

数据构建管理系统 LAB1--为LevelDB增加TTL

远程仓库地址: <https://gitea.shuishan.net.cn/10224602413/levelDB-LAB1.git>

协作者:

曹可心-10223903406

朴祉燕-10224602413

目录

[一、LevelDB的简要介绍](#)

[二、TTL功能介绍](#)

[三、设计思路](#)

[四、实现过程](#)

[4.1 修改数据结构--PUT](#)

[4.2 修改读取流程--GET](#)

[4.3 修改合并流程--Compaction](#)

[五、测试用例](#)

[六、特点/出现的问题及解决](#)

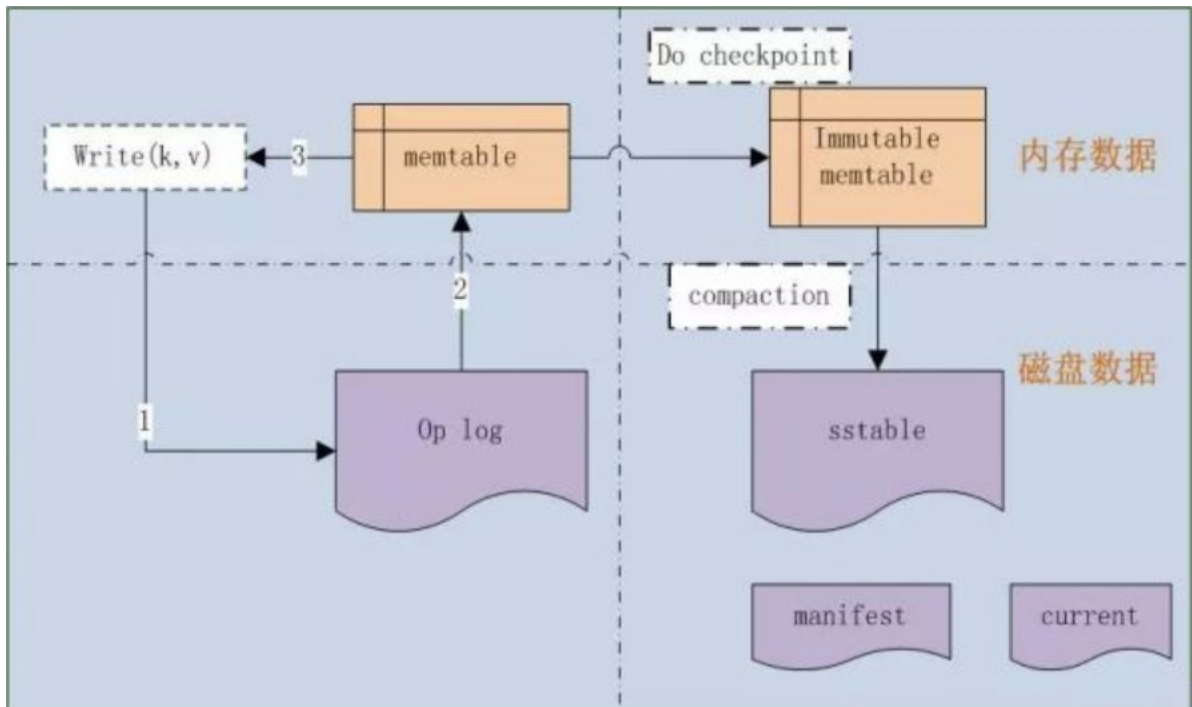
[七、总结](#)

一、LevelDB的简要介绍

LevelDB, 是一种基于operation log的文件系统, 是Log-Structured-Merge Tree的典型实现。是Google开发的单机K-V存储系统, 涉及到了skip list、内存KV table、LRU cache管理、table文件存储、operation log系统等。

1.1 LevelDB的框架

Leveldb的基本框架, 几大关键组件, 如下图所示:



由于采用了op log，它就可以把随机的磁盘写操作，变成了对op log的append操作，因此提高了IO效率，最新的数据则存储在内存memtable中。

当op log文件大小超过限定值时，就定时做check point。Leveldb会生成新的Log文件和Memtable，后台调度会将Immutable Memtable的数据导出到磁盘，形成一个新的SSTable文件。SSTable就是由内存中的数据不断导出并进行Compaction操作后形成的，而且SSTable的所有文件是一种层级结构，第一层为Level 0，第二层为Level 1，依次类推，层级逐渐增高，这也是为何称之为LevelDb的原因。

