

LSM-Tree

2020 年 4 月 19 日

设计

文件结构

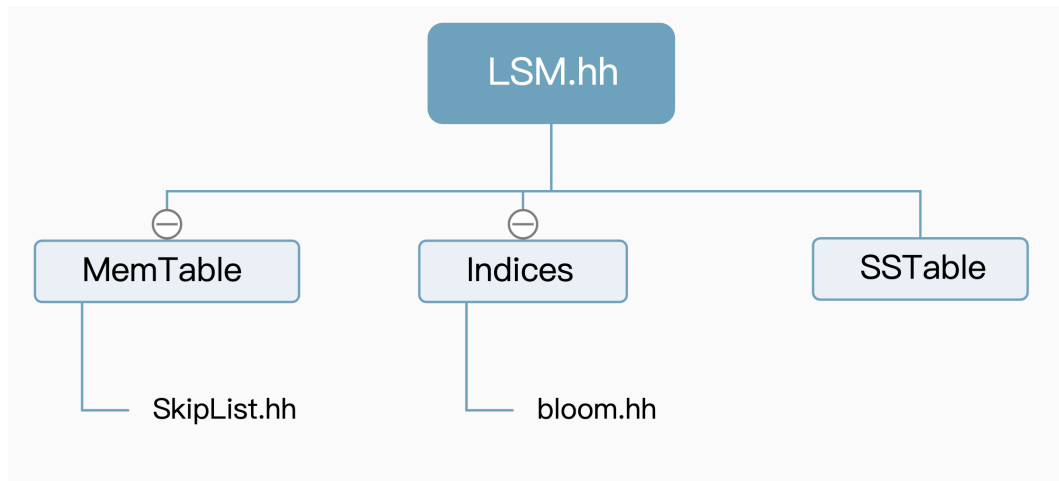


图 1: File Structure

上图为 LSM-Tree 的功能模块以及依赖的文件。LSM 为接口类,负责管理内存中的 SkipList(键值对的内存存储) 以及 Indices(SSTable 在内存中的缓存), 以及实现了对磁盘中 SSTable 的操作(如 dump,read,compact 等)。

MemTable

内存中的存储由跳表实现。我实现了字节统计和导出数据, 便于 LSM 调用。除此之外与上课讲的大致相同, 在此便不再赘述。

SSTable

Size{ 4 }	dataSegBias{ 4 }	indiceSeg{ key{ 8 } + dataBias{ 4 } + dataLen{ 4 } }	dataSeg{?}
-----------	------------------	--	------------

图 2: SSTable Structure

上图为 SSTable 的结构。首先是 32 位无符号整数 size(冗余, uint32_t), 接着是存储了数据段 (dataSeg, uint32_t) 关于文件开头的字节偏移量的 dataSegBias, 再接着是存储了每个数据的键值 (key, uint64_t)、关于数据段开头的偏置 (dataBias, uint32_t) 以及数据的字节长度 (dataLen, uint32_t), 最后是存储了所有值的数据段 dataSeg(sum(dataLen))。

SSTable 我将其主要功能设计为实现文件二进制数据 (Struct Bin) 与内存中键值对 (vector<Entry<uint64_t, string> >) 的转换由两个构造函数实现, 另外还实现了一些接口用来方便调用。(以下代码仅为示意)

```
1| template<K, V>
2| class SST {
3|     explicit SST(const vector<Entry<K, V> > &_data, uint32_t _dataBytes);
4|     explicit SST(char *_bin)
5|     explicit SST(const SST<K, V> & ano);
6|     ~SST();
7|     Bin toBin();
8|     Bin toIndexBin();
9|     vector<Entry<K, V> > &vecData();
10|    uint32_t getSize() const;
11|    uint32_t getDataSegBias() const;
12|    K getLowBound() const;
13|    K getHighBound() const;
14| };
```

Indices

内存中的索引我设计为以查找以及根据二进制文件自动初始化为主要功能, 由 Indices(单个 SSTable 的索引) 和 IndicesTab(Indices 的集合, 包含在 Disk 结构中的层次信息)。

在硬盘中, 我使用文件名来区分文件所处的层次结构 (levelinLevel.bin), 如 ("00.bin", "315.bin")。我设定第 0 层的 MAX_SIZE 为 4, 其余层数为 $2^n \times 4$ 。

