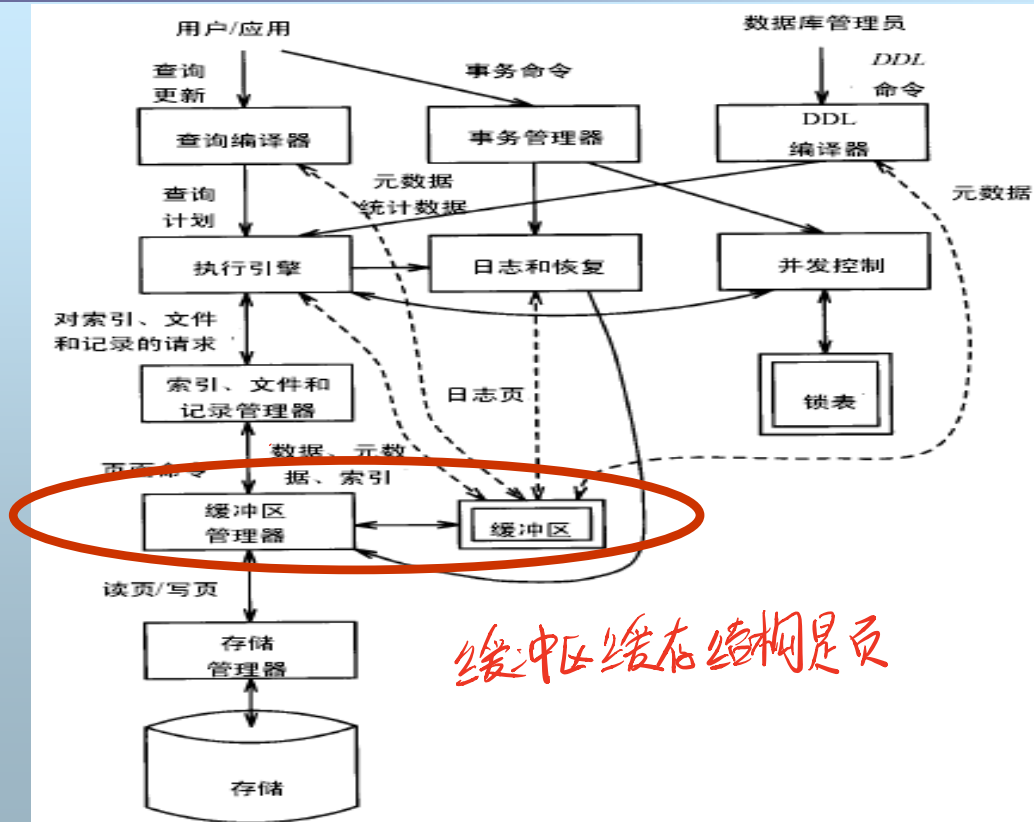


# Buffer Management



**Idea:** Minimize the count of disk I/Os by keeping likely-to-be-requested pages in memory (buffer frames)

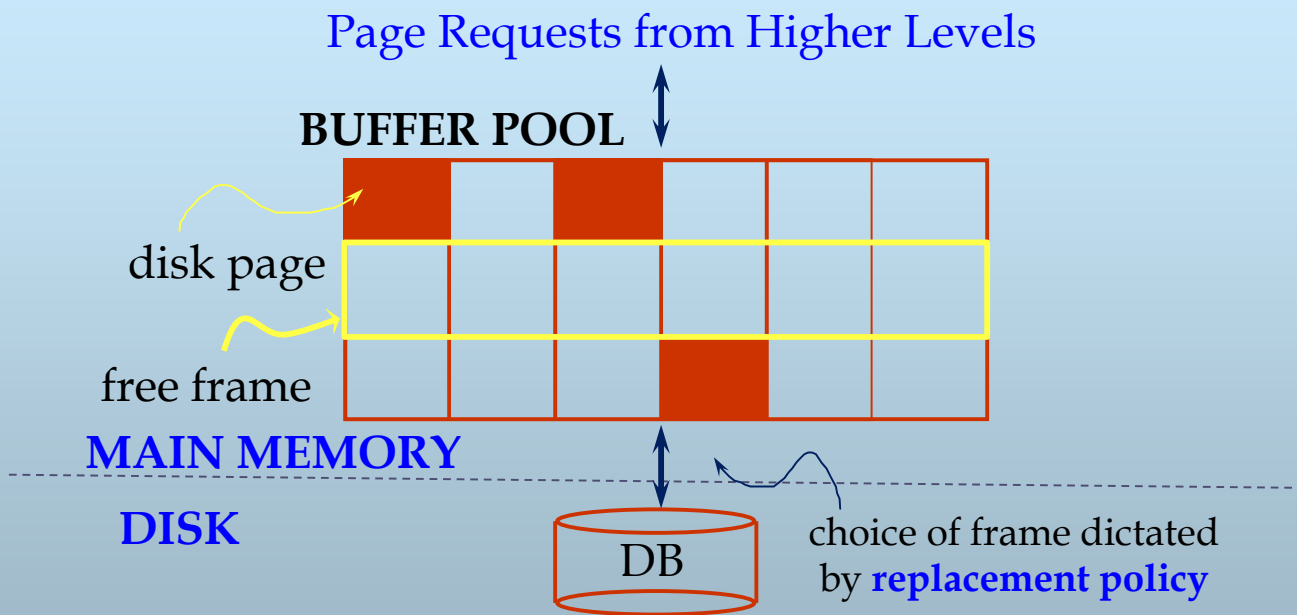


缓冲区缓存结构是页

# 主要内容

- 缓冲区结构
- 缓冲区置换算法
- 缓冲区管理的实现

# 一、缓冲区结构



- *Data must be in RAM for DBMS to operate on it!*
- *Buffer Mgr hides the fact that not all data is in RAM*

# 1、frame的参数

## ■ Dirty

- **Frame**中的块是否已经被修改

## ■ Pin-count

- **Frame**的块的已经被请求并且还未释放的计数，即当前的用户数

## ■ \*Others

- **Latch**: 是否加锁

## 2、当请求块时

- If the requested block is not in the pool:
  - Choose a frame for replacement
  - If the frame is *dirty (some blocks are modified and haven't been written to the disk)*, write it to the disk
  - Read the requested block into the chosen frame
- *Pin (increment the pin-count of the frame)* the block and return its address.

## 2、当释放块时

- Requestor must *unpin* the frame containing the block
- Requestor must indicate whether block has been modified:
  - *dirty* bit is used for this.

## 二、缓冲区替换策略

- Frame is chosen for replacement by a *replacement policy*:
  - Least-recently-used (LRU), Clock, FIFO, MRU (Most-recently-used) etc.
- Only frames whose pin-count=0 are candidates
- Policy can have big impact on # of I/O's; depends on the *access pattern*.