1.2 LevelDB的数据结构

1.2.1 Slice

- 包括length和一个指向外部字节数组的指针。
- 和string一样，允许字符串中包含'\0'。

提供一些基本接口，可以把const char和string转换为Slice；把Slice转换为string，取得数据指针const char。

1.2.2 Status

Leveldb 中的返回状态，将错误号和错误信息封装成Status类，统一进行处理。并定义了几种具体的返回状态，如成功或者文件不存在等。

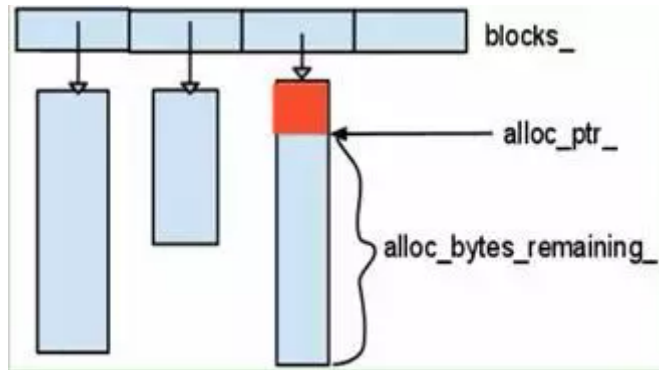
为了节省空间Status并没有用std::string来存储错误信息，而是将返回码(code), 错误信息message及长度打包存储于一个字符串数组中。

成功状态OK 是NULL state，否则state 是一个包含如下信息的数组：

```
state_[0..3] == 消息message长度
state_[4]    == 消息code
state_[5..] ==消息message
```

1.2.3 Arena

Leveldb的简单的内存池，它所作的工作十分简单，申请内存时，将申请到的内存块放入std::vector blocks_中，在Arena的生命周期结束后，统一释放掉所有申请到的内存，内部结构如下图所示。



Arena主要提供了两个申请函数：其中一个直接分配内存，另一个可以申请对齐的内存空间。

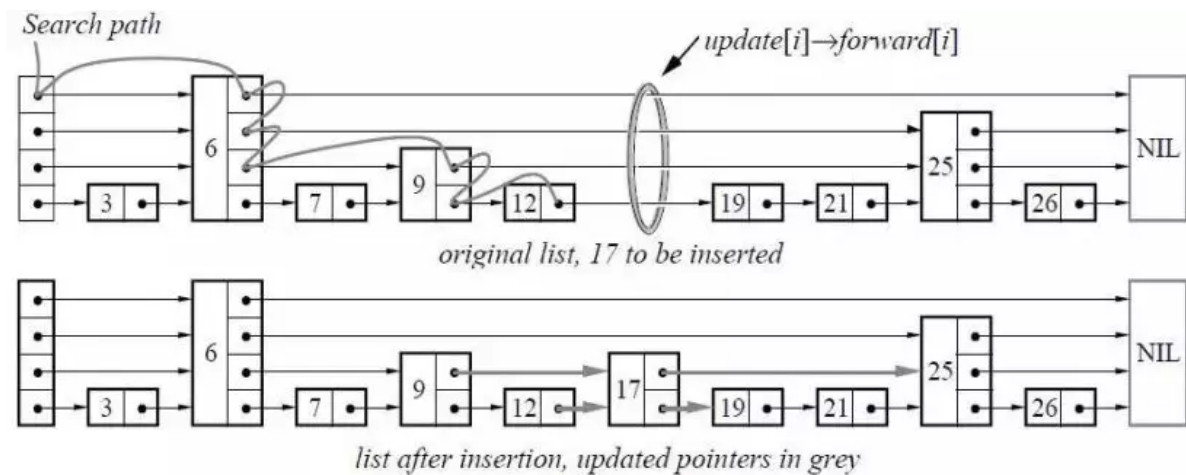
Arena没有直接调用delete/free函数，而是由Arena的析构函数统一释放所有的内存。

这应该是和leveldb特定的应用场景相关的，比如一个memtable使用一个Arena，当memtable被释放时，由Arena统一释放其内存。

1.2.4 Skip list

Skip list(跳跃表)是一种可以代替平衡树的数据结构。Skip lists应用概率保证平衡，平衡树采用严格的旋转（比如平衡二叉树有左旋右旋）来保证平衡，因此Skip list比较容易实现，而且相比平衡树有着较高的运行效率。