基于以下逻辑, 该命名法还可以用来保证 lazy 方式下不使用时间戳时操作的正确性。在第 0 层中, inLevel 越大越新 (03 > 02 > 01 > 00); 在其他层中, 因为同层键值对不会相交, 故层内不必区分时间顺序, 而又因为归并是向下的, 上层一定新于下层。故归纳可得

$$MemTable > 03.bin > 02.bin > 01.bin > 00.bin > level1 > level2 > ...!$$

```

1| template<K>
2| class IndicesTab {
3|     explicit IndicesTab(const string &_dir);
4|     bool insert(const Indices<K> &idx, string &filename);
5|     vector<Indices<K>> *rLevel(uint32_t levelNum);
6|     uint32_t getHeight() const;
7|     bool find(const K &key,
8|              string *filename = nullptr, uint32_t *dataSegBias= nullptr,
9|              uint32_t *bias = nullptr, uint32_t *length = nullptr);
10|     void addNewLevel();
11|     void clear();
12| }
13|
14| template<K>
15| class Indices {
16|     explicit Indices(const Bin &bin, uint32_t __size, uint32_t __dataSegBias);
17|     bool find(const K &k, uint32_t *b = nullptr, uint32_t *l = nullptr);
18|     uint32_t getSize() const;
19|     uint32_t getDataSegBias() const;
20|     uint32_t getLowBound() const;
21|     uint32_t getHighBound() const;
22| }

```

Bonus

除了基本功能之外，我还使用了布隆过滤器用来快速判断某个 key 是否有可能在索引中，且 LSM-Tree 是以模板的方式实现的，能够存储包括 <uint64_t, string> 在内的其他键值对 (因为模板的实现与声明不能分开，故我的所有实现都在.hh 文件中)。

测试

测试环境

机器: MacBook Pro (13-inch, 2018, Four Thunderbolt 3 Ports)

系统: macOS Catalina 10.15.3

硬盘: APPLE SSD AP0512M

CPU: Intel i5-8259U

Correctness Test

对一定范围内的 key，我使用长度为 1-100 的随机字符串进行赋值，并同时随机地对 LSM 和 SkipList 进行操作，最后进行比对测试两者结果是否相同。测试结果：

