



浙江大學
ZHEJIANG UNIVERSITY

数据库系统

minisql 个人报告

3220103422 刘韬

2024 年 6 月 14 日

目录

1	DISK AND BUFFER POOL MANAGER	3
1.1	Disk Manager summary	3
1.2	Bitmap page	3
1.3	Disk Manager	5
1.4	Buffer pool manager	8
1.4.1	Buffer pool replacement algorithm	8
1.4.2	buffer pool manager implementation	10
1.5	test	14
2	RECORD MANAGER	16
2.1	Record design	16
2.2	Record Persistence	16
2.2.1	Column	16
2.2.2	Schema	18
2.2.3	Row	19
2.3	table heap	20
2.3.1	table heap implementation	20
2.3.2	table heap iterator implementation	24
2.3.3	table heap iterator	27
2.4	test	28
3	INDEX MANAGER	31
3.1	summary	31
3.2	BPlusTree	31
3.2.1	BPlusTreePage	31
3.2.2	BPlusTreeInternalPage	33
3.2.3	BPlusTreeLeafPage	38
3.2.4	BPlusTree implementation	42
3.3	IndexIterator	53
3.3.1	IndexIterator implementation	53
3.3.2	BPlusTree IndexIterator	54
3.4	Test	55
4	Summary	57

个人完成部分

Lab1 及 bonus, Lab2, Lab3, 系统整体功能调试。完善了 Lab2, Lab3 的测试用例。

1 DISK AND BUFFER POOL MANAGER

1.1 Disk Manager summary

为了管理数据库中的文件，我们需要实现一个对磁盘进行读写的类 DiskManager，它的功能是对数据中的 page（基本数据单元）进行 allocation 和 de allocation，执行从磁盘中进行读写 page 的操作。除了存储数据的数据页，我们还需要一些元数据来帮助和加速数据页的管理，于是位图页就自然的出现了。我们将一个位图页和它所管理的一段连续的数据页称作数据库文件中的一个分区（extent）。为了提高存储的容量，我们使用一个额外的元信息页来存储各个分区的信息。这种套娃的方式可以使磁盘文件能够维护更多的数据页信息。下面我们依次实现 DiskManager 中的各个部分，并对其进行测试。



图 1: Disk data

1.2 Bitmap page

首先是位图页的实现。位图页的大小与数据页的大小一致，如2所示，位图页由 Bitmap Page Meta 和 Bitmap Content 两部分组成。Bitmap Page Meta 就是位图页的元数据，含当前已经分配的页的数量（page_allocated_）以及下一个空闲的数据页（next_free_page_）。Bitmap Content 使用 char 数组来记录数据，其容量大小是 $\text{MAX_CHARS} = \text{PageSize} - 2 * \text{sizeof}(\text{uint32_t})$ ，每个 bit 的含义是对应的页是否被分配，一个 char 的大小是 8bits，位图页对外提供分配、回收数据页和判断数据页是否可用的接口，内部实现是对上述数组的管理，以下是各个接口的具体实现。

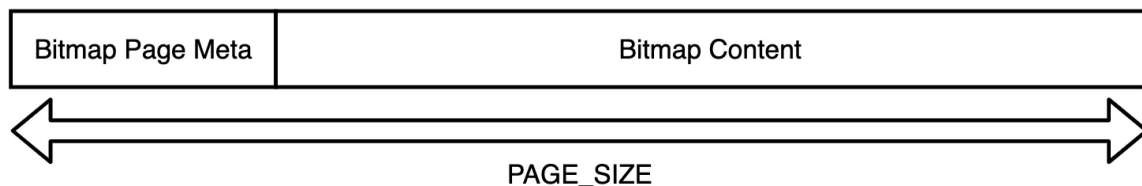


图 2: Bitmap page

AllocatePage 设置相应的next_free_page，并更新free_index来利用离散空间，返回一个 bool 值来说明是否分配成功，并将分配的页号返回。

```

1  template <size_t PageSize>
2  bool BitmapPage<PageSize>::AllocatePage(uint32_t &page_offset) {
3      if(IsPageFree(next_free_page_)){
4          bytes[next_free_page_ / 8] |= (1 << (next_free_page_ % 8)); // 把指定位设为 1
5          page_offset = next_free_page_;
6          page_allocated++;
7          uint32_t free_index = 0;
8          while (!IsPageFree(free_index) && free_index < GetMaxSupportedSize() - 1)
9              {
10                 free_index++;
11             }
12             next_free_page_ = free_index;
13             return true;
14         }
15         return false;
16     }

```

DeAllocatePage DeAllocatePage 与 AllocatePage 相对，将输入的 page_offset 释放，更新 bitmap page 的信息，并更新 next_free_page。

```

1  template <size_t PageSize>
2  bool BitmapPage<PageSize>::DeAllocatePage(uint32_t page_offset) {
3      if(IsPageFree(page_offset)){
4          return false;
5      }
6      bytes[page_offset / 8] &= ~(1 << (page_offset % 8)); // 把指定位设为 0
7      page_allocated--;
8      if(page_offset < next_free_page_){
9          next_free_page_ = page_offset;
10     }
11     return true;
12 }

```

IsPageFree 访问位图页信息，找到对应的 bit 位，根据信息判断是否空闲。需要注意的是应当判断 page_offset 的有效性。

```

1  template <size_t PageSize>
2  bool BitmapPage<PageSize>::IsPageFree(uint32_t page_offset) const {

```

```

3   if (page_offset >= GetMaxSupportedSize())
4   {
5       LOG(ERROR) << "page_offset is out of range";
6       return false;
7   }
8   return IsPageFreeLow(page_offset / 8, page_offset % 8);
9 }
10
11 template <size_t PageSize>
12 bool BitmapPage<PageSize>::IsPageFreeLow(uint32_t byte_index, uint8_t bit_index)
13     ↪ const {
14     return (bytes[byte_index] & (1 << bit_index)) == 0;
15 }

```

1.3 Disk Manager

DiskManager 具有向磁盘读写物理页，分配和回收页的功能。为了更好地管理数据，我们在前面设置了元数据页，但对于上层的程序而言，数据管理器提供的应当是连续的数据页，即逻辑页号。因此我们需要在这里对逻辑页和物理页进行映射，使得上层能够无感知地获得数据页。映射示意图3

物理页号	0	1	2	3	4	5	6	...
职责	磁盘元数据	位图页	数据页	数据页	数据页	位图页	数据页	
逻辑页号	/	/	0	1	2		3	

图 3: phy-log

MapPageId 对输入的逻辑页号返回其物理页号，按照图3的映射完成。这个函数将在物理页的读写中被调用。

```

1   page_id_t DiskManager::MapPageId(page_id_t logical_page_id) {
2       return 1 + logical_page_id / BITMAP_SIZE * (BITMAP_SIZE + 1) + 1 + logical_page_id
3           ↪ % BITMAP_SIZE;
4   }

```

IsPageFree 根据逻辑页号找出该数据页对应的位图页的物理页号，调用位图页的函数来判断是否 free, 需要考虑计算出的物理页号是否有效，如果是无效的超出已经分配的范围，那直接返回 true 即可。

```

1  bool DiskManager::IsPageFree(page_id_t logical_page_id) {
2      BitmapPage<PAGE_SIZE>* bitmap_page = new BitmapPage<PAGE_SIZE>();
3      uint32_t extent_id = logical_page_id / BITMAP_SIZE;
4      uint32_t page_offset = logical_page_id % BITMAP_SIZE;
5      page_id_t bitmap_page_id = 1 + extent_id * (BITMAP_SIZE + 1);
6      // if the extent_id is larger than the number of extents, then the page is free
7      if (extent_id >= reinterpret_cast<DiskFileMetaPage*>(meta_data_)->num_extents_)
8      {
9          return true;
10     }
11     ReadPhysicalPage(bitmap_page_id, reinterpret_cast<char*>(bitmap_page));
12     bool is_free = bitmap_page->IsPageFree(page_offset);
13     delete bitmap_page;
14     return is_free;
15 }

```

AllocatePage 这是磁盘管理器的分配页实现，思路是先找到一个空闲的页，在磁盘元数据页的定义中我们保留了已经分配的页数，总共的 extent 数以及每个 extent 中已经被使用的页数。那么在读取磁盘中的元数据页后，我们通过访问每个 extent 中已经被使用页数的数组可以找到目前已经分配好的 extent 中是否存在空闲页，然后调用位图页的 allocate page 函数进行分配。如果目前所有的页都被分配，说明我们需要再新建一个 extent，然后分配页。同时需要更新元数据页中的数据，并写回磁盘，具体实现如下：

```

1  page_id_t DiskManager::AllocatePage() {
2      // read meta page
3      DiskFileMetaPage* meta_page = reinterpret_cast<DiskFileMetaPage*>(meta_data_);
4      if(meta_page->num_allocated_pages_ == MAX_VALID_PAGE_ID){
5          return INVALID_PAGE_ID;
6      }
7      // find the first free page
8      // if there is free page in existed page
9      for(uint32_t i = 0; i < meta_page->num_extents_; i++){
10         if(meta_page->GetExtentUsedPage(i) < BITMAP_SIZE){
11             BitmapPage<PAGE_SIZE>* bitmap_page = new BitmapPage<PAGE_SIZE>();
12             ReadPhysicalPage(1 + i * (BITMAP_SIZE +
13                 ↪ 1), reinterpret_cast<char*>(bitmap_page));
14             uint32_t page_offset;
15             bool res = bitmap_page->AllocatePage(page_offset);
16             ASSERT(res, "Allocate page failed");
17             if(res){
18                 meta_page->num_allocated_pages_++;

```

```

18     meta_page->extent_used_page_[i]++;
19     WritePhysicalPage(1 + i * (BITMAP_SIZE +
    ↪ 1), reinterpret_cast<char*>(bitmap_page));
20     return i * BITMAP_SIZE + page_offset;
21 }
22 }
23 }
24 // create a new extent
25 BitmapPage<PAGE_SIZE>* bitmap_page = new BitmapPage<PAGE_SIZE>();
26 uint32_t page_offset = 0;
27 bitmap_page->AllocatePage(page_offset);
28 WritePhysicalPage(1 + meta_page->num_extents_ * (BITMAP_SIZE +
    ↪ 1), reinterpret_cast<char*>(bitmap_page));
29 meta_page->num_allocated_pages++;
30 meta_page->extent_used_page_[meta_page->num_extents_]++;
31 meta_page->num_extents_++;
32 return (meta_page->num_extents_ - 1) * BITMAP_SIZE + page_offset;
33 }

```

DeAllocatePage 这部分要做的事情是找到相应的数据页, 调用 Bitmap page 的 DeAllocatePage, 更新元数据, 写回磁盘, 具体实现如下:

```

1 void DiskManager::DeAllocatePage(page_id_t logical_page_id) {
2     uint32_t extent_id = logical_page_id / BITMAP_SIZE;
3     page_id_t bitmap_page_id = 1 + extent_id * (BITMAP_SIZE + 1);
4
5     DiskFileMetaPage* meta_page = reinterpret_cast<DiskFileMetaPage*>(meta_data_);
6     BitmapPage<PAGE_SIZE>* bitmap_page = new BitmapPage<PAGE_SIZE>();
7
8     ReadPhysicalPage(bitmap_page_id, reinterpret_cast<char*>(bitmap_page));
9     uint32_t page_offset = logical_page_id % BITMAP_SIZE;
10    bool res = bitmap_page->DeAllocatePage(page_offset);
11    if(res){
12        meta_page->num_allocated_pages--;
13        meta_page->extent_used_page_[extent_id]--;
14        WritePhysicalPage(bitmap_page_id, reinterpret_cast<char*>(bitmap_page));
15    }else{
16        LOG(ERROR) << "Deallocate page failed." << logical_page_id;
17    }
18 }

```

1.4 Buffer pool manager

Buffer pool manager 是提供给上层调用，用于数据读写管理的接口。功能需求是

1. `Page *FetchPage(page_id_t page_id);` 根据逻辑页号获取对应的数据页，如果该数据页不在内存中，则需要从磁盘中进行读取；
2. `BufferPoolManager::NewPage(&page_id):` 分配一个新的数据页，并将逻辑页号于 `page_id` 中返回；
3. `BufferPoolManager::UnpinPage(page_id, is_dirty):` 取消固定一个数据页；
4. `BufferPoolManager::FlushPage(page_id):` 将数据页转储到磁盘中；
5. `BufferPoolManager::DeletePage(page_id):` 释放一个数据页；

1.4.1 Buffer pool replacement algorithm

为了保证在 buffer pool 没有空闲页的情况下可以正常的使用，我们需要一个替换单元，这里使用 LRU 替换算法来决定哪一页被替换，同时还实现了 Bonus: Clock Replacer 并为其实现测试用例。以下是两种替换算法的实现：

