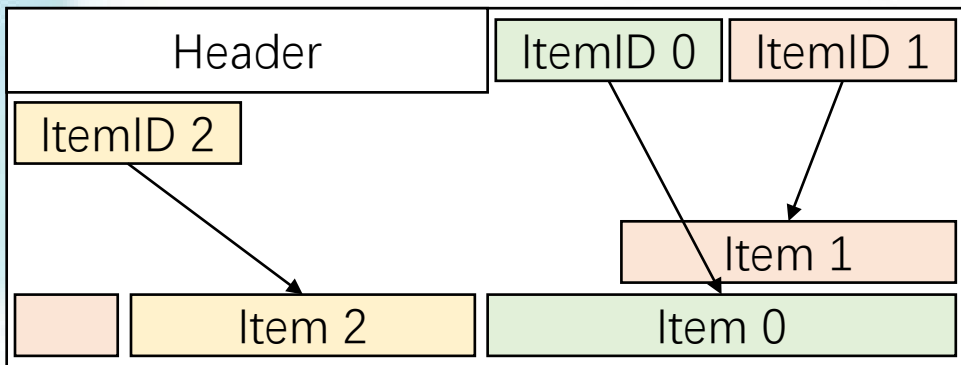


- `ItemPtr get_item(Block *block, short idx);`
  - 返回block第idx项Item的起始地址。
- `short new_item(Block *block, ItemPtr item, short item_size);`
  - 向block插入大小为item\_size, 起始地址为item的Item, 返回插入后的项号 (idx)。
- `void delete_item(Block *block, short idx);`
  - 删除block第idx项Item。

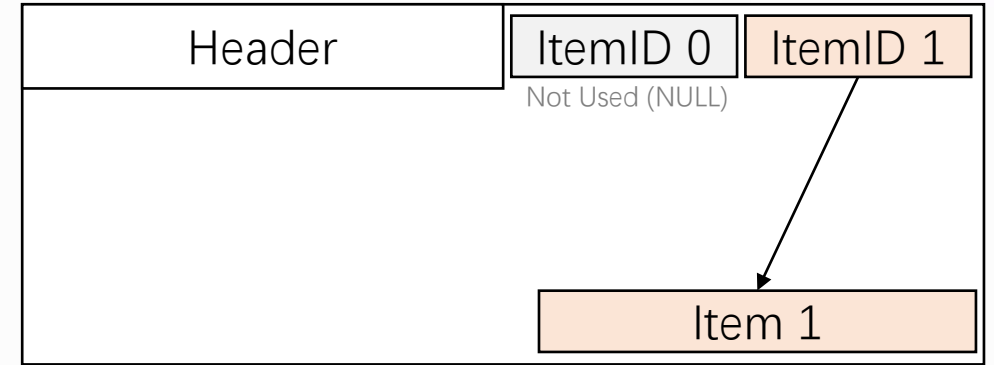
# block.h



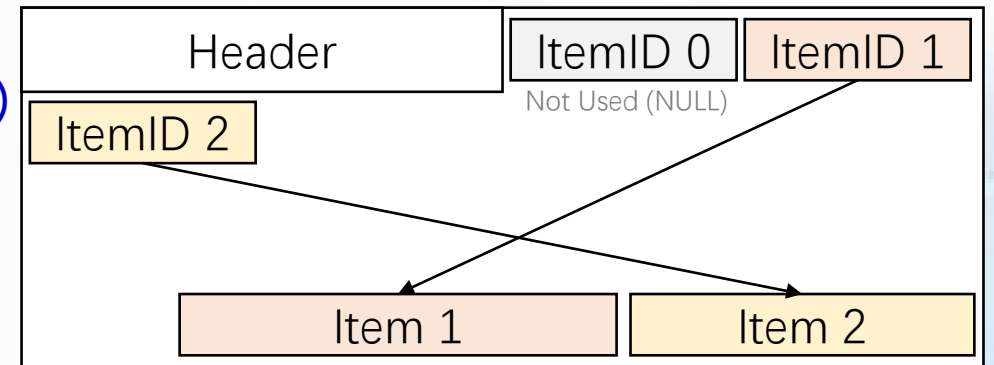
new\_item



delete\_item(0)



delete\_item(2)