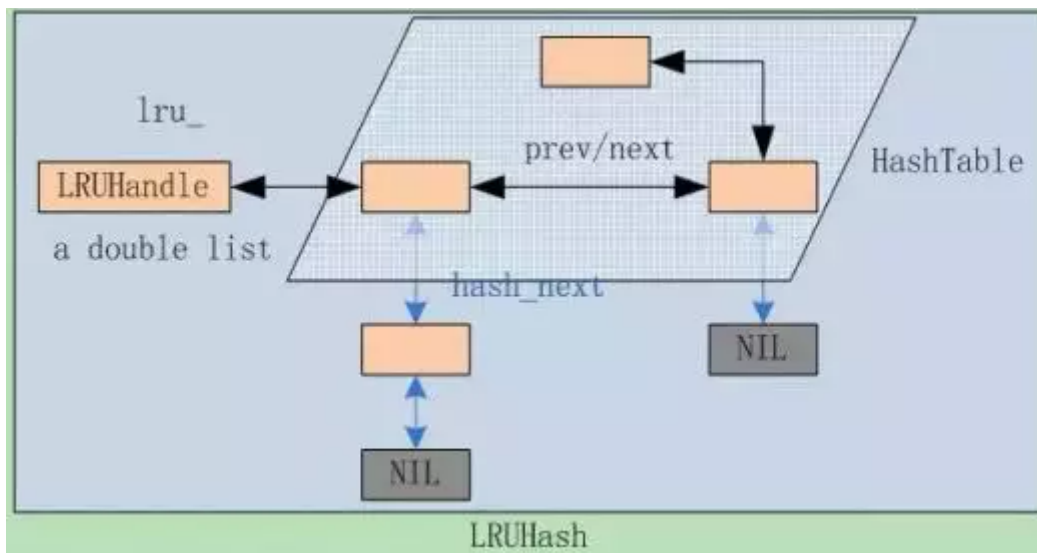
从概率上保持数据结构的平衡比显式的保持数据结构平衡要简单的多。对于大多数应用，用skip list要比用树更自然，算法也会相对简单。由于skip list比较简单，实现起来会比较容易，虽然和平衡树有着相同的时间复杂度($O(\log n)$)，但是skip list的常数项相对小很多。skip list在空间上也比较节省。一个节点平均只需要1.333个指针（甚至更少），并且不需要存储保持平衡的变量。



在Leveldb中，skip list是实现memtable的核心数据结构，memtable的KV数据都存储在skip list中。

1.2.5 Cache

Leveldb内部通过双向链表实现了一个标准版的LRUCache，先上个示意图，看看几个数据之间的关系：



1.3 LevelDB的主要流程

1.3.1 open

1. 基本检查

- a. 根据传入的 db 路径，对 LOCK 文件做 flock 来判断是否已经有 db 实例启动，一份数据同时只能有一个 db 实例操作。
- b. 根据 option 内的 `create_if_missing/error_if_exists`，来确定当数据目录已经存在时要做的处理。

2. db 元信息检查 (VersionSet::recover())

- a. 从 CURRENT 文件中读取当前的 MANIFEST 文件。
- b. 从 MANIFEST 文件中依次读取每个 record (`VersionEdit::DecodeFrom`)，检查 Comparator 是否一致，然后依次 replay。
- c. 检查解析 MANIFEST 的最终状态中的基本的信息是否完整 (log number, FileNumber, SequenceNumber)，将其生效成 db 当前的状态。此时，整个 db 的各种元信息 (FileNumber, SequenceNumber, 各 level 的文件数目, size, range, 下一次 compact 的 start_key 等等) 均 load 完成，db 恢复成上一次退出前的状态。

3. 从 log 中恢复上一次可能丢失的数据 (RecoverLogFile)

4. 生成新的 log 文件。更新 db 的元信息 (`VersionSet::LogAndApply()`，生成最新的 MANIFEST 文件)，删除无用文件 (`DeleteObsoleteFiles()`)，尝试 compact (`MaybeScheduleCompaction()`)。

5. 启动完毕

1.3.2 put

leveldb 中的写操作不是瓶颈，但可能出现过量写影响读的效率（比如 level-0 中文件过多，查找某个 key 可能会造成过量的 io），所以有一系列策略主动去限制写。