Least Recently Used Replacer LRU 替换算法是将最近最少被访问的元素作为被替代的对象，因此实现时我们可以通过一个有序列表，将每次使用后的都换到列表的头部，那么列表的尾部就是我们所需要的对象。我们使用了 `list` 来存保证有序性，用 `unordered_map` 来加快查询。

Victim 替换（即删除）与所有被跟踪的页相比最近最少被访问的页将其页帧号（即数据页在 Buffer Pool 的 Page 数组中的下标）存储在输出参数 `frame_id` 中输出并返回 `true`，如果当前没有可以替换的元素则返回 `false`；

```

1  bool LRUReplacer::Victim(frame_id_t *frame_id) {
2      if (lru_list_.empty()) {
3          return false;
4      }
5      *frame_id = lru_list_.back();
6      lru_list_.pop_back();
7      lru_list_map_.erase(*frame_id);
8      return true;
9  }
```


Pin 将数据页固定使之不能被 Replacer 替换，即从 lru_list_ 中移除该数据页对应的页帧。

```

1 void LRUReplacer::Pin(frame_id_t frame_id) {
2     auto it = lru_list_map_.find(frame_id);
3     if(it != lru_list_map_.end()){
4         lru_list_.erase(it->second);
5         lru_list_map_.erase(it);
6     }
7 }

```

UnPin 将数据页解除固定，放入 lru_list_ 中，使之可以在必要时被 Replacer 替换掉。

```

1 void LRUReplacer::Unpin(frame_id_t frame_id) {
2     auto it = lru_list_map_.find(frame_id);
3     if(it == lru_list_map_.end()){
4         lru_list_.push_front(frame_id);
5         lru_list_map_[frame_id] = lru_list_.begin();
6     }
7 }

```

Clock Replacer 这种替换方式采用标记，对每个元素进行标记，然后每次选择替换时进行一次遍历，对找到的每个元素查看标记，如果标记是 1 就变为 0，如果标记是 0 就替换出，每次查询是上次结束的位置开始，对新入队的元素标记是 1。类的成员变量有 list 用于存储替换页和 map 用于存储标记状态。

Victim 替换（即删除）与所有被跟踪的页相比最近最少被访问的页将其页帧号（即数据页在 Buffer Pool 的 Page 数组中的下标）存储在输出参数 frame_id 中输出并返回 true，如果当前没有可以替换的元素则返回 false；

```

1 bool CLOCKReplacer::Victim(frame_id_t *frame_id) {
2     // if the clock_list is empty, return false
3     if (clock_list.empty()) {
4         return false;
5     }
6     // find the first frame that is not pinned
7     while (clock_status[clock_list.front()] == 1) {
8         clock_list.push_back(clock_list.front());
9         clock_list.pop_front();
10        clock_status[clock_list.back()] = 0;

```

```

11     }
12     // return the frame_id
13     *frame_id = clock_list.front();
14     clock_list.pop_front();
15     clock_status.erase(*frame_id);
16     return true;
17 }

```

Pin 将数据页固定使之不能被 Replacer 替换，即从 lru_list_ 中移除该数据页对应的页帧。

```

1 void CLOCKReplacer::Pin(frame_id_t frame_id) {
2     // if the frame is in the clock_list, then remove it from the clock_list
3     if (clock_status.find(frame_id) != clock_status.end()) {
4         clock_list.remove(frame_id);
5         clock_status.erase(frame_id);
6     }
7 }
8

```

1.4.2 buffer pool manager implementation

UnPin 将数据页解除固定，放入 lru_list_ 中，使之可以在必要时被 Replacer 替换掉。

```

1 void CLOCKReplacer::Unpin(frame_id_t frame_id) {
2     // if the frame is not in the clock_list, then add it to the clock_list
3     if (clock_status.find(frame_id) == clock_status.end()) {
4         clock_list.push_back(frame_id);
5         clock_status[frame_id] = 1;
6     }else{
7         clock_status[frame_id] = 1;
8     }
9 }

```

FetchPage 首先尝试在 page table 中找到要求的 page(P)，如果存在就固定它并立刻返回这个 Page；如果不存在则需要从 free list（优先）或者 replacer 中找到一个替换页(R)。如果 R 是脏页，就应当先写回磁盘保证数据的一致性，然后把 R 从 page table 中删掉并且插入 P，并更新 P 的元数据，从磁盘中读出数据即可。

```

1 Page *BufferPoolManager::FetchPage(page_id_t page_id) {
2     // 1.    Search the page table for the requested page (P).
3     frame_id_t frame_id = INVALID_FRAME_ID;
4     if(page_table_.find(page_id) != page_table_.end()) {
5         // 1.1    If P exists, pin it and return it immediately.
6         frame_id = page_table_[page_id];
7         pages_[frame_id].pin_count_++;
8         replacer_>Pin(frame_id);
9         return &pages_[frame_id];
10    }else{
11        // 1.2    If P does not exist, find a replacement page (R) from either the free
        ↪ list or the replacer.
12        //        Note that pages are always found from the free list first.
13        if(!free_list_.empty()){
14            frame_id = free_list_.front();
15            free_list_.pop_front();
16        }else{
17            // get from replacer
18            if(!replacer_>Victim(&frame_id)){
19                return nullptr;
20            }
21        }
22        // 2.    If R is dirty, write it back to the disk.
23        if(pages_[frame_id].IsDirty()){
24            FlushPage(pages_[frame_id].GetPageId());
25        }
26        // 3.    Delete R from the page table and insert P.
27        page_table_.erase(pages_[frame_id].GetPageId());
28        page_table_[page_id] = frame_id;
29        // 4.    Update P's metadata, read in the page content from disk, and then
        ↪ return a pointer to P.
30        pages_[frame_id].pin_count_ = 1;
31        pages_[frame_id].page_id_ = page_id;
32        pages_[frame_id].ResetMemory();
33        disk_manager_>ReadPage(page_id,pages_[frame_id].GetData());
34        replacer_>Pin(frame_id);
35        return &pages_[frame_id];
36    }
37 }

```

NewPage 如果已经没有空闲页（free list 和 replacer 都没有了），就返回空指针；如果还有就取出一个，然后更新元数据，清空页面的原数据，固定并返回 page。

```

1 Page *BufferPoolManager::NewPage(page_id_t &page_id) {
2     // 0. Make sure you call AllocatePage!
3     // 1. If all the pages in the buffer pool are pinned, return nullptr.
4     if(free_list_.empty() && replacer_>Size() == 0){
5         return nullptr;
6     }
7     // 2. Pick a victim page P from either the free list or the replacer. Always pick
8     ↪ from the free list first.
9     frame_id_t frame_id = INVALID_FRAME_ID;
10    if(!free_list_.empty()){
11        frame_id = free_list_.front();
12        free_list_.pop_front();
13    }else{
14        if(!replacer_>Victim(&frame_id)){
15            return nullptr;
16        }
17        // if the page get from replacer, we should erase it and get space for the new
18        ↪ page.
19        page_table_.erase(pages_[frame_id].GetPageId());
20        // 3. Update P's metadata, zero out memory and add P to the page table.
21        page_id = AllocatePage();
22        pages_[frame_id].page_id_ = page_id;
23        pages_[frame_id].pin_count_ = 1;
24        pages_[frame_id].is_dirty_ = false;
25        pages_[frame_id].ResetMemory();
26        page_table_[page_id] = frame_id;
27        replacer_>Pin(frame_id);
28        // 4. Set the page ID output parameter. Return a pointer to P.
29        return &pages_[frame_id];
30    }

```

DeletePage 找到需要删除的 page，如果不存在返回 true，如果存在但是仍然被固定，返回 false，除此以外这个 page 可以被删除，从 table 中删除，reset 数据，放回 free list。

```

1 bool BufferPoolManager::DeletePage(page_id_t page_id) {
2     // 0. Make sure you call DeallocatePage!
3     // 1. Search the page table for the requested page (P).
4     // 1. If P does not exist, return true.
5     if(page_table_.find(page_id) == page_table_.end()){
6         return true;

```

```

7   }else{
8   // 2.   If P exists, but has a non-zero pin-count, return false. Someone is using
        ↪ the page.
9       frame_id_t frame_id = page_table_[page_id];
10      if(pages_[frame_id].GetPinCount() != 0){
11          return false;
12      }else{
13      // 3.   Otherwise, P can be deleted. Remove P from the page table, reset its
        ↪ metadata and return it to the free list.
14          page_table_.erase(page_id);
15          DeallocatePage(page_id);
16          pages_[frame_id].ResetMemory();
17          pages_[frame_id].page_id_ = INVALID_PAGE_ID;
18          pages_[frame_id].pin_count_ = 0;
19          pages_[frame_id].is_dirty_ = false;
20          free_list_.push_back(frame_id);
21          replacer_>Pin(frame_id);
22          return true;
23      }
24  }
25  }

```

UnpinPage 在 pagetable 中找到这个页,对是否要写回磁盘进行标记,然后将 pin_count-1, 如果 pin 数为 0 就把它从 replacer 中 unpin 掉。

```

1  bool BufferPoolManager::UnpinPage(page_id_t page_id, bool is_dirty) {
2      if (page_table_.find(page_id) == page_table_.end()){
3          return true;
4      }
5      frame_id_t frame_id = page_table_[page_id];
6      if(pages_[frame_id].pin_count_ > 0){
7          pages_[frame_id].pin_count_ --;
8      }
9      pages_[frame_id].is_dirty_ |= is_dirty;
10     if(pages_[frame_id].pin_count_ == 0){
11         replacer_>Unpin(frame_id);
12     }
13     return true;
14 }

```

FlushPage 找到 page，将其写回磁盘，无论是否 pin。这是为了保证磁盘数据的一致性。

```
1 bool BufferPoolManager::FlushPage(page_id_t page_id) {
2     if(page_table_.find(page_id) == page_table_.end()){
3         return false;
4     }else{
5         frame_id_t frame_id = page_table_[page_id];
6         disk_manager_>WritePage(page_id,pages_[frame_id].GetData());
7         pages_[frame_id].is_dirty_ = false;
8         return true;
9     }
10 }
```

1.5 test

使用所给的测试用例进行测试，增加了 clock replacer 的测试，其测试与 LRU 测试相同。测试结果如下：

本节的测试全都通过，包括 bonus 的测试。

```
> ./test/disk_manager_test
[=====] Running 2 tests from 1 test suite.
[-----] Global test environment set-up.
[-----] 2 tests from DiskManagerTest
[ RUN      ] DiskManagerTest.BitMapPageTest
[      OK  ] DiskManagerTest.BitMapPageTest (162 ms)
[ RUN      ] DiskManagerTest.FreePageAllocationTest
[      OK  ] DiskManagerTest.FreePageAllocationTest (24502 ms)
[-----] 2 tests from DiskManagerTest (24664 ms total)

[-----] Global test environment tear-down
[=====] 2 tests from 1 test suite ran. (24664 ms total)
[ PASSED  ] 2 tests.
```

图 4: test1

```
> ./test/lru_replacer_test
[=====] Running 1 test from 1 test suite.
[-----] Global test environment set-up.
[-----] 1 test from LRURewriterTest
[ RUN      ] LRURewriterTest.SampleTest
[      OK  ] LRURewriterTest.SampleTest (0 ms)
[-----] 1 test from LRURewriterTest (0 ms total)

[-----] Global test environment tear-down
[=====] 1 test from 1 test suite ran. (0 ms total)
[ PASSED  ] 1 test.
```

图 5: test2

```
> ./test/clock_replacer_test
[=====] Running 1 test from 1 test suite.
[-----] Global test environment set-up.
[-----] 1 test from CLOCKReplacerTest
[ RUN      ] CLOCKReplacerTest.SampleTest
[      OK  ] CLOCKReplacerTest.SampleTest (0 ms)
[-----] 1 test from CLOCKReplacerTest (0 ms total)

[-----] Global test environment tear-down
[=====] 1 test from 1 test suite ran. (0 ms total)
[ PASSED  ] 1 test.
```

图 6: test3

```
> ./test/buffer_pool_manager_test
[=====] Running 1 test from 1 test suite.
[-----] Global test environment set-up.
[-----] 1 test from BufferPoolManagerTest
[ RUN      ] BufferPoolManagerTest.BinaryDataTest
[      OK  ] BufferPoolManagerTest.BinaryDataTest (7 ms)
[-----] 1 test from BufferPoolManagerTest (7 ms total)

[-----] Global test environment tear-down
[=====] 1 test from 1 test suite ran. (7 ms total)
[ PASSED  ] 1 test.
```

图 7: test4

2 RECORD MANAGER

2.1 Record design

我们的记录包含以下几个概念：

- **Column**：用于定义和表示数据表中的某一个字段，即包含了这个字段的字段名、字段类型、是否唯一等等；
- **Schema**：用于表示一个数据表或是一个索引的结构。一个 Schema 由一个或多个的 Column 构成；
- **Field**：它对应于一条记录中某一个字段的数据信息，如存储数据的数据类型，是否是空，存储数据的值等等；
- **Row**：用于表示一条记录，即一行数据，由多个 Field 组成；