# 1、LRU vs. CLOCK

## ■ LRU (Oracle, Sybase, Informix)

- 当**Pin-count**为**0**时，**frame**放入替换队列
- 选择队列头的**frame**替换

## ■ Clock (MS SQL Server)

- **N**个**frame**组成环形，**current**指针指向当前**frame**；每个**frame**有一个**referenced**位，初始为**1**；
- 当需要置换页时，从**current**开始检查，若**pin-count**>**0**，**current**增加**1**；若**referenced**已启动（=**1**），则关闭它（=**0**）并增加**current**（保证最近的不被替换）；若**pin-count**=**0**并且**referenced**关闭（=**0**），则替换

# 1、LRU vs. CLOCK

## ■ Clock算法的一个示例

- 4 frames, 10 requests

1	2	3	4	1	2	5	1	2	3	4	5
1 1	1 1	1 1	1 1	1 1	1 1	5 1	5 1	5 1	5 1	4 1	4 1
	2 1	2 1	2 1	2 1	2 1	2 0	1 1	1 1	1 1	1 0	5 1
		3 1	3 1	3 1	3 1	3 0	3 0	2 1	2 1	2 0	2 0
			4 1	4 1	4 1	4 0	4 0	4 0	3 1	3 0	3 0

10 page faults

## 2、为何不使用OS缓冲区管理？

- **DBMS经常能预测访问模式(Access Pattern)**
  - 可以使用更专门的缓冲区替换策略
  - 有利于**pre-fetch**策略的有效使用
- **DBMS需要强制写回磁盘能力（如WAL）**  
， **OS**的缓冲写回一般通过记录写请求来实现（来自不同应用），实际的磁盘修改推迟，因此不能保证写顺序

# 三、缓冲区管理器的实现



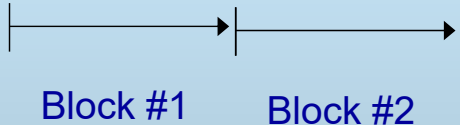
# 1、错误的记录操作实现例子

- 例如，插入记录  
**int insert\_record(DBFILE\*, DBRECORD)**
  - .....
  - **fopen()**
  - **fseek()**
  - **fwrite()**
  - .....
- 没有**DBMS**自己的缓冲区管理和存储管理
- 直接基于文件系统，使用了**FS**的缓冲管理
  - 不能保证**WAL**
  - 不利于查询优化
  - 不适应应用需求

## 2、Block vs. Disk File

### ■ Disk File

0101001001001001011110100111010100101101111.....