1. 将 key value 封装成 WriteBatch.
2. 循环检查当前 db 状态，确定策略(`DBImpl:: MakeRoomForWrite()`):
3. 设置 WriteBatch 的 SequenceNumber.
4. 先将 WriteBatch 中的数据记 log(`Log::AddRecord()`).
5. 将 WriteBatch 应用在 memtable 上。 (`WriteBatchInternal::InsertInto()`)，即遍历 decode 出 WriteBatch 中的 key/value/ValueType，根据 ValueType 对 memetable 进行 put/delete 操作。
6. 更新 SequenceNumber (`last_sequence + WriteBatch::count()`)。

1.3.3 get

总体来说, get 即是找到 userkey 相同, 并且 SequenceNumber 最大 (最新) 的数据。leveldb 支持对特定 Snapshot 的 get, 只是简单的将 Snapshot 的 SequenceNumber 作为最大的 SequenceNumber 即可。

1. 如果 ReadOption 指定了 snapshot, 则将指定 snapshot 的 SequenceNumber 作为最大 SequenceNumber, 否则, 将当前最大 SequenceNumber (VersionSet::last_sequence_number) 作为最大值。
2. 在 memtable 中查找(MemTable::Get())
3. 如果 memtable 中未找到, 并且存在 immutable memtable, 就在 immutable memtable 中查找 (Memtable::Get())。
4. 仍未找到, 在 sstable 中查找(VersionSet::Get())。从 level-0 开始, 每个 level 上依次进行查找, 一旦找到, 即返回。

1.3.4 delete

delete 相比于 put 操作, 只在构造 WriteBatch 时, 设置 ValueType 为 kTypeDeletion, 其他流程和 put 相同。

1.3.5.snapshot

1. 取得当前的 SequenceNumber
2. 构造出 Snapshot, 插入到已有链表中。

1.3.6.NewIterator

构造 DBIter, 做 Seek() 即可。参见 Iterator。

1.3.7.compact

leveldb 中有且仅有一个后台进程 (第一次 compact 触发时 create 出来) 单独做 compact。当主线程主动触发 compact 时 (MaybeScheduleCompaction()), 做以下流程:

1. 如果 compact 已经运行或者 db 正在退出, 直接返回。
2. 检查当前的运行状态, 确定是否需要进行 compact, 如果需要, 则触发后台调度 compact (Env::Schedule()), 否则直接返回。
3. 做实际的 compact 逻辑 (BackgroundCompaction()), 完成后, 再次主动触发 compact (主线程将任务入队列即返回, 不会有递归栈溢出的问题)。

二、TTL功能介绍

2.1 TTL功能

TTL (Time To Live), 即生存时间, 是指数据在存储系统中的有效期。设置 TTL 可以使得过期的数据自动失效, 减少存储空间占用, 提高系统性能。

为什么需要 TTL 功能:

- 数据自动过期: 无需手动删除过期数据, 简化数据管理。
- 节省存储空间: 定期清理无效数据, 优化资源利用。
- 提高性能: 减少无效数据的干扰, 提升读写效率。

2.2 实验要求

- 在LevelDB中实现键值对的TTL功能，使得过期的数据在读取时自动失效，并在适当的时候被合并清理。
- 修改LevelDB的源码，实现对TTL的支持，包括数据的写入、读取和过期数据的清理。
- 编写测试用例，验证TTL功能的正确性和稳定性。

三、设计思路

3.1 数据结构组织

WriteBatch结构由一个header和若干个entry组成，每个entry由kTypeValue、key和value（delete类型的entry中没有value）组成。时间戳可以被加入到key或value中，也可以单独将时间戳的长度和值追加在entry尾部。

考虑到key的使用频率较高，如果在key中封装时间戳，每当在查询阶段比较key值大小时都需要对key进行解析抽取键值，时间效率较低；如果将时间戳单独封装在entry的尾部，一方面时间戳的长度也需要单独存储，导致空间浪费，另一方面在memtable和sstable中的数据结构也需要重新组织；综合考虑，我们选择在值中加入过期时间的信息。

3.2 过期数据读取

读取时，依次从memtable、immutable memtable和sstable中查找目标key的思路不变，如果目标key值没有被找到，则结果的处理方式也不变。如果目标key值被找到，则需要对相应的value进行解析，将过期时间戳与当前时间进行比较，如果数据已经过期需要返回NotFound。

question:如果目标key值第一次被找到时发现已经过期，是否需要继续寻找？

answer:如果目标key是在sstable的level0中第一次被找到，我们需要考虑到是否会在Level0层的其他sstable中含有该键值对的更新版本。但由于查找操作是在level0层的sstable被排序后依次进行的，因此可以保证第一次找到的键值对就是最新版本。