这是我们数据在数据库中的组织方式，但是为了将其持久化地存储在硬盘中，我们需要使用序列化和反序列化地方法，将其转化为字节流。下面的章节介绍了我们如何进行序列化和反序列化。

2.2 Record Persistence

2.2.1 Column

Column 的结构是 |COLUMN_MAGIC_NUM 4 | name_ length 4 | name.strlen() | | type_ 4 | len_ 4 | table_ind_ 4 | nullable_ 1 | unique_ 1 |

Serialize 逐个写入每个字段，然后将其转化为字节流。

```

1  uint32_t Column::SerializeTo(char *buf) const {
2      uint32_t offset = 0;
3      memcpy(buf + offset, &COLUMN_MAGIC_NUM, sizeof(COLUMN_MAGIC_NUM));
4      offset += sizeof(COLUMN_MAGIC_NUM);
5      uint32_t name_len = name_.length();
6      memcpy(buf + offset, &name_len, sizeof(name_len));
7      offset += sizeof(name_len);
8      memcpy(buf + offset, name_.c_str(), name_len);
9      offset += name_len;
10     memcpy(buf + offset, &type_, sizeof(type_));
11     offset += sizeof(type_);
12     memcpy(buf + offset, &len_, sizeof(len_));
13     offset += sizeof(len_);
14     memcpy(buf + offset, &table_ind_, sizeof(table_ind_));

```



```

15     offset += sizeof(table_ind_);
16     memcpy(buf + offset, &nullable_, sizeof(nullable_));
17     offset += sizeof(nullable_);
18     memcpy(buf + offset, &unique_, sizeof(unique_));
19     offset += sizeof(unique_);
20     return offset;
21 }

```

Deserialize 逐个读取每个字段，然后将其转化为 Column。

```

1  uint32_t Column::DeserializeFrom(char *buf, Column *&column) {
2      if (column != nullptr) {
3          LOG(WARNING) << "Pointer to column is not null in column deserialize." <<
4              ↪ std::endl;
5      }
6      uint32_t offset = 0;
7      uint32_t magic_num;
8      memcpy(&magic_num, buf + offset, sizeof(magic_num));
9      offset += sizeof(magic_num);
10     ASSERT(magic_num == COLUMN_MAGIC_NUM, "Invalid magic number.");
11     if (magic_num != COLUMN_MAGIC_NUM) {
12         return 0;
13     }
14     uint32_t name_len;
15     memcpy(&name_len, buf + offset, sizeof(name_len));
16     offset += sizeof(name_len);
17     std::string name(buf + offset, name_len);
18     offset += name_len;
19     TypeId type;
20     memcpy(&type, buf + offset, sizeof(type));
21     offset += sizeof(type);
22     uint32_t len;
23     memcpy(&len, buf + offset, sizeof(len));
24     offset += sizeof(len);
25     uint32_t table_ind;
26     memcpy(&table_ind, buf + offset, sizeof(table_ind));
27     offset += sizeof(table_ind);
28     bool nullable;
29     memcpy(&nullable, buf + offset, sizeof(nullable));
30     offset += sizeof(nullable);
31     bool unique;
32     memcpy(&unique, buf + offset, sizeof(unique));
33     offset += sizeof(unique);

```

```

33     if (type == TypeId::kTypeChar) {
34         column = new Column(name, type, len, table_ind, nullable, unique);
35     } else {
36         column = new Column(name, type, table_ind, nullable, unique);
37     }
38     return offset;
39 }

```

2.2.2 Schema

Schema 的结构是 SCHEMA_MAGIC_NUM 4 | column size 4 | column1 | column2 | ... ||

Serialize 逐个写入每个字段，然后将其转化为字节流。

```

1  uint32_t Schema::SerializeTo(char *buf) const {
2      uint32_t offset = 0;
3      memcpy(buf + offset, &SCHEMA_MAGIC_NUM, sizeof(SCHEMA_MAGIC_NUM));
4      offset += sizeof(SCHEMA_MAGIC_NUM);
5      uint32_t column_size = columns_.size();
6      memcpy(buf + offset, &column_size, sizeof(column_size));
7      offset += sizeof(column_size);
8      for (auto column : columns_) {
9          offset += column->SerializeTo(buf + offset);
10     }
11     return offset;
12 }

```

Deserialize 逐个读取每个字段，然后将其转化为 Schema。

```

1  uint32_t Schema::DeserializeFrom(char *buf, Schema *&schema) {
2      uint32_t offset = 0;
3      uint32_t magic_num;
4      memcpy(&magic_num, buf + offset, sizeof(magic_num));
5      offset += sizeof(magic_num);
6      ASSERT(magic_num == SCHEMA_MAGIC_NUM, "Invalid magic number.");
7      if (magic_num != SCHEMA_MAGIC_NUM) {
8          return 0;
9      }
10     uint32_t column_size;
11     memcpy(&column_size, buf + offset, sizeof(column_size));

```

```

12  offset += sizeof(column_size);
13  std::vector<Column *> columns;
14  for (uint32_t i = 0; i < column_size; i++) {
15      Column *temp = nullptr;
16      offset += Column::DeserializeFrom(buf + offset, temp);
17      columns.emplace_back(temp);
18  }
19  schema = new Schema(columns, true);
20  return offset;
21  }

```

2.2.3 Row

对于 Row 类型对象的序列化，对于为 null 的 Field，可以通过位图的方式标记 (即 Null Bitmaps)；对于 Row 类型对象的反序列化，在反序列化每一个 Field 时，反序列化出来的 Field 的内存都由该 Row 对象维护。

Serialize 写入长度 num，写入位图页，然后写入每一个 Field。

```

1  uint32_t Schema::SerializeTo(char *buf) const {
2      // | SCHEMA_MAGIC_NUM 4 | column size 4 | column1 | column2 | ... |
3      uint32_t offset = 0;
4      memcpy(buf + offset, &SCHEMA_MAGIC_NUM, sizeof(SCHEMA_MAGIC_NUM));
5      offset += sizeof(SCHEMA_MAGIC_NUM);
6      uint32_t column_size = columns_.size();
7      memcpy(buf + offset, &column_size, sizeof(column_size));
8      offset += sizeof(column_size);
9      for (auto column : columns_) {
10         offset += column->SerializeTo(buf + offset);
11     }
12     return offset;
13 }

```

Deserialize 读取长度 num，读取位图页，然后读取每一个 Field。

```

1  uint32_t Schema::DeserializeFrom(char *buf, Schema *&schema) {
2      uint32_t offset = 0;
3      uint32_t magic_num;
4      memcpy(&magic_num, buf + offset, sizeof(magic_num));
5      offset += sizeof(magic_num);
6      ASSERT(magic_num == SCHEMA_MAGIC_NUM, "Invalid magic number.");

```

```

7   if (magic_num != SCHEMA_MAGIC_NUM) {
8       return 0;
9   }
10  uint32_t column_size;
11  memcpy(&column_size, buf + offset, sizeof(column_size));
12  offset += sizeof(column_size);
13  std::vector<Column*> columns;
14  for (uint32_t i = 0; i < column_size; i++) {
15      Column *temp = nullptr;
16      offset += Column::DeserializeFrom(buf + offset, temp);
17      columns.emplace_back(temp);
18  }
19  schema = new Schema(columns, true);
20  return offset;
21 }

```

2.3 table heap

对于数据表中的每一行记录，都有一个唯一标识符 RowId 与之对应。RowId 同时具有逻辑和物理意义，在物理意义上，它是一个 64 位整数，是每行记录的唯一标识；RowId 的作用主要体现在两个方面：一是在索引中，叶结点中存储的键值对是索引键 Key 到 RowId 的映射，通过索引键 Key，沿着索引查找，我们能够得到该索引键对应记录的 RowId，也就能够在堆表中定位到该记录；二是在堆表中，借助 RowId 中存储的逻辑信息（page_id 和 slot_num），可以快速定位到其对应的记录位于物理文件的哪个位置。而在逻辑意义上，它的高 32 位存储的是该 RowId 对应记录所在数据页的 page_id，低 32 位存储的是该 RowId 在 page_id 对应的数据页中对应的是第几条记录。我们可以通过 RowId 的逻辑信息，得到该 Row 的数据。

我们需要一个 table heap 来存储数据表中的记录，table heap 的结构如下：堆表中的每个数据页都由表头（Table Page Header）、空闲空间（Free Space）和已经插入的数据（Inserted Tuples）三部分组成，表头在页中从左往右扩展，记录了 PrevPageId、NextPageId、FreeSpacePointer 以及每条记录在 TablePage 中的偏移和长度；插入的记录在页中从右向左扩展，每次插入记录时会将 FreeSpacePointer 的位置向左移动。

2.3.1 table heap implementation

实现堆表的插入、删除、查询等操作，具体实现如下：

InsertTuple 向堆表中插入一条记录，插入记录后生成的 RowId 需要通过 row 对象返回（即 row.rid_）；我们首先判断插入的记录是否能够放入当前页中，如果超出当前页的空间限制，我们需要将记录插入到下一个页中。

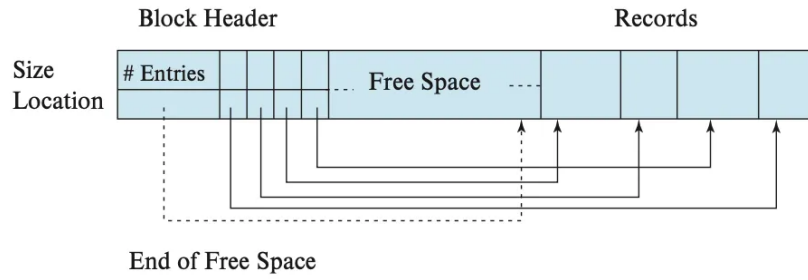


Figure 13.6 Slotted-page structure.

图 8: Table Page

```

1  bool TableHeap::InsertTuple(Row &row, Txn *txn) {
2      if(row.GetSerializedSize(schema_) >= PAGE_SIZE) return false;
3      auto page = reinterpret_cast<TablePage>
4      ↪ *(>(buffer_pool_manager_>FetchPage(GetFirstPageId())));
5      if (page == nullptr) {
6          return false;
7      }
8      while(!page->InsertTuple(row, schema_, txn, lock_manager_, log_manager_)){
9          buffer_pool_manager_>UnpinPage(page->GetTablePageId(), false);
10         page_id_t next_page_id = page->GetNextPageId();
11         if(next_page_id == INVALID_PAGE_ID){
12             auto new_page = reinterpret_cast<TablePage>
13             ↪ *(>(buffer_pool_manager_>NewPage(next_page_id)));
14             page->SetNextPageId(next_page_id);
15             new_page->Init(next_page_id, page->GetTablePageId(), log_manager_, txn);
16             buffer_pool_manager_>UnpinPage(page->GetPageId(), true);
17             page = new_page;
18         }else{
19             auto next_page = reinterpret_cast<TablePage>
20             ↪ *(>(buffer_pool_manager_>FetchPage(next_page_id)));
21             buffer_pool_manager_>UnpinPage(page->GetPageId(), false);
22             page = next_page;
23         }
24     }
25     buffer_pool_manager_>UnpinPage(page->GetTablePageId(), true);
26     return true;
27 }

```

ApplyDelete 删除记录是采用先标记后统一执行输出的方式进行删除，这里实现的是执行删除的部分。

```

1 void TableHeap::ApplyDelete(const RowId &rid, Txn *txn) {
2     // Step1: Find the page which contains the tuple.
3     auto page = reinterpret_cast<TablePage>
4         ↳ *(buffer_pool_manager_>FetchPage(rid.GetPageId()));
5     ASSERT(page != nullptr, "The page could not be found.");
6     // Step2: Delete the tuple from the page.
7     page->ApplyDelete(rid, txn, log_manager_);
8     buffer_pool_manager_>UnpinPage(page->GetTablePageId(), true);
9 }

```

GetTuple 根据 Row 找到 page，然后根据 page 读出数据

```

1 bool TableHeap::GetTuple(Row *row, Txn *txn) {
2     page_id_t page_id = row->GetRowId().GetPageId();
3     auto page = reinterpret_cast<TablePage>
4         ↳ *(buffer_pool_manager_>FetchPage(page_id));
5     if (page == nullptr) {
6         return false;
7     }
8     bool result = page->GetTuple(row, schema_, txn, lock_manager_);
9     buffer_pool_manager_>UnpinPage(page_id, false);
10    return result;
11 }