文件存储在磁盘上的物理形式是  
**bits/bytes**，**block**是由**OS**或**DBMS**软件  
对文件所做的抽象，这一抽象是通过控制  
数据在文件中的起止**offset**来实现的

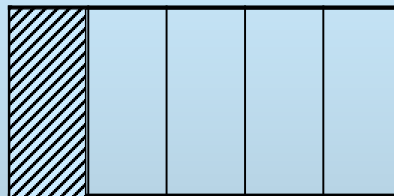
### 3、Buffer vs. Disk File

CPU

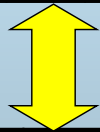
refers to

frame

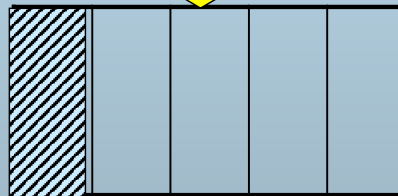
缓冲区管理器



Buffer = set of frames



存储管理器



File = set of pages

page/block

通常,

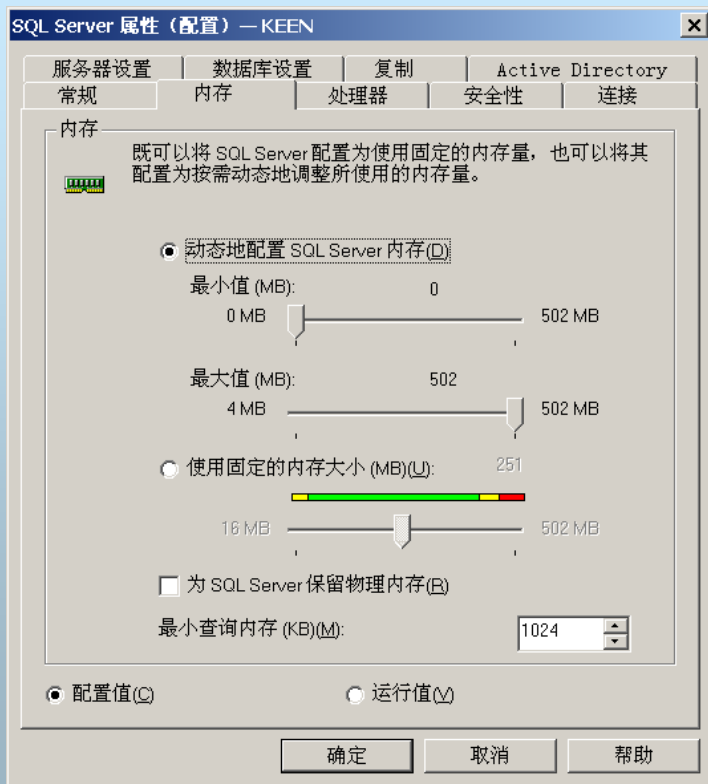
frame大小 =  
page大小

## 4、Buffer Size

- 设计**DBMS**时应是一个可变的输入参数
- 通常**DBMS**允许用户自行配置



# 4、Buffer Size



配置  
实例

## 5、Buffer的存储结构

- **Buffer**是一个**frame**的列表，每个**frame**用于表示和存放一个磁盘块

Buffer的存储结构定义示例

```
#define FRAMESIZE 4096
struct bFrame
{
    Char field [FRAMESIZE ];
};
```