如果目标key在其他地方被找到，线然仅需根据第一次找到的value的过期情况判断数据是否有效。

3.3 过期数据删除

leveldb中的合并方式有minor compaction、seek compaction、size compaction和manual compaction四种，我们计划在无论哪种合并中一旦检查到过期的数据，都将其从数据库中直接删除，在四种合并中检查和删除的逻辑是相同的，因此我们希望找到四种合并共有的代码段进行修改以减少重复工作。

3.4 分工合作

在分配任务阶段，重点在于确定能否将读取数据与合并数据两个部分交由两名组员分别实现。

本质上来说两部分所需的工作都是在合适的地方增加解析出时间戳的操作，并相应地修改返回结果或数据组织的逻辑。只要在插入数据时固定了时间戳的格式，两部分的代码就是相互独立的。但是，如果负责合并部分的同学先于负责读取部分的同学完成任务，是否会影响合并部分的测试呢？

我们仔细研究了ttl_test.cc中对于合并部分的测试，得出两部分的实现顺序没有先后性的结论。在TEST(TestTTL, CompactionTTL)中，db->CompactRange(nullptr, nullptr)触发全范围的手动合并后，通过db->GetApproximateSizes(ranges, 1, sizes)来计算数据库中键值在ranges[0]范围内的文件大小，最后检查返回的sizes，这个过程中并不调用读取数据相关的接口，而是有一套独立的查询通道。

四、实现过程

4.1 修改数据结构--PUT

在插入数据的batch.Put()参数列表中增加ttl，将过期时间戳封装到batch中。

```
// 假设增加一个新的Put接口，包含TTL参数，单位（秒）
Status DB::Put(const WriteOptions& opt, const Slice& key, const Slice& value,
uint64_t ttl){
    WriteBatch batch;
    batch.Put(key, value, ttl);
    return Write(opt, &batch);
}
```

具体的封装过程：

计算出过期时间，将过期时间转化成格式化的固定长度的字符串，并与原本的value值拼接在一起，转化为slice数据类型插入到batch.rep_中。

```
void WriteBatch::Put(const Slice& key, const Slice& value, std::uint64_t ttl) {
    WriteBatchInternal::SetCount(this, WriteBatchInternal::Count(this) + 1);
    rep_.push_back(static_cast<char>(kTypeValue));
    PutLengthPrefixedSlice(&rep_, key);

    // 获取当前时间
    auto now = std::chrono::system_clock::now();

    // 加上ttl
    auto future_time = now + std::chrono::seconds(ttl);

    // 转换为 time_t
    std::time_t future_time_t = std::chrono::system_clock::to_time_t(future_time);

    // 将 time_t 转换为 tm 结构
    std::tm* local_tm = std::localtime(&future_time_t);

    // 格式化为字符串
    char buffer[20]; // 格式化字符串的缓冲区
    std::strftime(buffer, sizeof(buffer), "%Y-%m-%d %H:%M:%S", local_tm);
    std::string future_time_str(buffer);

    // 拼接原本的值和时间字符串
    std::string combined_str = value.ToString() + future_time_str;

    PutLengthPrefixedSlice(&rep_, Slice(combined_str));
}
```