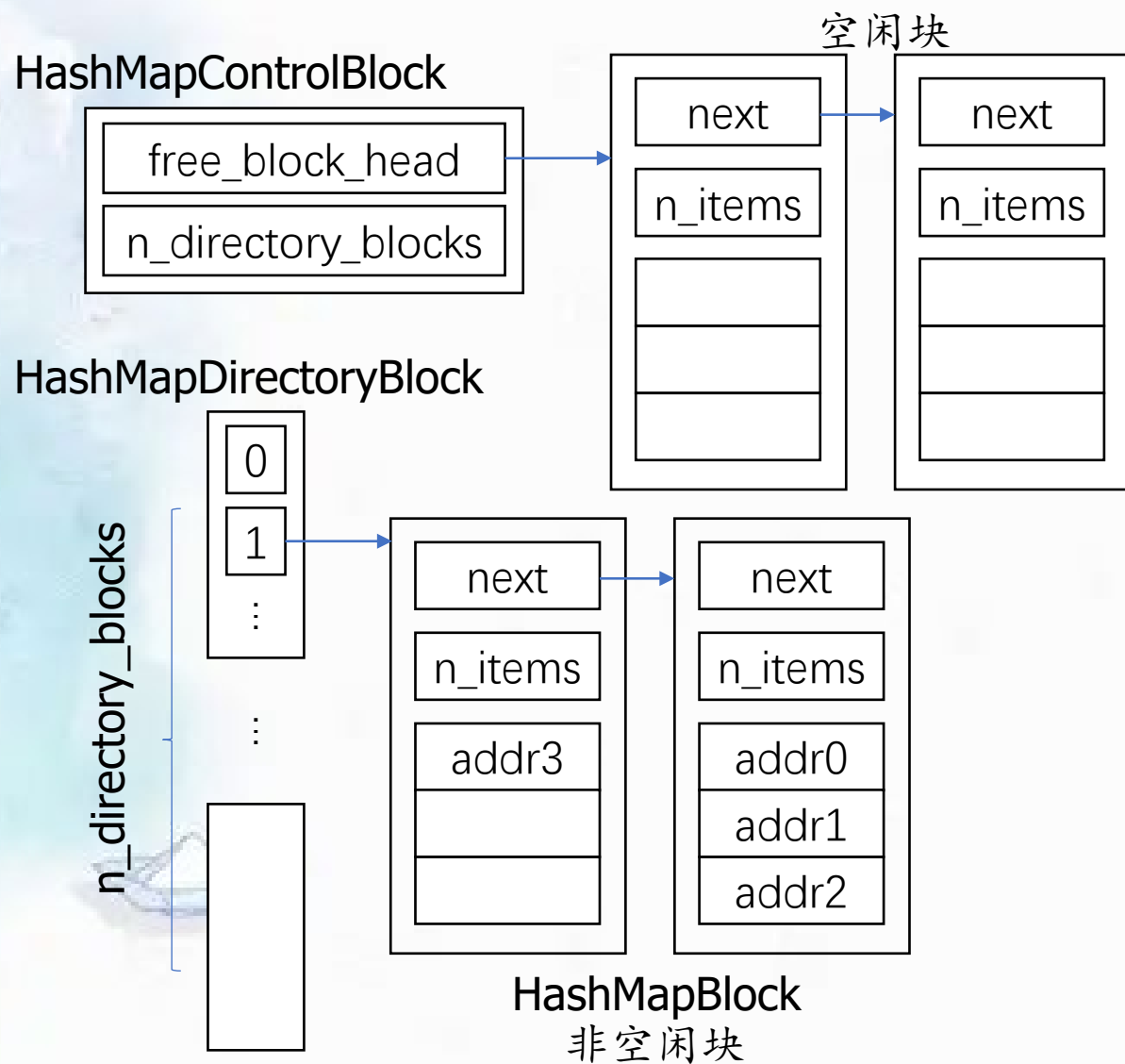


# hash\_map.h

- 管理块的空闲空间。
- 重要函数：
- `void hash_table_insert(BufferPool *pool, short size, off_t addr);`
  - 标记地址为`addr`的块有`size`的空闲空间。
- `off_t hash_table_pop_lower_bound(BufferPool *pool, short size);`
  - 返回至少包含`size`空闲空间的块地址。
- `void hash_table_pop(BufferPool *pool, short size, off_t addr);`
  - 删除地址为`addr`的块的记录，该块有`size`的空闲空间。