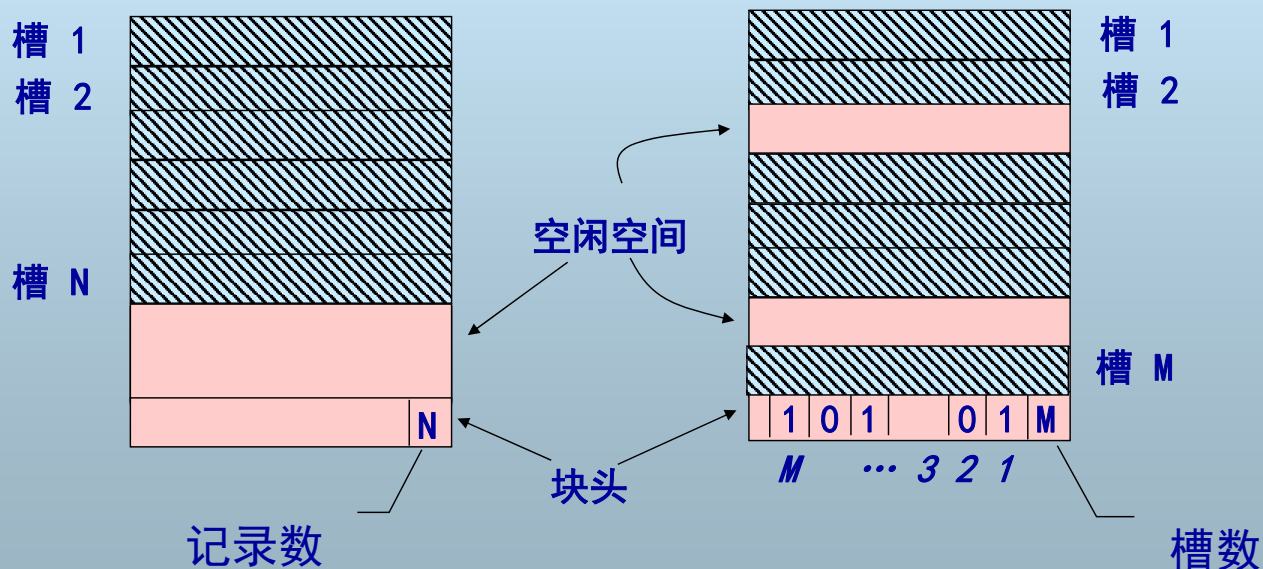
```
#define BUFSIZE 1024 // frame数目

bFrame buf[BUFSIZE];
//也可以为用户配置的值
```