```

UpdateTuple 这一部分要根据 page 的 UpdateTuple 的返回值来判读，我在原来 page 的函数基础上修改了返回值，如果 slot_num 大于 tuple_count，说明 slot_num 是无效的，返回 1；如果 tuple_size 是删除的，返回 2；如果没有足够的空间更新，返回 3；如果成功，返回 0。修改后的函数如下：

```

1 uint32_t TablePage::UpdateTuple(Row &new_row, Row *old_row, Schema *schema, Txn *txn,
2     ↳ LockManager *lock_manager,
3     LogManager *log_manager) {
4     ASSERT(old_row != nullptr && old_row->GetRowId().Get() != INVALID_ROWID.Get(),
5         ↳ "invalid old row.");
6     uint32_t serialized_size = new_row.GetSerializedSize(schema);
7     ASSERT(serialized_size > 0, "Can not have empty row.");
8     uint32_t slot_num = old_row->GetRowId().GetSlotNum();
9     // If the slot number is invalid, abort.
10    if (slot_num >= GetTupleCount()) {
11        return 1;
12    }

```

```

11  uint32_t tuple_size = GetTupleSize(slot_num);
12  // If the tuple is deleted, abort.
13  if (IsDeleted(tuple_size)) {
14      return 2;
15  }
16  // If there is not enough space to update, we need to update via delete followed by
17  ↪ an insert (not enough space).
18  if (GetFreeSpaceRemaining() + tuple_size < serialized_size) {
19      return 3;
20  }
21  ...
22  // if success
23  return 0;

```

在上面的三种失败情况后，如果是情况 3，我们需要先删除再插入，下面是具体的实现：

```

1  bool TableHeap::UpdateTuple(Row &row, const RowId &rid, Txn *txn) {
2      // Find the page which contains the tuple.
3      auto page = reinterpret_cast<TablePage>
4      ↪ *(buffer_pool_manager_>FetchPage(rid.GetPageId()));
5      if (page == nullptr) {
6          return false;
7      }
8      Row old_row(rid);
9      // Update the tuple in the page.
10     int update_res = page->UpdateTuple(row, &old_row, schema_, txn, lock_manager_,
11     ↪ log_manager_);
12     if(update_res != 0){
13         if(update_res == 3){
14             if(InsertTuple(row, txn)){
15                 MarkDelete(rid, txn);
16                 buffer_pool_manager_>UnpinPage(page->GetTablePageId(), true);
17                 return true;
18             }else{
19                 buffer_pool_manager_>UnpinPage(page->GetTablePageId(), false);
20                 return false;
21             }
22         }else{
23             buffer_pool_manager_>UnpinPage(page->GetTablePageId(), false);
24             return false;
25         }
26     }
27 }

```

```

25     else{
26         row.SetRowId(rid);
27         buffer_pool_manager_>UnpinPage(page->GetTablePageId(), true);
28         return true;
29     }
30 }

```

2.3.2 table heap iterator implementation

实现迭代器，用于遍历 table heap 中的所有记录。在上层中提供 begin 和 end 函数，返回迭代器，然后通过迭代器的 next 函数来遍历所有记录。迭代器的成员变量是 TableHeap *table_heap_; RowId rid_;，其中 table_heap 是指向 table heap 的指针，rid 是当前记录的 RowId。下面给出迭代器的实现：

TableIterator ctor 构造函数，初始化迭代器的成员变量。

```

1 TableIterator::TableIterator(TableHeap *table_heap, RowId rid, Txn *txn) {
2     table_heap_ = table_heap;
3     rid_ = rid;
4     txn_ = txn;
5 }

```

TableIterator copy ctor 拷贝构造函数，用于拷贝迭代器。

```

1 TableIterator::TableIterator(const TableIterator &other) {
2     table_heap_ = other.table_heap_;
3     rid_ = other.rid_;
4     txn_ = other.txn_;
5 }

```

TableIterator dtor 析构函数，用于释放迭代器的资源。

```

1 TableIterator::~TableIterator() {
2     table_heap_ = nullptr;
3     rid_ = RowId();
4     txn_ = nullptr;
5 }

```


operator== == 运算符，用于判断两个迭代器是否相等。

```

1  bool TableIterator::operator==(const TableIterator &itr) const {
2      return this->rid_ == itr.rid_;
3  }

```

operator!= != 运算符，用于判断两个迭代器是否不相等。

```

1  bool TableIterator::operator!=(const TableIterator &itr) const {
2      return !(*this == itr);
3  }

```

operator* * 运算符，用于获取迭代器指向的记录。根据 RowId 找到 page，然后根据 page 读出数据，返回 Row。

```

1  const Row &TableIterator::operator*() {
2      auto page = reinterpret_cast<TablePage>
3      ↪ *(table_heap_>buffer_pool_manager_>FetchPage(rid_.GetPageId()));
4      if (page == nullptr) {
5          table_heap_>buffer_pool_manager_>UnpinPage(rid_.GetPageId(), false);
6          Row *row = new Row(INVALID_ROWID);
7          return *row;
8      }
9      Row *row = new Row(rid_);
10     bool res = page->GetTuple(row, table_heap_>schema_, txn_,
11     ↪ table_heap_>lock_manager_);
12     if(res){
13         table_heap_>buffer_pool_manager_>UnpinPage(page->GetPageId(), false);
14         return *row;
15     }else{
16         table_heap_>buffer_pool_manager_>UnpinPage(page->GetPageId(), false);
17         Row *row = new Row(INVALID_ROWID);
18         return *row;
19     }
20 }

```

operator-> -> 运算符，用于获取迭代器指向的记录，返回 Row 的指针。

```

1 Row *TableIterator::operator->() {
2     auto page = reinterpret_cast<TablePage>
        ↳ *(table_heap_>buffer_pool_manager_>FetchPage(rid_.GetPageId()));
3     if (page == nullptr) {
4         table_heap_>buffer_pool_manager_>UnpinPage(rid_.GetPageId(), false);
5         return nullptr;
6     }
7     Row *row = new Row(rid_);
8     bool res = page->GetTuple(row, table_heap_>schema_, txn_,
        ↳ table_heap_>lock_manager_);
9     if(res){
10         table_heap_>buffer_pool_manager_>UnpinPage(page->GetPageId(), false);
11         return row;
12     }else{
13         table_heap_>buffer_pool_manager_>UnpinPage(page->GetPageId(), false);
14         return nullptr;
15     }
16 }

```

operator= = 运算符，用于赋值和初始化。

```

1 TableIterator &TableIterator::operator=(const TableIterator &itr) noexcept {
2     table_heap_ = itr.table_heap_;
3     rid_ = itr.rid_;
4     txn_ = itr.txn_;
5     return *this;
6 }

```

operator++/+/++iter ++ 运算符，用于迭代器自增，从 tableheap 中找到下一个记录，需要注意的是可能存在下一条记录在下一页的情况。

```

1 TableIterator &TableIterator::operator++() {
2     auto page = reinterpret_cast<TablePage>
        ↳ *(table_heap_>buffer_pool_manager_>FetchPage(rid_.GetPageId()));
3     if (page == nullptr) {
4         table_heap_>buffer_pool_manager_>UnpinPage(rid_.GetPageId(), false);
5         this->table_heap_ = nullptr;
6         this->rid_ = RowId();
7         this->txn_ = nullptr;
8         return *this;
9     }

```

```

10  bool res = page->GetNextTupleRid(rid_, &rid_);
11  if(res){
12      // 不用翻下一页
13      table_heap_->buffer_pool_manager_->UnpinPage(page->GetPageId(), false);
14      return *this;
15  }else{
16      // 翻下一页
17      page_id_t next_page_id = page->GetNextPageId();
18      if(next_page_id == INVALID_PAGE_ID){
19          this->table_heap_ = nullptr;
20          this->rid_ = RowId();
21          this->txn_ = nullptr;
22          return *this;
23      }else{
24          auto next_page = reinterpret_cast<TablePage>
25              ↳ *(table_heap_->buffer_pool_manager_->FetchPage(next_page_id));
26          table_heap_->buffer_pool_manager_->UnpinPage(page->GetPageId(), false);
27          page = next_page;
28          page->GetFirstTupleRid(&rid_);
29      }
30      table_heap_->buffer_pool_manager_->UnpinPage(page->GetPageId(), false);
31      return *this;
32  }

```

operator++(int)//**iter++** ++ 运算符，用于迭代器自增，从 tableheap 中找到下一个记录，需要注意的是可能存在下一条记录在下一页的情况。

```

1  TableIterator TableIterator::operator++(int) {
2      TableIterator p(*this);
3      ++(*this);
4      return TableIterator{p};
5  }

```

2.3.3 table heap iterator

这是迭代器给上层的接口，包括 begin 和 end 函数，用于返回迭代器。

begin 返回一个指向第一条记录的迭代器，使用 page 实现的 GetFirstTuple 函数。如果没找到，返回一个空迭代器。

```

1 TableIterator TableHeap::Begin(Txn *txn) {
2     auto page = reinterpret_cast<TablePage>
3     ↪ *(>(buffer_pool_manager_>FetchPage(GetFirstPageId()));
4     if (page == nullptr) {
5         buffer_pool_manager_>UnpinPage(page->GetTablePageId(), false);
6         return TableIterator(this, RowId(), nullptr);
7     }
8     RowId rid;
9     if(page->GetFirstTupleRid(&rid)){
10         buffer_pool_manager_>UnpinPage(page->GetTablePageId(), false);
11         return TableIterator(this, rid, txn);
12     }else{
13         buffer_pool_manager_>UnpinPage(page->GetTablePageId(), false);
14         return TableIterator(this, RowId(), nullptr);
15     }
16 }

```

end 返回 end 迭代器，应当是最后一个的下一个，即返回一个空迭代器。

```

1 TableIterator TableHeap::End() { return TableIterator(nullptr, RowId(), nullptr); }

```

2.4 test

这部分框架提供的测试用例太少了，只对 insert tuple 和 get tuple 进行了测试，为了全面地测试我们实现的 table heap 和 table heap iterator，我们自己编写了测试用例，测试了 iterator，update tuple and free 测试用例如下：

这部分代码是自己添加的测试代码，完善了框架的测试，希望能够加入到下一年的测试框架中。

```

1     static string db_file_name = "table_heap_test.db";
2     using Fields = std::vector<Field>;
3
4     TEST(TableHeapTest, TableHeapSampleTest) {
5         // init testing instance
6         vector<RowId> row_ids;
7         remove(db_file_name.c_str());
8         auto disk_mgr_ = new DiskManager(db_file_name);

```

```

9  auto bpm_ = new BufferPoolManager(DEFAULT_BUFFER_POOL_SIZE, disk_mgr_);
10 const int row_nums = 10000;
11 // create schema
12 std::vector<Column *> columns = {new Column("id", TypeId::kTypeInt, 0, false,
13     ↪ false),
14     new Column("name", TypeId::kTypeChar, 64, 1, true,
15     ↪ false),
16     new Column("account", TypeId::kTypeFloat, 2, true,
17     ↪ false)};
18
19 auto schema = std::make_shared<Schema>(columns);
20 // create rows
21 std::unordered_map<int64_t, Fields *> row_values;
22 uint32_t size = 0;
23 TableHeap *table_heap = TableHeap::Create(bpm_, schema.get(), nullptr, nullptr,
24     ↪ nullptr);
25
26 for (int i = 0; i < row_nums; i++) {
27     int32_t len = RandomUtils::RandomInt(0, 64);
28     char *characters = new char[len];
29     RandomUtils::RandomString(characters, len);
30     Fields *fields =
31         new Fields{Field(TypeId::kTypeInt, i), Field(TypeId::kTypeChar,
32     ↪ const_cast<char *>(characters), len, true),
33             Field(TypeId::kTypeFloat, RandomUtils::RandomFloat(-999.f,
34     ↪ 999.f))};
35
36     Row row(*fields);
37     ASSERT_TRUE(table_heap->InsertTuple(row, nullptr));
38     if (row_values.find(row.GetRowId().Get()) != row_values.end()) {
39         std::cout << row.GetRowId().Get() << std::endl;
40         ASSERT_TRUE(false);
41     } else {
42         row_values.emplace(row.GetRowId().Get(), fields);
43         row_ids.push_back(row.GetRowId());
44         size++;
45     }
46     delete[] characters;
47 }
48
49 ASSERT_EQ(row_nums, row_values.size());
50 ASSERT_EQ(row_nums, size);
51 for (auto row_kv : row_values) {
52     size--;
53     Row row(RowId(row_kv.first));
54     table_heap->GetTuple(&row, nullptr);
55     ASSERT_EQ(schema.get()->GetColumnCount(), row.GetFields().size());
56     for (size_t j = 0; j < schema.get()->GetColumnCount(); j++) {