CorrectnessTestResult : 262144/262144 => 100%

Latency Test

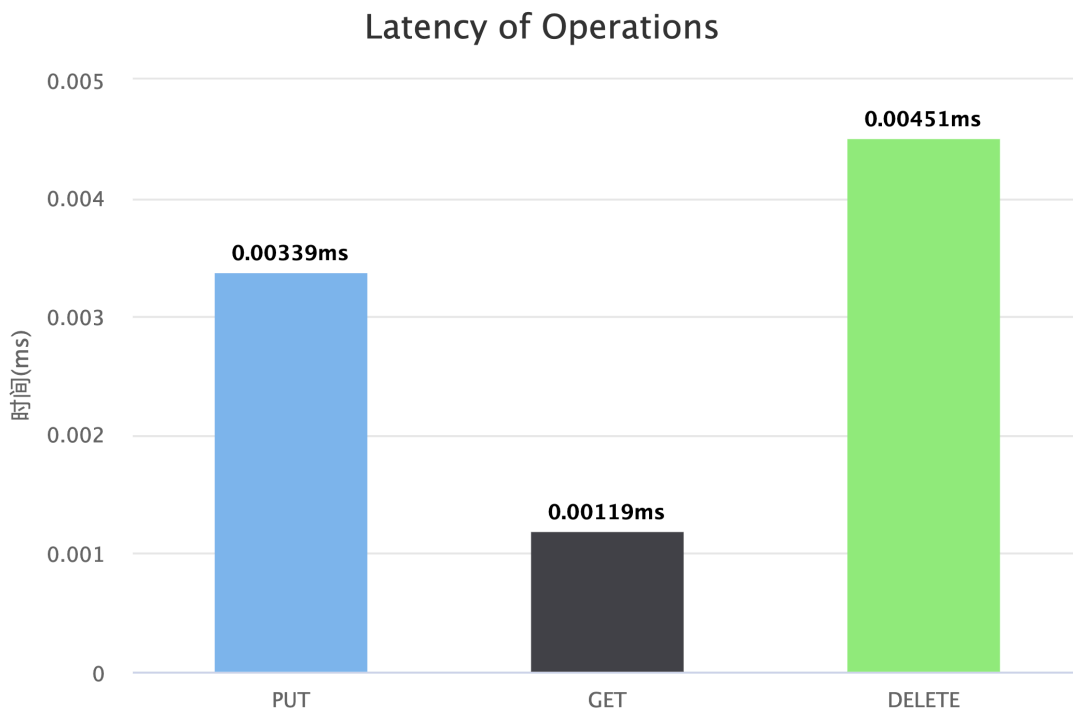


图 3: Latency Test

上图为非 Compact 情况下的 PUT, Get, Delete 延迟测试结果。可见 $DELETE > PUT > GET$ 且有 $DELETE \approx GET + PUT$ 。原因是在 DELETE 时先要进行与 GET 类似的查询操作，而我将 DELETE 的 lazy 操作用等价的 `put(key, "")` 替代了，故 DELETE 操作需经过约 PUT 和 GET 加和的时间。