## 6、Page/block的一般存储格式

### ■ 对于定长记录

- 记录地址rid通常使用 <块号, 槽号>表示

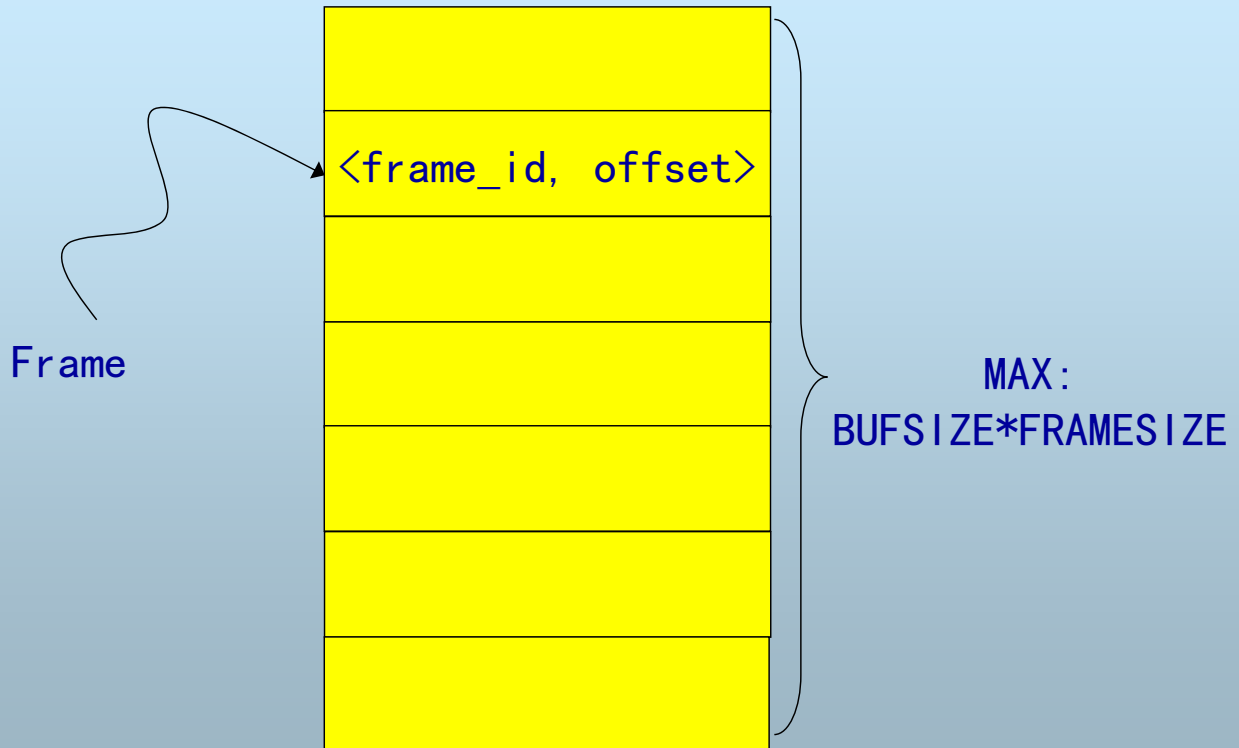


# 6、Page/block的一般存储格式

## ■ Record的存储结构

```
struct Record {  
    int page_id;  
    int slot_num;  
};
```

## 7、Buffer中的Frame存储结构



## 7、Buffer中的Frame存储结构

```
struct Frame{  
    int frame_id;  
    int offset;  
};
```

## 8、Buffer中Frame的查找

- 读磁盘块时：根据page\_id确定在Buffer中是否已经存在frame
- 写磁盘块时：要根据frame\_id快速找到文件中对应的page\_id

## 8、Buffer中Frame的查找

- 首先，要维护Buffer中所有frame的维护信息（**Buffer Control Blocks**），如

```
struct BCB
{
    BCB();
    int page_id;
    int frame_id;
    int count;
    int time;
    int dirty;
    BCB * next;
};
```



## 8、Buffer中Frame的查找

- 建立frame-page之间的索引
- 若用Hash Table, 需要建立2个
  - BCB hTable[BufferSize] //page 2 frame
  - int hTable[BufferSize] //frame 2 page

一个简单的Hash Function例子

$$H(k) = (\text{page\_id}) \% (\text{buffersize})$$

## 9、Buffer Manager的基本功能

- **FixPage(int page\_id)**
  - 将对应`page_id`的`page`读入到`buffer`中。如果`buffer`已满，则需要选择换出的`frame`。
- **FixNewPage()**
  - 在插入数据时申请一个新`page`并将其放在`buffer`中
- **SelectVictim()**
  - 选择换出的`frame_id`
- **FindFrame(int page\_id)**
  - 查找给定的`page`是否已经在某个`frame`中了
- **SetDirty(int frame\_id)**

# 10、数据库文件的一般组成

## ■ 数据文件

- 首块在**Insert\_Record**时创建（调用**Buffer Manager**的**FixNewPage**），一般块号从0开始

## ■ 系统目录文件

- 首块一般**Create\_Table**时创建（调用**Buffer Manager**的**FixNewPage**）

Note:

所有数据和元数据操作都唯一通过  
Buffer Manager来请求page

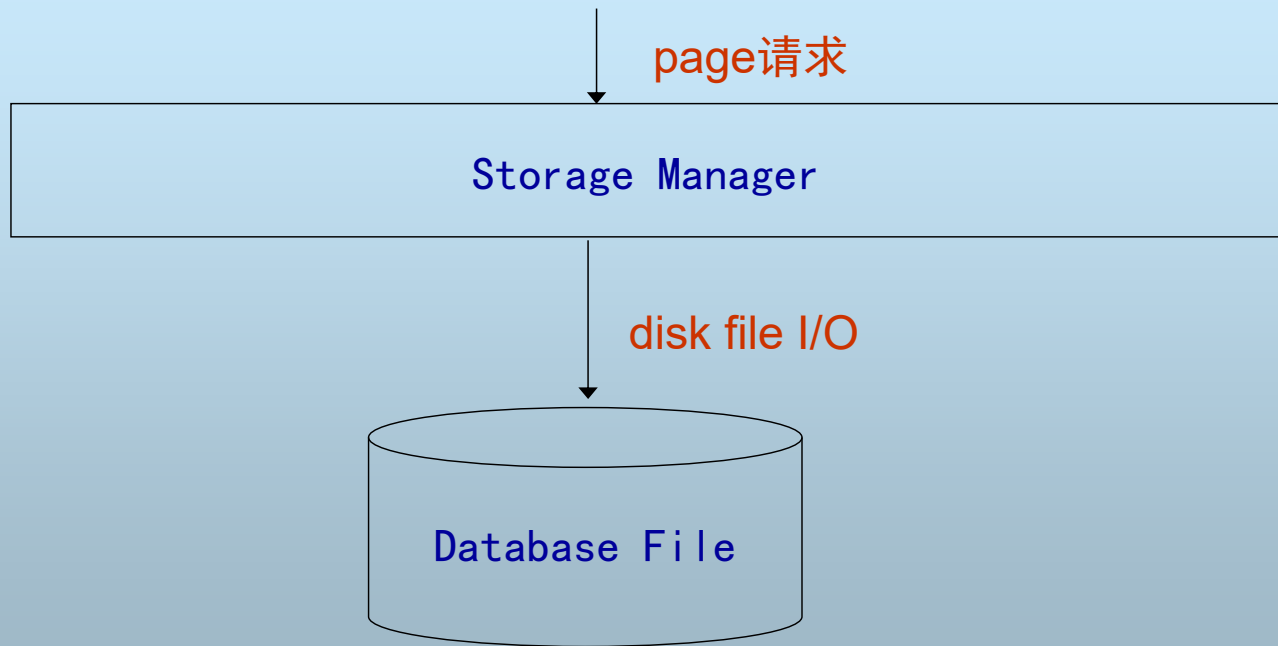
# 11、文件记录操作与Buffer Manager



# 12、存储管理器

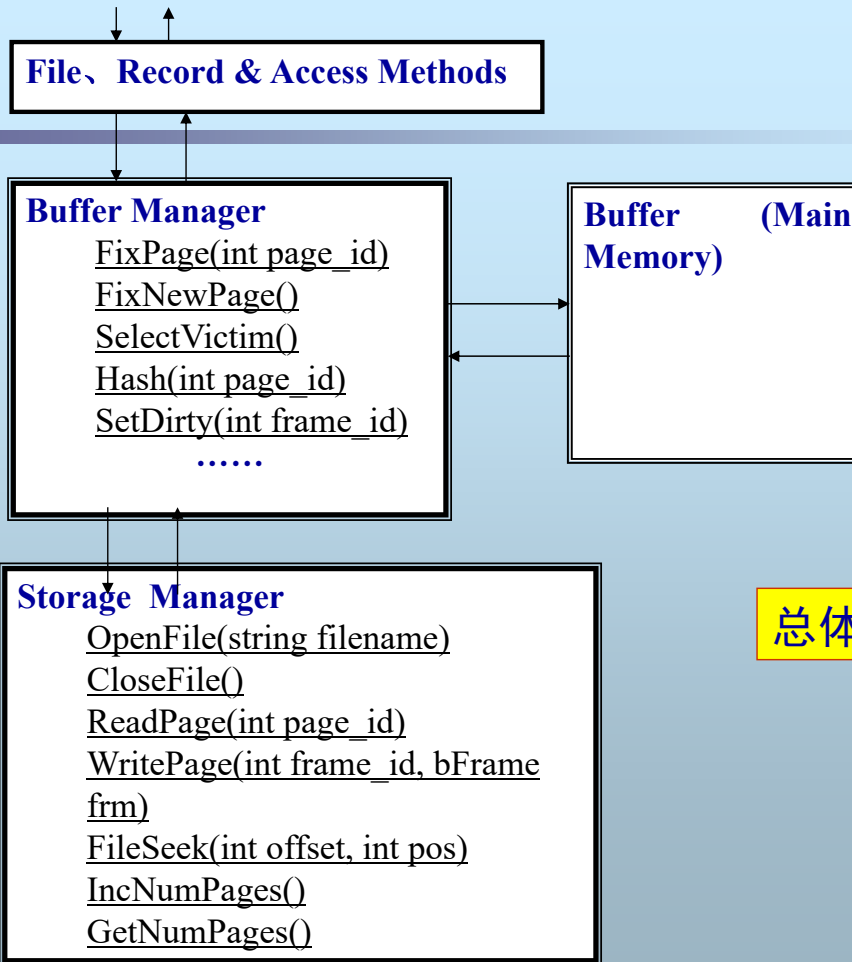


# 12、存储管理器



# 13、存储管理器功能

- 从磁盘中存取物理数据，为**Buffer Manager**提供**Page**抽象
  - **OpenFile/CloseFile**
  - **ReadPage/WritePage**
  - **FileSeek**
  - **GetNumPages**
  - **IncreaseNumPages**
  - .....



总体结构图



# 总结

- 缓冲区结构
- 缓冲区置换算法
- 缓冲区管理的实现