```

```

47     ASSERT_EQ(CmpBool::kTrue,
48         ↪ row.GetField(j)->CompareEquals(row_kv.second->at(j)));
49     }
50     // free spaces
51     delete row_kv.second;
52 }
53 // test iterator
54 ASSERT_EQ(size, 0);
55 int count = 0;
56 for(auto iter = table_heap->Begin(nullptr); iter != table_heap->End(); ++iter){
57     Row row(iter->GetRowId());
58     table_heap->GetTuple(&row, nullptr);
59     ASSERT_EQ(schema.get()->GetColumnCount(), row.GetFields().size());
60     ASSERT_EQ(schema.get()->GetColumnCount(), iter->GetFields().size());
61     for(size_t j = 0; j < schema.get()->GetColumnCount(); j++){
62         ASSERT_EQ(CmpBool::kTrue, row.GetField(j)->CompareEquals(*iter->GetField(j)));
63     }
64     count++;
65 }
66 ASSERT_EQ(count, row_nums);
67
68 // test update
69 std::unordered_map<int64_t, Fields *> new_row_values;
70 std::set<page_id_t> used_pages;
71 for(int i = 0; i < row_nums; i++){
72     int32_t len = RandomUtils::RandomInt(0, 64);
73     char *characters = new char[len];
74     RandomUtils::RandomString(characters, len);
75     Fields *fields = new Fields{Field(TypeId::kTypeInt, i), Field(TypeId::kTypeChar,
76         ↪ const_cast<char *>(characters), len, true),
77         Field(TypeId::kTypeFloat,
78         ↪ RandomUtils::RandomFloat(-999.f, 999.f))};
79     Row row(*fields);
80     ASSERT_EQ(true, table_heap->UpdateTuple(row, row_ids.at(i), nullptr));
81     ASSERT_EQ(false, row.GetRowId().GetPageId() == INVALID_PAGE_ID);
82     new_row_values[row.GetRowId().Get()] = fields;
83     delete[] characters;
84 }
85 // check the update result
86 for(auto row_kv2 : new_row_values){
87     Row row(RowId(row_kv2.first));
88     ASSERT_EQ(true, table_heap->GetTuple(&row, nullptr));
89     used_pages.insert(row.GetRowId().GetPageId());
90     ASSERT_EQ(schema.get()->GetColumnCount(), row.GetFields().size());

```

```

88     for(size_t j = 0; j < schema.get()->GetColumnCount(); j++){
89         ASSERT_EQ(CmpBool::kTrue,
90             ↪ row.GetField(j)->CompareEquals(row_kv2.second->at(j)));
91     }
92     table_heap->ApplyDelete(row.GetRowId(), nullptr);
93     ASSERT_EQ(false, table_heap->GetTuple(&row, nullptr));
94     delete row_kv2.second;
95 }
96 // test free
97 table_heap->FreeTableHeap();
98 auto *disk = new DiskManager(db_file_name);
99 auto *bpm = new BufferPoolManager(100, disk);
100 ASSERT_EQ(true, bpm->IsPageFree(*used_pages.begin()));
101 for(auto page_id : used_pages){
102     ASSERT_EQ(true, bpm->IsPageFree(page_id));
103 }
104 delete disk;
105 delete bpm;
106 delete disk_mgr_;
107 delete bpm_;
108 }

```

本节的 test 全部通过。

3 INDEX MANAGER

3.1 summary

Index Manager 负责数据表索引的实现和管理，包括：索引的创建和删除，索引键的等值查找，索引键的范围查找（返回对应的迭代器），以及插入和删除键值等操作，并对外提供相应的接口。为了提高搜索的效率，我们实现基于磁盘的 B+ 树动态索引结构。首先完成 b+ 树数据结构的实现，然后建立在此基础上建立索引。

3.2 BPlusTree

首先我们要实现 B+ 树的数据结构，在这里是基于磁盘的数据页实现的，因此我们需要实现 BPlusTreeInternal 和 BPlusTreeLeaf，为 BPlusTree 的实现提供接口。

3.2.1 BPlusTreePage

我们知道，b+ 树存在两种节点，一种是内部节点，一种是叶子节点。内部节点存储的是索引键和指向子节点的指针，叶子节点存储的是索引键和指向记录的指针。这二

```

L> ./test/tuple_test
[=====] Running 2 tests from 1 test suite.
[-----] Global test environment set-up.
[-----] 2 tests from TupleTest
[ RUN     ] TupleTest.FieldSerializeDeserializeTest
[         OK ] TupleTest.FieldSerializeDeserializeTest (0 ms)
[ RUN     ] TupleTest.RowTest
[         OK ] TupleTest.RowTest (0 ms)
[-----] 2 tests from TupleTest (0 ms total)

[-----] Global test environment tear-down
[=====] 2 tests from 1 test suite ran. (0 ms total)
[ PASSED ] 2 tests.

```

图 9: Tuple Test

```

L> ./test/table_heap_test
[=====] Running 1 test from 1 test suite.
[-----] Global test environment set-up.
[-----] 1 test from TableHeapTest
[ RUN     ] TableHeapTest.TableHeapSampleTest
[         OK ] TableHeapTest.TableHeapSampleTest (1399 ms)
[-----] 1 test from TableHeapTest (1399 ms total)

[-----] Global test environment tear-down
[=====] 1 test from 1 test suite ran. (1399 ms total)
[ PASSED ] 1 test.

```

图 10: Table Heap Test

者具有部分相同的成员变量，因此设计上可以通过一个共同父类继承得到两个子类，这里的BPlusTreepage是BPlusTreeInternalPage和BPlusTreeLeafPage的父类。包含了中间结点和叶子结点的共同成员变量。

- **page_type_**: 用于标识当前页是叶子页还是中间页。
- **key_size_**: 索引键的大小。
- **lsn_**: 数据页的日志序列号，该模块中不会用到；
- **size_**: 当前结点中存储 Key-Value 键值对的数量；
- **max_size_**: 当前结点中最多可以存储的 Key-Value 键值对的数量；
- **parent_page_id_**: 父结点对应数据页的 page_id；
- **page_id_**: 当前数据页的 page_id；

这个类的实现较为简单，主要是一些 getter 和 setter 函数，需要注意的点就是最小 size 应当是最大 size 的一半，这部分实现较为简单，这里就不再赘述。

3.2.2 BPlusTreeInternalPage

内部节点的实现，继承自 BPlusTreePage，内部节点存储的是索引键和指向子节点的指针，因此需要实现插入、删除、查找等操作，除此以外，BPlusTree 还会做 merge 和 split 操作，需要一些将子节点的 key 和 value 移动到另一个节点的操作。以下是各个接口的具体实现：

Init 初始化 page_id, parent_page_id, page_type, key_size, max_size, size, lsn 等成员变量。

```

1 void InternalPage::Init(page_id_t page_id, page_id_t parent_id, int key_size, int
  ↳ max_size) {
2     SetPageType(IndexPageType::INTERNAL_PAGE);
3     SetSize(0);
4     SetPageId(page_id);
5     SetParentPageId(parent_id);
6     SetKeySize(key_size);
7     SetMaxSize(max_size);
8     SetLSN(INVALID_LSN);
9 }

```

Search

Lookup 这里要使用二分查找来提高搜索的速度，这里调用了 KeyManager 的 CompareKeys 函数来比较两个 key 的大小，我们应当是无感知地获得比较的结果。

```

1 page_id_t InternalPage::Lookup(const GenericKey *key, const KeyManager &KM) {
2     int left = 1, right = GetSize() - 1;
3     while (left <= right)
4     {
5         int mid = (left + right) / 2;
6         if (KM.CompareKeys(KeyAt(mid), key) <= 0)
7             left = mid + 1;
8         else
9             right = mid - 1;
10    }
11    return ValueAt(right);
12 }

```

Insert

实现插入操作需要使用的几个函数。

PopulateNewRoot 给建立的新的根节点进行初始化，将两个子节点的 page_id 和 key_id 插入到新的根节点中。

```

1 void InternalPage::PopulateNewRoot(const page_id_t &old_value, GenericKey *new_key,
  ↳ const page_id_t &new_value) {
2     SetSize(2);
3     SetValueAt(0, old_value);
4     SetKeyAt(1, new_key);
5     SetValueAt(1, new_value);
6 }

```

InsertNodeAfter 在给定的位置后插入一个新的 key 和 value，这里需要将原来的 key 和 value 向后移动一个位置，然后插入新的 key 和 value，返回当前的 size。

```

1 int InternalPage::InsertNodeAfter(const page_id_t &old_value, GenericKey *new_key,
  ↳ const page_id_t &new_value) {
2     int index = ValueIndex(old_value);
3     if (index == -1)
4         return -1;
5     int size = GetSize();
6     for (int i = size; i > index + 1; --i)
7     {
8         PairCopy(pairs_off + i * pair_size, pairs_off + (i - 1) * pair_size);
9     }
10    SetKeyAt(index + 1, new_key);
11    SetValueAt(index + 1, new_value);
12    IncreaseSize(1);
13    return size + 1;
14 }

```

split

MoveHalfTo 将当前节点的一半 key 和 value 移动到另一个节点，这个函数调用了 CopyNfrom 这个复制函数，将当前节点的一半 key 和 value 复制到另一个节点中，然后将当前节点的一半 key 和 value 删除。

```

1 void InternalPage::MoveHalfTo(InternalPage *recipient, BufferPoolManager
  ↳ *buffer_pool_manager) {
2     int size = GetSize();
3     int left_size = size / 2, right_size = size - left_size;
4     recipient->CopyNFrom(pairs_off + left_size * pair_size, right_size,
  ↳ buffer_pool_manager);

```

```

5     IncreaseSize(-right_size);
6 }

```

CopyNFrom 向自身拷贝另一个节点的一部分 key 和 value，这里调用了 PairCopy 函数，将另一个节点的 key 和 value 复制到自身中。此时是中间结点还要对儿子节点的 parent 进行更新。

```

1 void InternalPage::CopyNFrom(void *src, int size, BufferPoolManager
  ↪ *buffer_pool_manager) {
2     // copy entries
3     int old_size = GetSize();
4     PairCopy(pairs_off + old_size * pair_size, src, size);
5     IncreaseSize(size);
6     // update parent page id and get persisted
7     for(int i = old_size; i < size; ++i)
8     {
9         page_id_t page_id = ValueAt(i);
10        auto *page = buffer_pool_manager->FetchPage(page_id);
11        ASSERT(page != nullptr, "page is nullptr");
12        auto *internal_page = reinterpret_cast<BPlusTreePage *>(page->GetData());
13        internal_page->SetParentPageId(GetPageId());
14        buffer_pool_manager->UnpinPage(page_id, true);
15    }
16 }

```

Remove

Remove 在中间结点中删除特定的 key 和 value。

```

1 void InternalPage::Remove(int index) {
2     int size = GetSize();
3     for (int i = index; i < size - 1; ++i)
4     {
5         PairCopy(pairs_off + i * pair_size, pairs_off + (i + 1) * pair_size);
6     }
7     IncreaseSize(-1);
8 }

```

RemoveAndReturnOnlyChild 删除当前节点的 key 和 value，返回当前节点的唯一子节点的 page_id。这个函数只会在 AdjustRoot() 中调用。

```

1 page_id_t InternalPage::RemoveAndReturnOnlyChild() {
2     page_id_t value = ValueAt(0);
3     IncreaseSize(-1);
4     return value;
5 }

```

merge

MoveAllTo 将自身的全部结点移动到另一个接收结点中，需要更新子节点的 `parent_page_id`。

```

1 void InternalPage::MoveAllTo(InternalPage *recipient, GenericKey *middle_key,
   ↪ BufferPoolManager *buffer_pool_manager) {
2     int size = GetSize(), recipient_size = recipient->GetSize();
3     recipient->CopyNFrom(pairs_off, size, buffer_pool_manager);
4     recipient->SetKeyAt(recipient_size, middle_key);
5     // update parent page id and get persisted
6     for(int i = 0; i < size; ++i)
7     {
8         page_id_t page_id = ValueAt(i);
9         auto *page = buffer_pool_manager->FetchPage(page_id);
10        ASSERT(page != nullptr, "page is nullptr");
11        auto *internal_page = reinterpret_cast<BPlusTreePage *>(page->GetData());
12        internal_page->SetParentPageId(recipient->GetPageId());
13        buffer_pool_manager->UnpinPage(page_id, true);
14    }
15    SetSize(0);
16 }

```

Redistribute

MoveFirstToEndOf 将另一个节点的第一个 key 和 value 移动到自身的最后一个位置。

```

1 void InternalPage::MoveFirstToEndOf(InternalPage *recipient, GenericKey *middle_key,
   ↪ BufferPoolManager *buffer_pool_manager) {
2     recipient->CopyLastFrom(middle_key, ValueAt(0), buffer_pool_manager);
3     Remove(0);
4 }

```

CopyLastFrom 将另一个节点的最后一个 key 和 value 复制到自身的最后一个位置。

```

1 void InternalPage::CopyLastFrom(GenericKey *key, const page_id_t value,
  ↳ BufferPoolManager *buffer_pool_manager) {
2     int size = GetSize();
3     SetKeyAt(size, key);
4     SetValueAt(size, value);
5     IncreaseSize(1);
6     auto *page = buffer_pool_manager->FetchPage(value);
7     ASSERT(page != nullptr, "page is nullptr");
8     auto *internal_page = reinterpret_cast<BPlusTreePage *>(page->GetData());
9     internal_page->SetParentPageId(GetPageId());
10    buffer_pool_manager->UnpinPage(value, true);
11 }