4.2 修改读取流程--GET

成功将过期时间戳与值一起存储后，需要在Get操作中，解析存储的值，判断是否过期。

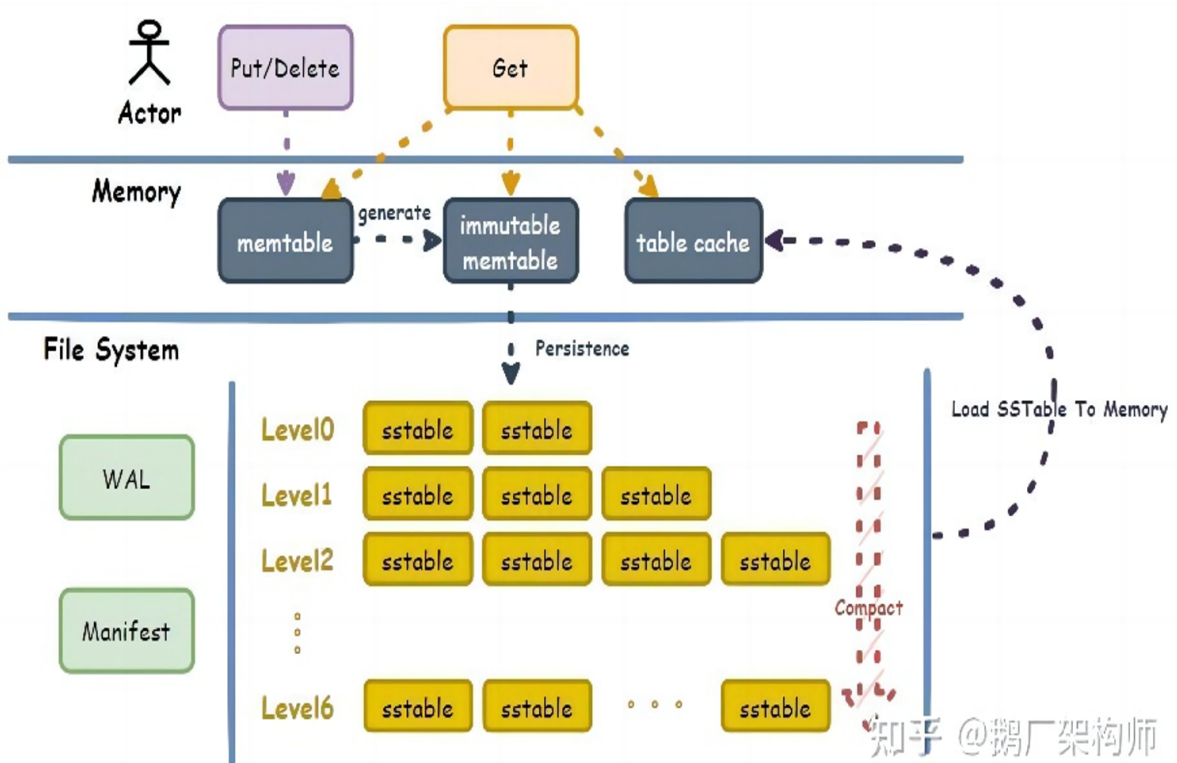
LevelDB的整个读取流程有很多个函数嵌套来完成的。具体如下：

```

DBImpl::Get()
• MemTable::Get() //跳表中读取键值对
• Version::Get() //SSTable中读取键值对
  ◦ Version::ForEachOverlapping()
    § FindFile()
    § TableCache::Get()
      □ TableCache::FindTable()
      □ Table::InternalGet()
        ® FilterBlockReader::KeyMayMatch()
        ® Block::Iter::Seek()

```

也就是说。LevelDB会在调用 `DBImpl::Get()` 时先通过 `MemTable::Get()` 去memtable查看我们要查找的数据在不在memtable里面，如果有，就返回；没有就继续通过 `Version::Get()` 去sstable里查找。



那么，我们增加TTL后的数据读取的相关修改就要分为两部分去进行：一个是 `MemTable::Get()`，一个是 `Version::Get()`。所以，`DBImpl::Get` 的主要代码就是下面条件判断这一部分：

```

Status DBImpl::Get(const ReadOptions& options, const Slice& key,
                  std::string* value) {
  Status s;
  MutexLock l(&mutex_);
  .....

  if (mem->Get(lkey, value, &s)) { //尝试从当前 memtable 中获取键的值。如果找到，则
    直接返回!! 燕
    // Done
  } else if (imm != nullptr && imm->Get(lkey, value, &s)) { //memtable没有，去
    imm找!! 燕
    // Done
  } else {

```



```

    s = current->Get(options, lkey, value, &stats); //最后去sstable找!! 燕
    have_stat_update = true;
}
mutex_.Lock();
.....
return s; //返回状态对象 s, 告知调用方是否成功找到值
}

```

4.2.1 MemTable::Get()的修改

从 `if (mem->Get(lkey, value, &s))` 跳转至 `db/memtable.cc` 中的 `MemTable::Get()`, 发现具体 `get` 功能就是在这个函数里实现的, 没有在里面调用的其他函数, 那么在为 `MemTable::Get` 函数增加 TTL 功能时, 只要通过在读取数据时解析过期时间, 并与当前时间进行比较, 从而决定数据是否有效就可以了。

1. 获取存储格式中的原始值和时间戳

由于在写入时将数据与过期时间拼接为一个字符串格式存储, 所以在读取时需首先将该组合字符串解包, 分离出实际的值部分和过期时间:

- 使用 `slice v = GetLengthPrefixedSlice(key_ptr + key_length);` 获取存储在 `MemTable` 中的完整组合字符串。
- 通过 `std::string combined_str(v.data(), v.size());` 将 `Slice` 转换为 `std::string`, 便于进一步分离数据。
- 在 `combined_str` 中, 提取前部分作为原始值 `actual_value`, 即 `combined_str.substr(0, combined_str.size() - 20);`。
- 由于过期时间占最后 19 个字符, 将这部分取出作为过期时间字符串 `time_str`。

```

// 获取存储的值和时间戳
Slice v = GetLengthPrefixedSlice(key_ptr + key_length);
std::string combined_str(v.data(), v.size());

// 根据存储格式分离原始值和时间戳
std::string actual_value = combined_str.substr(0, combined_str.size() - 20);
std::string time_str = combined_str.substr(combined_str.size() - 19, 19);

```

2. 解析当前时间字符串

为了便于直接与存储的过期时间比较, 获取系统的当前时间并格式化成字符串:

- 使用 `std::chrono::system_clock::now()` 得到当前时间。
- 使用 `std::strftime` 将当前时间格式化为 `"%Y-%m-%d %H:%M:%S"`, 与存储格式一致的时间字符串 `current_time_str`, 便于后续比较。

```

// 获取当前时间 (字符串格式)
auto now = std::chrono::system_clock::now();
auto now_time_t = std::chrono::system_clock::to_time_t(now);
std::tm* now_tm = std::localtime(&now_time_t);
char buffer[20];
std::strftime(buffer, sizeof(buffer), "%Y-%m-%d %H:%M:%S", now_tm);
std::string current_time_str(buffer);

```

3. 过期时间检查

通过比较存储的 `time_str` 和当前时间 `current_time_str`，判断数据是否过期：

- 如果 `time_str` 小于或等于 `current_time_str`，说明数据已经过期，此时设置 `Status` 为 `NotFound`，并返回已过期的信息。
- 如果未过期，表示数据有效，将 `actual_value` 赋值给 `value`，返回实际存储的值。

```
// 检查过期
if (time_str <= current_time_str) {
    *s = Status::NotFound("key has expired"); // 已过期
    return true;
}

// 未过期，返回实际值
value->assign(actual_value);
return true;
```

4. 返回数据

通过上述步骤，最终可以确保只返回未过期的有效数据。

```
case kTypeDeletion:
    *s = Status::NotFound(Slice());
    return true;
```

4.2.2 Version::Get()的修改

从 `else {s = current->Get(options, lkey, value, &stats);` 跳转至 `db/version_set.cc` 中的 `Version::Get()`，发现具体 `get` 功能还要继续往里查。

如果直接在 `Version::Get()` 里实现 TTL 过期判断的逻辑，可能会造成一层一层的函数嵌套进去，获取到了具体数据，再一层一层传出来，最后发现过期返回没找到。这样可能会比较泄气，但是代码实现也许比较简单，所以一开始是想直接在 `Version::Get()` 里面添加 TTL 过期判断逻辑代码的。

```
Status Version::Get(const ReadOptions& options, const LookupKey& k,
                   std::string* value, GetStats* stats) {
    stats->seek_file = nullptr;
    stats->seek_file_level = -1;

    struct State {
        Saver saver;
        GetStats* stats;
        .....
    };

    static bool Match(void* arg, int level, FileMetaData* f) {
        State* state = reinterpret_cast<State*>(arg);

        if (state->stats->seek_file == nullptr &&
            state->last_file_read != nullptr) { //如果 seek_file 为空且
last_file_read 不为空
            // We have had more than one seek for this read.  Charge the 1st file.
            // 记录第一次的查找文件信息
            state->stats->seek_file = state->last_file_read; //则记录 seek_file 以标记
第一次读取的文件和层级。
            state->stats->seek_file_level = state->last_file_read_level;
```

```

}

state->last_file_read = f; //更新 last_file_read 和 last_file_read_level
state->last_file_read_level = level;
//调用 table_cache->Get: 从缓存中获取指定文件, 并在查找键时使用 savevalue 回调,
state->s = state->vset->table_cache->Get(*state->options, f->number,
                                         f->file_size, state->ikey,
                                         &state->saver, SaveValue);

// 使用TableCache::Get并传递自定义的SaveValue函数以进行TTL检查, 燕改
auto ttl_save_value = [](void* arg, const Slice& key, const Slice& value)
{
    Saver* saver = reinterpret_cast<Saver*>(arg);
    if (value.size() < 19) {
        saver->state = kNotFound;
        return;
    }

    // 解析时间戳并检查过期情况
    uint64_t expire_time = DecodeFixed64(value.data() + value.size() - 19);
    if (expire_time > Env::Default()->NowMicros()) {
        SaveValue(arg, key, value); // 未过期, 保存数据
    } else {
        saver->state = kNotFound; // 数据已过期
    }
};

// 调用Get并传递自定义的ttl_save_value, 燕改
state->s = state->vset->table_cache->Get(*state->options, f->number,
                                         f->file_size, state->ikey,
                                         &state->saver, ttl_save_value);

if (!state->s.ok()) {
    return false;
}
.....
};
}}