Throughput Test

可能是因为固态硬盘随机读取性能较好的缘故, 1s 内难以体现出 Compaction 对吞吐的影响, 于是我把时间设定为 0.2ms 进行测试, 结果如下。

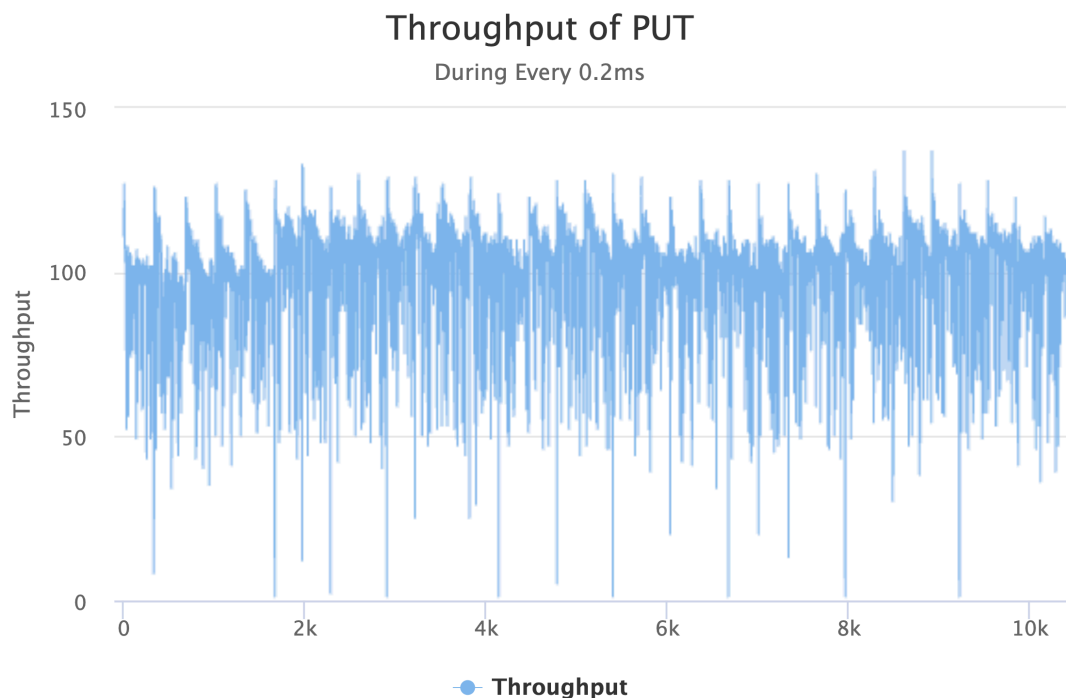


图 4: Throughput Test

上图可见随着 put 的元素越来越多, LSM 的吞吐量呈现周期性的波动, 但总体稳定在 100/0.2ms 左右

Code

```
1| #define TEST_SIZE (1 << 20)
2|
3| char *randstr(char *str, const int len) {
4|     int i;
5|     for (i = 0; i < len; ++i) {
6|         switch ((rand() % 3)) {
7|             case 1:
8|                 str[i] = 'A' + rand() % 26;
9|                 break;
10|            case 2:
11|                str[i] = 'a' + rand() % 26;
12|                break;
13|            default:
```

```

14|         str[i] = '0' + rand() % 10;
15|         break;
16|     }
17| }
18| return str;
19| }
20|
21| void doSomething(LSM<uint64_t, string> &lsm, uint64_t key, SkipList<uint64_t, string> *memTab
= nullptr) {
22|     switch(rand() % 2) {
23|     case 0: {
24|         char ranStr[100]; int len = rand() % 99;
25|         randstr(ranStr, len);
26|         lsm.put(key, string(ranStr, len));
27|         if (memTab) { memTab->put(key, string(ranStr, len)); }
28|         break;
29|     }
30|     case 1: {
31|         lsm.remove(key);
32|         if (memTab) { memTab->remove(key); }
33|         break;
34|     }
35|     }
36| }
37|
38| void correctnessTest(LSM<uint64_t, string> &lsm, uint64_t size) {
39|     lsm.reset();
40|     SkipList<uint64_t, string> memTab;
41|     for (uint64_t i = 0; i < size; ++i) {
42|         doSomething(lsm, i, &memTab);
43|     }
44|     for (uint64_t i = 0; i < size; ++i) {
45|         doSomething(lsm, i, &memTab);
46|     }
47|
48|     uint64_t cnt = 0;
49|     auto data = memTab.data();
50|     for (uint64_t i = 0; i < size; ++i) {
51|         string lsmGet = lsm.get(i), *memGet = memTab.get(i);
52|         if (memGet) {
53|             if (lsmGet != *memGet) {
54|                 cout << "LSM get: " << lsmGet << endl << "MenTable get: " << *memGet << endl;
55|             }
56|             else { ++cnt; }
57|         }
58|         else {
59|             if (!lsmGet.empty()) { cout << "LSM get: " << lsmGet << endl << "MenTable get: " <<
endl; }
60|             else { ++cnt; }
61|         }
62|     }
63|     cout << "Correctness Test Result: " << cnt << '/' << size << " => " << double(cnt) / size
* 100 << '%' << endl;
64| }
65|
66| struct Lat {
67|     double putLat;
68|     double getLat;
69|     double delLat;
70|     uint64_t size;
71| };
72|
73| void latencyTest(LSM<uint64_t, string> &lsm, uint64_t size) {
74|     lsm.reset();
75|     vector<Lat> latRecd;
76|     clock_t start;
77|     for (uint64_t i = 0; i < size; ++i) {
78|         Lat lat;
79|         char ranStr[100]; int len = rand() % 100;
80|         randstr(ranStr, len);
81|         string input = string(ranStr, len);
82|         start = clock(); bool flag = lsm.put(i, input); lat.putLat = double(clock() - start) /
CLOCKS_PER_SEC * 1000;
83|         if (flag) {
84|             start = clock(); lsm.get(i); lat.getLat = double(clock() - start) / CLOCKS_PER_SEC
85|

```

```

|      * 1000;
87|      start = clock(); lsm.remove(i); lat.delLat = double(clock() - start) /
|      CLOCKS_PER_SEC * 1000;
88|      latRecd.push_back(lat);
89|  }
90|  }
91|
92|  double avePut = 0, aveGet = 0, aveDel = 0;
93|  for (auto i = latRecd.begin(); i != latRecd.end(); ++i) {
94|      avePut += i->putLat;
95|      aveGet += i->getLat;
96|      aveDel += i->delLat;
97|  }
98|
99|  avePut /= latRecd.size(); aveGet /= latRecd.size(); aveDel /= latRecd.size();
100|
101|  cout << "Latency Test Result: " << "Put: " << avePut << "ms " << "Get: " << aveGet << "ms
|  " << "Delete: " << aveDel << "ms" << endl;
102|  }
103|
104|  void throughputTest(LSM<uint64_t, string> &lsm, uint64_t size) {
105|      lsm.reset();
106|      uint64_t i = 0;
107|      vector<int> t;
108|      char data[100][100];
109|      for (int i = 0; i < 100; ++i) { randstr(data[i], 100); }
110|
111|      clock_t start = clock(); uint64_t cnt = 0;
112|      while (i < size) {
113|          if (clock() - start > 200) {
114|              t.push_back(cnt);
115|              cnt = 0;
116|              start = clock();
117|          }
118|          lsm.put(i, string(data[i % 100], rand() % 100)); ++cnt; ++i;
119|      }
120|      cout << "Throughput Test Result:" << endl;
121|      cout << '[';
122|      uint64_t c = 0;
123|      for (auto &i:t) {
124|          cout << '[' << c++ << ', ' << i << ']' << ', ';
125|      }
126|      cout << ']';
127|  }
128|
129|  int main() {
130|      srand(time(0));
131|      LSM<uint64_t, string> lsm("./data");
132|      correctnessTest(lsm, TEST_SIZE);
133|      latencyTest(lsm, TEST_SIZE);
134|      throughputTest(lsm, TEST_SIZE);
135|      return 0;
136|  }

```