```

MoveLastToFrontOf 将另一个节点的最后一个 key 和 value 移动到自身的第一个位置。

```

1 void InternalPage::MoveLastToFrontOf(InternalPage *recipient, GenericKey
  ↳ *middle_key, BufferPoolManager *buffer_pool_manager) {
2     recipient->CopyFirstFrom(ValueAt(GetSize() - 1), buffer_pool_manager);
3     recipient->SetKeyAt(0, KeyAt(GetSize() - 1));
4     recipient->SetKeyAt(1, middle_key);
5     Remove(GetSize() - 1);
6 }

```

CopyFirstFrom 将另一个节点的第一个 key 和 value 复制到自身的第一个位置，修改儿子节点的 parent_page_id。

```

1 void InternalPage::CopyFirstFrom(const page_id_t value, BufferPoolManager
  ↳ *buffer_pool_manager) {
2     int size = GetSize();
3     for(int i = size; i > 0; i--){
4         PairCopy(pairs_off + i * pair_size, pairs_off + (i - 1) * pair_size);
5     }
6     SetValueAt(0, value);
7     IncreaseSize(1);
8     auto *page = buffer_pool_manager->FetchPage(value);
9     ASSERT(page != nullptr, "page is nullptr");
10    auto *internal_page = reinterpret_cast<BPlusTreePage *>(page->GetData());
11    internal_page->SetParentPageId(GetPageId());

```

```

12     buffer_pool_manager->UnpinPage(value, true);
13 }

```

3.2.3 BPlusTreeLeafPage

叶子节点的实现，继承自 BPlusTreePage，叶子节点存储的是索引键和指向记录的指针，大体的实现与内部节点类似，只是叶子节点不需要存储子节点的 page_id，而需要存储 next_page_id，用于在范围查找时返回下一个叶子节点的 page_id。以下是各个接口的具体实现：

Init 初始化 page_id, parent_page_id, page_type, key_size, max_size, size, lsn 等成员变量。

```

1 void LeafPage::Init(page_id_t page_id, page_id_t parent_id, int key_size, int
   ↪ max_size) {
2     SetPageType(IndexPageType::LEAF_PAGE);
3     SetKeySize(key_size);
4     SetLSN(INVALID_LSN);
5     SetSize(0);
6     SetPageId(page_id);
7     SetParentPageId(parent_id);
8     SetMaxSize(max_size);
9     SetNextPageId(INVALID_PAGE_ID);
10 }

```

Search

KeyIndex 在叶子节点中查找特定的 key，返回对应的 value。

```

1 int LeafPage::KeyIndex(const GenericKey *key, const KeyManager &KM) {
2     int left = 0, right = GetSize() - 1;
3     while (left <= right)
4     {
5         int mid = (left + right) / 2;
6         if (KM.CompareKeys(KeyAt(mid), key) >= 0)
7             right = mid - 1;
8         else
9             left = mid + 1;
10    }
11    return left;
12 }

```

Lookup 对给定的 key，查找是否存在。

```

1 bool LeafPage::Lookup(const GenericKey *key, RowId &value, const KeyManager &KM) {
2     int index = KeyIndex(key, KM);
3     if (index < GetSize() && KM.CompareKeys(KeyAt(index), key) == 0)
4     {
5         value = ValueAt(index);
6         return true;
7     }
8     return false;
9 }

```

Insert 在叶子节点中插入特定的 key 和 value

```

1 int LeafPage::Insert(GenericKey *key, const RowId &value, const KeyManager &KM) {
2     int index = KeyIndex(key, KM);
3     if (index < GetSize() && KM.CompareKeys(KeyAt(index), key) == 0) {
4         return GetSize();
5     }
6     int size = GetSize();
7     for (int i = size; i > index; i--)
8     {
9         PairCopy(pairs_off + i * pair_size, pairs_off + (i - 1) * pair_size);
10    }
11    SetKeyAt(index, key);
12    SetValueAt(index, value);
13    SetSize(size + 1);
14    return size + 1;
15 }

```

split

MoveHalfTo 将当前节点的一半 key 和 value 移动到另一个节点，这个函数调用了 CopyNfrom 这个复制函数，将当前节点的一半 key 和 value 复制到另一个节点中，然后将当前节点的一半 key 和 value 删除。

```

1 void LeafPage::MoveHalfTo(LeafPage *recipient) {
2     int size = GetSize();
3     int half = size / 2;
4     recipient->CopyNFrom(pairs_off + half * pair_size, size - half);

```

```

5   IncreaseSize(-size + half);
6 }

```

CopyNFrom 向自身拷贝另一个节点的一部分 key 和 value，这里调用了 PairCopy 函数，将另一个节点的 key 和 value 复制到自身中。

```

1 void LeafPage::CopyNFrom(void *src, int size) {
2     PairCopy(pairs_off + GetSize() * pair_size, src, size);
3     IncreaseSize(size);
4 }

```

Remove

RemoveAndDeleteRecord 在叶子节点中删除特定的 key 和 value。

```

1 int LeafPage::RemoveAndDeleteRecord(const GenericKey *key, const KeyManager &KM) {
2     int index = KeyIndex(key, KM);
3     if (index < GetSize() && KM.CompareKeys(KeyAt(index), key) == 0)
4     {
5         int size = GetSize();
6         for (int i = index; i < size - 1; i++)
7         {
8             PairCopy(pairs_off + i * pair_size, pairs_off + (i + 1) * pair_size);
9         }
10        IncreaseSize(-1);
11    }
12    return GetSize();
13 }

```

merge

MoveAllTo 将自身的全部结点移动到另一个接收结点中，需要更新 next_page_id。

```

1 void LeafPage::MoveAllTo(LeafPage *recipient) {
2     recipient->CopyNFrom(pairs_off, GetSize());
3     recipient->SetNextPageId(GetNextPageId());
4     SetSize(0);
5 }

```


Redistribute

MoveFirstToEndOf 将另一个节点的第一个 key 和 value 移动到自身的最后一个位置。

```

1 void LeafPage::MoveFirstToEndOf(LeafPage *recipient) {
2     int size = GetSize();
3     recipient->CopyLastFrom(KeyAt(0), ValueAt(0));
4     for (int i = 0; i < size - 1; i++)
5     {
6         PairCopy(pairs_off + i * pair_size, pairs_off + (i + 1) * pair_size);
7     }
8     IncreaseSize(-1);
9 }

```

CopyLastFrom 将另一个节点的最后一个 key 和 value 复制到自身的最后一个位置。

```

1 void LeafPage::CopyLastFrom(GenericKey *key, const RowId value) {
2     int size = GetSize();
3     SetKeyAt(size, key);
4     SetValueAt(size, value);
5     IncreaseSize(1);
6 }

```

MoveLastToFrontOf 将另一个节点的最后一个 key 和 value 移动到自身的第一个位置。

```

1 void LeafPage::MoveLastToFrontOf(LeafPage *recipient) {
2     int size = GetSize();
3     recipient->CopyFirstFrom(KeyAt(size - 1), ValueAt(size - 1));
4     IncreaseSize(-1);
5 }

```

CopyFirstFrom 将另一个节点的第一个 key 和 value 复制到自身的第一个位置。

```

1 void LeafPage::CopyFirstFrom(GenericKey *key, const RowId value) {
2     int size = GetSize();
3     for (int i = size; i > 0; i--)
4     {
5         PairCopy(pairs_off + i * pair_size, pairs_off + (i - 1) * pair_size);

```

```

6   }
7   SetKeyAt(0, key);
8   SetValueAt(0, value);
9   IncreaseSize(1);
10  }

```

3.2.4 BPlusTree implementation

在前面的部分，我们已经完成了叶子结点和中间结点的实现，接下来我们要实现 BPlusTree 的接口，完成 B+ 树的创建、插入、删除、查找和释放等操作。这里 BPlusTree 的成员变量包括：

- `index_id_t index_id_;`
- `page_id_t root_page_id_INVALID_PAGE_ID;`
- `BufferPoolManager *buffer_pool_manager_;`
- `KeyManager processor_;`
- `int leaf_max_size_;` 即 $((\text{PAGE_SIZE} - \text{LEAF_PAGE_HEADER_SIZE}) / (\text{processor_}.\text{GetKeySize}() + \text{sizeof}(\text{RowId})) - 1)$
- `int internal_max_size_;` 即 $((\text{PAGE_SIZE} - \text{INTERNAL_PAGE_HEADER_SIZE}) / (\text{processor_}.\text{GetKeySize}() + \text{sizeof}(\text{page_id_t})) - 1)$

Init 首先创建一个 BPlusTree 的实例，初始化一些成员变量，包括 `key_size`, `internal_page_size`, `leaf_page_size`, `root_page_id`, `buffer_pool_manager` 等。除此以外，我们还需要修改元数据页即 `IndexRootsPage`，将 `root_page_id` 写入到元数据页中，调用了 `UpdateRootPageId` 函数。

```

1  BPlusTree::BPlusTree(index_id_t index_id, BufferPoolManager *buffer_pool_manager,
   ↪   const KeyManager &KM,
2      int leaf_max_size, int internal_max_size)
3      : index_id_(index_id),
4        buffer_pool_manager_(buffer_pool_manager),
5        processor_(KM),
6        leaf_max_size_(leaf_max_size),
7        internal_max_size_(internal_max_size) {
8      root_page_id_ = INVALID_PAGE_ID;
9      Page *page = buffer_pool_manager_>FetchPage(INDEX_ROOTS_PAGE_ID);
10     IndexRootsPage *index_roots_page = reinterpret_cast<IndexRootsPage>
   ↪     *(page->GetData());

```

```

11     if (!index_roots_page->GetRootId(index_id, &root_page_id_)) {
12         root_page_id_ = INVALID_PAGE_ID;
13         UpdateRootPageId(1);
14     }
15     buffer_pool_manager_->UnpinPage(INDEX_ROOTS_PAGE_ID, false);
16     leaf_max_size_ = LEAF_PAGE_SIZE;
17     internal_max_size_ = INTERNAL_PAGE_SIZE;
18 }

```

UpdateRootPageId 修改 index_roots_page 中的当前 index 的 root_page_id, 用于索引查找 b+ 树的根节点。

```

1 void BPlusTree::UpdateRootPageId(int insert_record) {
2     auto* page = buffer_pool_manager_->FetchPage(INDEX_ROOTS_PAGE_ID);
3     auto* index_roots_page = reinterpret_cast<IndexRootsPage*>(page->GetData());
4     if (insert_record) {
5         index_roots_page->Insert(index_id_, root_page_id_);
6     } else {
7         index_roots_page->Update(index_id_, root_page_id_);
8     }
9     buffer_pool_manager_->UnpinPage(INDEX_ROOTS_PAGE_ID, true);
10 }

```

destroy

我认为这里的 destroy 设计默认值为 INVALID_PAGE_ID 时实际上要求是删除全树这件事情十分不合理，这样的表意不明确，在文档里面也没有提及，造成了许多的困扰。

销毁以传入的结点为根结点的 b+ 树，这里的销毁是递归的，首先销毁子节点，然后销毁自身，最后释放自身的 page。如果是默认情况就删除全树，还要更新元数据页中的 root_page_id。

```

1 void BPlusTree::Destroy(page_id_t current_page_id) {
2     if (current_page_id == INVALID_PAGE_ID) {
3         current_page_id = root_page_id_;
4     }
5     Page* page = buffer_pool_manager_->FetchPage(current_page_id);

```

```

6   BPlusTreePage* node = reinterpret_cast<BPlusTreePage*>(page->GetData());
7   if(node->IsLeafPage()) {
8       buffer_pool_manager->DeletePage(current_page_id);
9       buffer_pool_manager->UnpinPage(current_page_id, true);
10  }else{
11      InternalPage* internal = reinterpret_cast<InternalPage*>(node);
12      for(int i = 0; i < internal->GetSize(); i++) {
13          Destroy(internal->ValueAt(i));
14      }
15      buffer_pool_manager->DeletePage(current_page_id);
16      buffer_pool_manager->UnpinPage(current_page_id, true);
17  }
18  if(current_page_id == root_page_id_) {
19      auto head = buffer_pool_manager->FetchPage(INDEX_ROOTS_PAGE_ID);
20      IndexRootsPage* index_roots_page =
21      ↪ reinterpret_cast<IndexRootsPage*>(head->GetData());
22      index_roots_page->Delete(index_id_);
23      buffer_pool_manager->UnpinPage(INDEX_ROOTS_PAGE_ID, true);
24  }

```

IsEmpty 判断 b+ 树是否为空，如果根结点 page_id 为 INVALID_PAGE_ID，那么 b+ 树为空。

```

1  bool BPlusTree::IsEmpty() const {
2      return root_page_id_ == INVALID_PAGE_ID;
3  }