[参考链接：操作系统文件存储空间管理（存储空间的划分与初始化）](#)

# hash\_map.h



- `*.fsm`文件由若干块构成。
- 第0块是控制块 (`HashMapControlBlock`)。
- `free_block_head`块保存空闲块链表头指针。
- 第1~`n_directory_blocks`块是哈希表目录块 (`HashMapDirectoryBlock`)，构成哈希表目录，按空闲空间大小分类保存 (`*.data`文件中) 块的地址。
- 其余均为哈希表块 (`HashMapBlock`)，以链表形式组织，保存哈希表块链表头指针。



# table.h



- 保存一个表需要有两个文件，数据文件（**\*.data**）和空闲空间映射文件（**Free Space Map, \*.fsm**，这里用哈希表实现）。
- 重要函数：
- **void table\_read(Table \*table, RID rid, ItemPtr dest);**
  - 根据**rid**，将数据读入**dest**，需要确保**dest**拥有适当的大小。
- **RID table\_insert(Table \*table, ItemPtr src, short size);**
  - 插入大小为**size**，起始地址为**src**的**Item**，返回**rid**。
- **void table\_delete(Table \*table, RID rid);**
  - 根据**rid**，删除相应**Item**。

# str.h

StringChunk:  
存储于文件

rid	size	data
下一块的标识	本块中存储的字符串长度	本块中存储的字符串(变长)

StringRecord:  
存储于内存

chunk	idx
当前块	下一个字符的位置

- `int has_next_char(StringRecord *record);`
  - 判断`record`是否还有下一个字符。
- `char next_char(Table *table, StringRecord *record);`
  - 获取`record`的下一个字符。

- `void read_string(Table *table, RID rid, StringRecord *record);`
  - 根据`rid`, 读取字符串至`record`。
- `RID write_string(Table *table, const char *data, off_t size);`
  - 将长度为`size`的字符串`data`写入`table`表, 并返回`rid`。
- `void delete_string(Table *table, RID rid);`
  - 根据`rid`, 删除对应字符串。

# b\_tree.h

- B-树索引存储于\*.idx文件，其组织方式与哈希表文件 (\*.fsm) 类似。
- 第0块为控制块，包含根节点地址与空闲块链表头指针。
- 其余块为B-树节点块。
- 重要函数：
- RID b\_tree\_search(BufferPool \*pool, void \*key, size\_t size, b\_tree\_ptr\_row\_cmp\_t cmp); /\* 搜索 \*/
- RID b\_tree\_insert(BufferPool \*pool, RID rid, b\_tree\_row\_row\_cmp\_t cmp); /\* 插入 \*/
- void b\_tree\_delete(BufferPool \*pool, RID rid, b\_tree\_row\_row\_cmp\_t cmp); /\* 删除 \*/

参考链接：[B-树](#)、[B+树](#)及其[持久化](#)

# myjql.h

- `size_t myjql_get(const char *key, size_t key_len, char *value, size_t max_size); /* get */`
- `void myjql_set(const char *key, size_t key_len, const char *value, size_t value_len); /* set */`
- `void myjql_del(const char *key, size_t key_len); /* del */`



# Grading

- 满分：10分
- ~~4组测试（每组测试2分，共8分）~~
  - 前三组测试代码公开，最后一组测试代码非公开
  - 1            test\_hash\_map            2分
  - 2            test\_str            2分
  - 3            test\_b\_tree            3分
  - 4            [?]            3分
- 实验报告：对缓冲池、块管理等组件进行有效改进（2分**附加分**，**加到满分10分为止**）
  - 有效改进是指：与原始实现相比，改进后的实现，在随机数据/某些极端情况下的性能有显著提升（平均耗时/空间占用情况）



Q. & A.