```

FindLeafPage 这是一个重要的函数，用于查找特定的 key 所在的叶子结点，在后面的插入、删除、查询中都会用到。

```

1  Page *BPlusTree::FindLeafPage(const GenericKey *key, page_id_t page_id, bool
2  ↪ leftMost) {
3      if(page_id == INVALID_PAGE_ID){
4          return FindLeafPage(key, root_page_id_, leftMost);
5      }
6      Page *page = buffer_pool_manager->FetchPage(page_id);
7      auto *node = reinterpret_cast<BPlusTreePage *>(page->GetData());
8      if(node->IsLeafPage()) {
9          return page;
10     }else{

```

```

10     auto *internal_node = reinterpret_cast<InternalPage *>(node);
11     page_id_t child_page_id = leftMost ? internal_node->ValueAt(0) :
        ↳ internal_node->Lookup(key, processor_);
12     buffer_pool_manager_->UnpinPage(page_id, false);
13     return FindLeafPage(key, child_page_id, leftMost);
14 }
15 return nullptr;
16 }

```

Search 查找特定的 key，返回对应的 value，这里调用了 FindLeafPage 函数，找到对应的叶子结点，然后调用叶子结点的 Lookup 函数，返回对应的 value。

```

1 bool BPlusTree::GetValue(const GenericKey *key, std::vector<RowId> &result, Txn
    ↳ *transaction) {
2     if(IsEmpty()) {
3         return false;
4     }
5     LeafPage *leaf_node = reinterpret_cast<LeafPage *>(FindLeafPage(key, root_page_id_,
        ↳ false)->GetData());
6     RowId value;
7     bool find = leaf_node->Lookup(key, value, processor_);
8     if(find) {
9         result.push_back(value);
10    }
11    buffer_pool_manager_->UnpinPage(leaf_node->GetPageId(), false);
12    return find;
13 }

```

Insertion

Insert 插入一个 key-value 键值对，如果是空树，那么直接插入到根结点中，否则调用插入到叶子结点的函数。返回值代表是否插入重复 key，我们是不允许插入重复 key 的。

```

1 bool BPlusTree::Insert(GenericKey *key, const RowId &value, Txn *transaction) {
2     if(IsEmpty()) {
3         StartNewTree(key, value);
4         return true;
5     }
6     return InsertIntoLeaf(key, value, transaction);
7 }

```

StartNewTree 如果是一个空树，那么直接插入到根结点中，这里需要创建一个新的叶子结点，然后将新的叶子结点的 `page_id` 写入到元数据页中。

```

1 void BPlusTree::StartNewTree(GenericKey *key, const RowId &value) {
2     Page *page = buffer_pool_manager->NewPage(root_page_id_);
3     ASSERT(page != nullptr, "out of memory");
4     LeafPage *leaf_page = reinterpret_cast<LeafPage *>(page->GetData());
5     leaf_page->Init(root_page_id_, INVALID_PAGE_ID, processor_.GetKeySize(),
6         ↪ leaf_max_size_);
7     leaf_page->Insert(key, value, processor_);
8     buffer_pool_manager->UnpinPage(root_page_id_, true);
9     UpdateRootPageId(0);
10 }

```

InsertIntoLeaf 插入一个 key-value 键值对到叶子结点中，如果叶子结点的 size 超过了最大值，那么需要进行 split 操作，将一半的 key 和 value 移动到另一个新的叶子结点中，然后将新的 key 和 value 插入到叶子结点中，然后将新的叶子插到 parent 中。

```

1 bool BPlusTree::InsertIntoLeaf(GenericKey *key, const RowId &value, Txn *transaction)
2 ↪ {
3     static int count = 0;
4     Page *leaf_page = FindLeafPage(key, root_page_id_, false);
5     LeafPage *leaf_node = reinterpret_cast<LeafPage *>(leaf_page->GetData());
6     RowId receive_value;
7     if(leaf_node->Lookup(key, receive_value, processor_)) {
8         cout << "Duplicate key" << endl;
9         buffer_pool_manager->UnpinPage(leaf_node->GetPageId(), false);
10        return false;
11    }
12    if(leaf_node->GetSize() < leaf_node->GetMaxSize()) {
13        leaf_node->Insert(key, value, processor_);
14        buffer_pool_manager->UnpinPage(leaf_node->GetPageId(), true);
15    }else{
16        // deal with split
17        LeafPage *new_leaf_node = Split(leaf_node, transaction);
18        if(processor_.CompareKeys(key, new_leaf_node->KeyAt(0)) < 0) {
19            leaf_node->Insert(key, value, processor_);
20        }else{
21            new_leaf_node->Insert(key, value, processor_);
22        }
23        InsertIntoParent(leaf_node, new_leaf_node->KeyAt(0), new_leaf_node, transaction);
24        buffer_pool_manager->UnpinPage(leaf_node->GetPageId(), true);
25        buffer_pool_manager->UnpinPage(new_leaf_node->GetPageId(), true);

```

```

25     }
26     return true;
27 }

```

Split 这一部分包含叶子结点和中间结点的 split 操作，首先创建一个新的叶子结点或者中间结点，然后将当前结点的一半 key 和 value 移动到新的结点中，然后返回新的结点。叶子结点需要更新 next_page_id，中间结点需要更新子节点的 parent_page_id。这里的实现依赖于前面结点页的实现，

```

1 BPlusTreeInternalPage *BPlusTree::Split(InternalPage *node, Txn *transaction) {
2     page_id_t new_page_id;
3     Page *page = buffer_pool_manager->NewPage(new_page_id);
4     ASSERT(page != nullptr, "out of memory");
5     InternalPage *new_internal_page = reinterpret_cast<InternalPage
6     ↪    *>(page->GetData());
7     new_internal_page->Init(new_page_id, node->GetParentPageId(),
8     ↪    processor_.GetKeySize(), internal_max_size_);
9     node->MoveHalfTo(new_internal_page, buffer_pool_manager_);
10    return new_internal_page;
11 }
12
13 BPlusTreeLeafPage *BPlusTree::Split(LeafPage *node, Txn *transaction) {
14     page_id_t new_page_id;
15     Page *page = buffer_pool_manager->NewPage(new_page_id);
16     ASSERT(page != nullptr, "out of memory");
17     LeafPage *new_leaf_page = reinterpret_cast<LeafPage *>(page->GetData());
18     new_leaf_page->Init(new_page_id, node->GetParentPageId(), processor_.GetKeySize(),
19     ↪    leaf_max_size_);
20     node->MoveHalfTo(new_leaf_page);
21     new_leaf_page->SetNextPageId(node->GetNextPageId());
22     node->SetNextPageId(new_page_id);
23     return new_leaf_page;
24 }

```

InsertIntoParent 插入一个 key 和 value 到父节点中，如果此时的结点已经是根节点了，那么就要新建一个根节点作为父节点，然后将原来的根节点和新的结点插入到新的根节点中，然后更新根节点的 page_id。否则就可以根据父节点的 size 判断是否可以插入，如果不能插入就需要进行 split 操作，然后递归调用 InsertIntoParent 函数，直到插入成功。

```

1 void BPlusTree::InsertIntoParent(BPlusTreePage *old_node, GenericKey *key,
  ↳ BPlusTreePage *new_node, Txn *transaction) {
2     if(old_node->IsRootPage()) {
3         page_id_t new_page_id;
4         Page *page = buffer_pool_manager->NewPage(new_page_id);
5         ASSERT(page != nullptr, "out of memory");
6         InternalPage *new_root_node = reinterpret_cast<InternalPage *>(page->GetData());
7         new_root_node->Init(new_page_id, INVALID_PAGE_ID, processor_.GetKeySize(),
  ↳ internal_max_size_);
8         new_root_node->PopulateNewRoot(old_node->GetPageId(), key,
  ↳ new_node->GetPageId());
9         old_node->SetParentPageId(new_page_id);
10        new_node->SetParentPageId(new_page_id);
11        root_page_id_ = new_page_id;
12        UpdateRootPageId(0);
13        buffer_pool_manager->UnpinPage(new_page_id, true);
14    }else{
15        page_id_t page_id = old_node->GetParentPageId();
16        Page *page = buffer_pool_manager->FetchPage(page_id);
17        InternalPage *parent_node = reinterpret_cast<InternalPage *>(page->GetData());
18        if(parent_node->GetSize() < parent_node->GetMaxSize()) {
19            parent_node->InsertNodeAfter(old_node->GetPageId(), key,
  ↳ new_node->GetPageId());
20            new_node->SetParentPageId(page_id);
21            buffer_pool_manager->UnpinPage(page_id, true);
22        }else{
23            InternalPage *new_internal_node = Split(parent_node, transaction);
24            if(processor_.CompareKeys(key, new_internal_node->KeyAt(0)) < 0) {
25                new_node->SetParentPageId(page_id);
26                parent_node->InsertNodeAfter(old_node->GetPageId(), key,
  ↳ new_node->GetPageId());
27            }else{
28                new_node->SetParentPageId(new_internal_node->GetPageId());
29                new_internal_node->InsertNodeAfter(old_node->GetPageId(), key,
  ↳ new_node->GetPageId());
30            }
31            InsertIntoParent(parent_node, new_internal_node->KeyAt(0), new_internal_node,
  ↳ transaction);
32            buffer_pool_manager->UnpinPage(parent_node->GetPageId(), true);
33            buffer_pool_manager->UnpinPage(new_internal_node->GetPageId(), true);
34        }
35    }
36 }

```


Deletion

插入的过程中需要分裂，删除的过程中则需要合并和重分配。下面是各个接口的具体实现：

Remove 如果是空树的话，直接返回 false，否则找到 leaf 后调用 RemoveAndDeleteRecord 函数，如果删除后叶子结点的 size 小于最小值，那么需要进行合并或者重分配操作。

```

1 void BPlusTree::Remove(const GenericKey *key, Txn *transaction) {
2     if(IsEmpty()) {
3         return;
4     }
5     Page* page = FindLeafPage(key, root_page_id_, false);
6     LeafPage *leaf_node = reinterpret_cast<LeafPage *>(page->GetData());
7     int new_size = leaf_node->RemoveAndDeleteRecord(key, processor_);
8     if(new_size < leaf_node->GetMinSize()) {
9         if(CoalesceOrRedistribute(leaf_node, transaction)) {
10             buffer_pool_manager_->UnpinPage(leaf_node->GetPageId(), false);
11             buffer_pool_manager_->DeletePage(leaf_node->GetPageId());
12         }else{
13             buffer_pool_manager_->UnpinPage(leaf_node->GetPageId(), true);
14         }
15     }else{
16         buffer_pool_manager_->UnpinPage(leaf_node->GetPageId(), true);
17     }
18 }

```

CoalesceOrRedistribute 这里使用模板实现了合并和重分配操作，如果叶子结点的 size 小于最小值，那么需要进行合并或者重分配操作。合并的条件是兄弟节点与自身的 size 之和小于最大值，我们的实现逻辑是优先实现重分配，如果重分配失败，那么就进行合并操作，合并操作需要递归调用 CoalesceOrRedistribute 函数，直到合并成功。这里返回值的意义是要删除的叶子节点所在的 page 是否要删除，即是否发生了合并。

```

1 template <typename N>
2 bool BPlusTree::CoalesceOrRedistribute(N *&node, Txn *transaction) {
3     if (node->IsRootPage()) {
4         return AdjustRoot(node);
5     }
6     auto parent = reinterpret_cast<InternalPage
7     ↳ *>(buffer_pool_manager_->FetchPage(node->GetParentPageId())->GetData());
8     int index = parent->ValueIndex(node->GetPageId());

```

```

8   page_id_t sibilings = index == 0 ? parent->ValueAt(1) : parent->ValueAt(index - 1);
9   auto sibilings_page = reinterpret_cast<N
    ↳ *>(buffer_pool_manager_->FetchPage(sibilings)->GetData());
10  if(node->GetSize() + sibilings_page->GetSize() > node->GetMaxSize()) {
11      Redistribute(sibilings_page, node, index);
12      buffer_pool_manager_->UnpinPage(sibilings_page->GetPageId(), true);
13      buffer_pool_manager_->UnpinPage(node->GetPageId(), true);
14      return false;
15  }else{
16      if(Coalesce(sibilings_page, node, parent, index, transaction)) {
17          if(parent->IsRootPage()) {
18              if(AdjustRoot(parent)) {
19                  buffer_pool_manager_->UnpinPage(parent->GetPageId(), false);
20                  buffer_pool_manager_->DeletePage(parent->GetPageId());
21              } else{
22                  buffer_pool_manager_->UnpinPage(parent->GetPageId(), true);
23              }
24          }else{
25              if(CoalesceOrRedistribute(parent, transaction)) {
26                  buffer_pool_manager_->UnpinPage(parent->GetPageId(), false);
27                  buffer_pool_manager_->DeletePage(parent->GetPageId());
28              } else{
29                  buffer_pool_manager_->UnpinPage(parent->GetPageId(), true);
30              }
31          }
32      }else{
33          buffer_pool_manager_->UnpinPage(parent->GetPageId(), true);
34      }
35      return true;
36  }
37  }

```

Coalesce 合并操作，将兄弟节点的 key 和 value 移动到自身，这里调用了 MoveAllTo 函数，根据 index 判断合并的方向，是合并到兄弟还是合并到自身，然后调用 MoveAllTo 函数。这里的操作也可以通过模板完成。

```

1  bool BPlusTree::Coalesce(LeafPage *&neighbor_node, LeafPage *&node, InternalPage
    ↳ *&parent, int index,
2                               Txn *transaction) {
3      if(index == 0){
4          neighbor_node->MoveAllTo(node);
5          LeafPage *temp = neighbor_node;

```

```

6     neighbor_node = node;
7     node = temp;
8     parent->Remove(1);
9 }else{
10     node->MoveAllTo(neighbor_node);
11     parent->Remove(index);
12 }
13 return parent->GetSize() < parent->GetMinSize();
14 }
15
16 bool BPlusTree::Coalesce(InternalPage *&neighbor_node, InternalPage *&node,
17 ↪ InternalPage *&parent, int index,
18                               Txn *transaction) {
19     if(index == 0){
20         neighbor_node->MoveAllTo(node, parent->KeyAt(1), buffer_pool_manager_);
21         node->MoveAllTo(neighbor_node, parent->KeyAt(0), buffer_pool_manager_);
22     }else{
23         node->MoveAllTo(neighbor_node, parent->KeyAt(index), buffer_pool_manager_);
24     }
25     parent->Remove(index);
26     return parent->GetSize() < parent->GetMinSize();
27 }

```

Redistribute 重分配操作, 将兄弟节点的 key 和 value 移动到自身, 这里调用了 MoveFirstToEndOf 和 MoveLastToFrontOf 函数, 将兄弟节点的第一个 key 和 value 移动到自身的最后一个位置, 或者将兄弟节点的最后一个 key 和 value 移动到自身的第一个位置。

```

1 void BPlusTree::Redistribute(LeafPage *neighbor_node, LeafPage *node, int index) {
2     auto parent = reinterpret_cast<InternalPage>
3     ↪ *>(buffer_pool_manager_->FetchPage(node->GetParentPageId())->GetData());
4     if(index == 0){
5         neighbor_node->MoveFirstToEndOf(node);
6         parent->SetKeyAt(1, neighbor_node->KeyAt(0));
7     }else{
8         neighbor_node->MoveLastToFrontOf(node);
9         parent->SetKeyAt(index, node->KeyAt(0));
10    }
11    buffer_pool_manager_->UnpinPage(parent->GetPageId(), true);
12 }
13 void BPlusTree::Redistribute(InternalPage *neighbor_node, InternalPage *node, int
14 ↪ index) {

```

```

13  auto parent = reinterpret_cast<InternalPage
    ↳ *(buffer_pool_manager_>FetchPage(node->GetParentPageId())->GetData());
14  if(index == 0){
15      neighbor_node->MoveFirstToEndOf(node, parent->KeyAt(1), buffer_pool_manager_);
16      parent->SetKeyAt(1, neighbor_node->KeyAt(0));
17  }else{
18      neighbor_node->MoveLastToFrontOf(node, parent->KeyAt(index),
    ↳ buffer_pool_manager_);
19      parent->SetKeyAt(index, node->KeyAt(0));
20  }
21  buffer_pool_manager_>UnpinPage(parent->GetPageId(), true);
22  }

```

AdjustRoot 在前面合并结点后，如果父节点是根节点的话，需要对下面两种情况进行处理：第一是删除了根节点的最后一个元素，但是根节点还有一个子节点；第二是根节点是最后一个节点了，那么就是删除整个 b+ 树。

```

1  bool BPlusTree::AdjustRoot(BPlusTreePage *old_root_node) {
2      // case 1: when you delete the last element in root page, but root page still has
    ↳ one last child
3  if(old_root_node->GetSize() == 1 && !old_root_node->IsLeafPage()){
4      InternalPage *old_page = reinterpret_cast<InternalPage *>(old_root_node);
5      page_id_t child_page_id = old_page->RemoveAndReturnOnlyChild();
6      Page *child_page = buffer_pool_manager_>FetchPage(child_page_id);
7      BPlusTreePage *child_node = reinterpret_cast<BPlusTreePage
    ↳ *(child_page->GetData());
8      child_node->SetParentPageId(INVALID_PAGE_ID);
9      root_page_id_ = child_page_id;
10     UpdateRootPageId(0);
11     buffer_pool_manager_>UnpinPage(child_page_id, true);
12     return true;
13 }
14 // case 2: delete the last element in whole b+ tree,
15 if(old_root_node->GetSize() == 0 && old_root_node->IsLeafPage() ){
16     root_page_id_ = INVALID_PAGE_ID;
17     UpdateRootPageId(0);
18     return true;
19 }
20 return false;
21 }

```

3.3 IndexIterator

前面我们已经实现了整个 b+ 树的创建和插入删除查找的功能。现在我们实现迭代器，用于遍历 b+ 树中的叶子节点的元素。

3.3.1 IndexIterator implementation

IndexIterator 这里使用 page_id, BufferPoolManager, index 创建的迭代器，默认的构造函数就是一个空的迭代器。

```

1 IndexIterator::IndexIterator() = default;
2
3 IndexIterator::IndexIterator(page_id_t page_id, BufferPoolManager *bpm, int index)
4     : current_page_id(page_id), item_index(index), buffer_pool_manager(bpm) {
5     if(current_page_id != INVALID_PAGE_ID)
6         page = reinterpret_cast<LeafPage
7             ↳ *>(buffer_pool_manager->FetchPage(current_page_id)->GetData());

```

IndexIterator 析构函数，用于释放迭代器的资源，这里需要释放 page，因为每个迭代器都在构造时候 fetch 了一个 page。

```

1 IndexIterator::~IndexIterator() {
2     if (current_page_id != INVALID_PAGE_ID)
3         buffer_pool_manager->UnpinPage(current_page_id, false);
4 }

```

operator* * 操作符，返回当前迭代器指向的元素，这里调用了 page 的 KeyAt 和 ValueAt 函数。

```

1 std::pair<GenericKey *, RowId> IndexIterator::operator*() {
2     return std::make_pair(page->KeyAt(item_index), page->ValueAt(item_index));
3 }

```

operator++ ++ 操作符，用于迭代器的自增，这里需要判断是否到达了叶子结点的末尾，如果到达了叶子结点的末尾，那么需要将当前的 page 释放，然后将 page 指向下一个叶子结点。

```

1 IndexIterator &IndexIterator::operator++() {
2     if (item_index + 1 < page->GetSize()) {
3         item_index++;
4     } else {
5         page_id_t next_page_id = page->GetNextPageId();
6         buffer_pool_manager->UnpinPage(current_page_id, false);
7         current_page_id = next_page_id;
8         if (current_page_id != INVALID_PAGE_ID) {
9             page = reinterpret_cast<LeafPage
10                 ↳ *>(buffer_pool_manager->FetchPage(current_page_id)->GetData());
11             item_index = 0;
12         } else{
13             page = nullptr;
14             current_page_id = INVALID_PAGE_ID;
15             item_index = 0;
16         }
17     }
18     return *this;
19 }

```

operator== == 操作符，用于判断两个迭代器是否相等，这里需要判断两个迭代器的 page_id 和 item_index 是否相等。

```

1 bool IndexIterator::operator==(const IndexIterator &itr) const {
2     return current_page_id == itr.current_page_id && item_index == itr.item_index;
3 }

```

operator!= != 操作符，用于判断两个迭代器是否不相等，这里调用了 == 操作符。

```

1 bool IndexIterator::operator!=(const IndexIterator &itr) const {
2     return !(*this == itr);
3 }

```

3.3.2 BPlusTree IndexIterator

Begin 返回一个迭代器，指向 b+ 树的第一个元素，这里需要找到 b+ 树的第一个叶子结点，然后返回一个迭代器，指向第一个叶子结点的第一个元素。

```

1 IndexIterator BPlusTree::Begin() {
2     page_id_t page_id = root_page_id_;
3     Page *page = buffer_pool_manager_>FetchPage(page_id);
4     BPlusTreePage *node = reinterpret_cast<BPlusTreePage *>(page->GetData());
5     while(!node->IsLeafPage()) {
6         InternalPage *internal_node = reinterpret_cast<InternalPage *>(node);
7         page_id = internal_node->ValueAt(0);
8         buffer_pool_manager_>UnpinPage(node->GetPageId(), false);
9         node = reinterpret_cast<BPlusTreePage
10         ↪ *>(buffer_pool_manager_>FetchPage(page_id)->GetData());
11     }
12     return IndexIterator(page_id, buffer_pool_manager_);
13 }

```

另外，Begin 函数有一个 overload 版本，可以指定一个 key，返回一个迭代器，指向第一个大于等于 key 的元素。

```

1 IndexIterator BPlusTree::Begin(const GenericKey *key) {
2     Page *page = FindLeafPage(key, root_page_id_, false);
3     if(page == nullptr){
4         return IndexIterator();
5     }
6     LeafPage *node = reinterpret_cast<LeafPage *>(page->GetData());
7     int index = node->KeyIndex(key, processor_);
8     buffer_pool_manager_>UnpinPage(node->GetPageId(), false);
9     return IndexIterator(node->GetPageId(), buffer_pool_manager_, index);
10 }

```

End 返回一个迭代器，根据迭代器的定义，应当是一个空的迭代器。

```

1 IndexIterator BPlusTree::End() {
2     return IndexIterator();
3 }

```

3.4 Test

使用框架提供的测试进行测试，其中 BPlusTreeTests 测试的 n 达到 1e6，测试结果如下：

本节测试全部通过。

```
L> ./test/b_plus_tree_test
[=====] Running 1 test from 1 test suite.
[-----] Global test environment set-up.
[-----] 1 test from BPlusTreeTests
[ RUN      ] BPlusTreeTests.SampleTest
[       OK ] BPlusTreeTests.SampleTest (105630 ms)
[-----] 1 test from BPlusTreeTests (105630 ms total)

[-----] Global test environment tear-down
[=====] 1 test from 1 test suite ran. (105630 ms total)
[ PASSED  ] 1 test.
```

图 11: test7

```
L> ./test/b_plus_tree_index_test
[=====] Running 2 tests from 1 test suite.
[-----] Global test environment set-up.
[-----] 2 tests from BPlusTreeTests
[ RUN      ] BPlusTreeTests.BPlusTreeIndexGenericKeyTest
[       OK ] BPlusTreeTests.BPlusTreeIndexGenericKeyTest (0 ms)
[ RUN      ] BPlusTreeTests.BPlusTreeIndexSimpleTest
[       OK ] BPlusTreeTests.BPlusTreeIndexSimpleTest (135 ms)
[-----] 2 tests from BPlusTreeTests (136 ms total)

[-----] Global test environment tear-down
[=====] 2 tests from 1 test suite ran. (136 ms total)
[ PASSED  ] 2 tests.
```

图 12: test8

```
L> ./test/index_iterator_test
[=====] Running 1 test from 1 test suite.
[-----] Global test environment set-up.
[-----] 1 test from BPlusTreeTests
[ RUN      ] BPlusTreeTests.IndexIteratorTest
[       OK ] BPlusTreeTests.IndexIteratorTest (115 ms)
[-----] 1 test from BPlusTreeTests (115 ms total)

[-----] Global test environment tear-down
[=====] 1 test from 1 test suite ran. (115 ms total)
[ PASSED  ] 1 test.
```

图 13: test9

所给的测试样例只进行了一半的删除，我添加了剩余一半的删除测试，否则 AdjustRoot 函数无法测试。

4 Summary

我还负责对系统全部实现后进行了测试，具体的测试在验收时都已经演示过了，这里就不再赘述。至此，我个人完成的部分已经全部实现并通过了测试。对于我完成的部分，b+ 树的部分是问题最大的，后续系统功能的测试中出问题的地方也几乎都是索引引发的，经过了漫长的调试，最终得到了一个可以使用的版本。对于这样的大型工程，单元测试的完善可以帮助我们更好地实现整体的功能，我认为目前框架内提供的单元测试还不足以完全测试代码的正确性。在 tableheap 这一部分可以说几乎没有测试，这些单元测试的不完备导致我们在后面的程序调试中需要往前寻找问题而这一过程是十分困难而且耗时的。我希望能够把我完善的部分测试（lab2,3）的单元测试加入到框架中，这样可以更好地帮助我们调试程序。也希望未来框架能够越来越好，现在有一些地方的说明也不是很清晰，比如 destroy 函数的默认值，这样的设计在文档中没有提及，导致了很大的困扰。报告的撰写也花费了不少的精力，感谢能看到这里的助教老师，希望我们数据库系统这门课能越来越好，也希望这样一个巨大工作量和难度的大程能在考核中获得更多的分数，而不是仅仅 15 分。