```

然后, 进行调试。无论怎么打断点调试, 始终都是:

```

Actual: false
Expected: true
[ FAILED ] TestTTL.ReadTTL (38035 ms)
[ RUN      ] TestTTL.CompactionTTL
open db failed
Aborted (core dumped)
jiyeon@judyPark:/mnt/e/aaa/课件/database/lab1/leveldb_base/build$ ./mt/e/aaa/课件/database/lab1/leveldb_base/build/ttl_test
[=====] Running 2 tests from 1 test suite.
[-----] Global test environment set-up.
[-----] 2 tests from TestTTL
[ RUN      ] TestTTL.ReadTTL
!!!!!!!!!!!!/mnt/e/aaa/课件/database/lab1/leveldb_base/test/ttl_test.cc:71: Failure
Value of: status.ok()
Actual: false
Expected: true
[ FAILED ] TestTTL.ReadTTL (40652 ms)
[ RUN      ] TestTTL.CompactionTTL
open db failed
Aborted (core dumped)
jiyeon@judyPark:/mnt/e/aaa/课件/database/lab1/leveldb_base/build$

```

于是, 继续往下看函数的具体内容, 发现其实是最下面的第二行:

```
ForEachOverlapping(state.saver.user_key, state.ikey, &state, &State::Match);

return state.found ? state.s : Status::NotFound(Slice());
```

其中的ForEachOverlapping中Match调用 `state->s = state->vset->table_cache->Get(*state->options, f->number, f->file_size, state->ikey, &state->saver, SaveValue);` 里 db/table_cache.cc的GET, 然后在table_cache.cc的GET里 `s = t->InternalGet(options, k, arg, handle_result);` 调用table/table.cc的 `InternalGet`。

于是, 选择在table/table.cc的 `InternalGet` 中添加具体TTL过期判断逻辑:

为 `InternalGet` 函数增加 TTL 功能的主要思路也就是通过在数据读取过程中解析并检查过期时间, 在数据已过期时返回“Not Found”状态, 而在数据未过期时返回实际值。

1. 获取存储的值和过期时间

在块迭代器 `block_iter` 中进行键查找, 通过 `seek(k)` 定位键, 并从中获取键对应的值。由于存储格式中将过期时间和原始值拼接为组合字符串, 接下来的任务是将其解包并分离出原始值和过期时间:

- 调用 `combined_value.ToString()` 将 `slice` 转换为 `std::string` 类型, 便于后续分离和解析。
- 假设存储格式的最后 19 个字符为过期时间, 使用 `combined_str.substr(combined_str.size() - 19, 19)` 获取 `expiration_time_str`, 并通过 `combined_str.substr(0, combined_str.size() - 19)` 提取实际存储值 `actual_value`。

```
if (block_iter->valid()) {
    // 这里获取存储的组合字符串
    s = Status::OK(); // 确保状态为 OK
    Slice combined_value = block_iter->value();

    // 获取实际的值和过期时间
    std::string combined_str = combined_value.ToString();

    // 假设过期时间是字符串的最后19个字符
    std::string expiration_time_str = combined_str.substr(combined_str.size() -
19, 19); // 获取过期时间字符串
    std::string actual_value = combined_str.substr(0, combined_str.size() - 19);
    // 获取实际值
```

2. 解析过期时间并与当前时间进行比较

过期时间从字符串格式解析为时间戳, 以便与当前时间进行比较:

- 使用 `std::istringstream` 和 `std::get_time` 解析 `expiration_time_str` 为 `tm` 结构, 然后通过 `std::mktime` 转换为时间戳 `expiration_time`。
- 获取系统当前时间戳 `now_time_t`, 并与 `expiration_time` 进行比较, 判断是否过期。

```

// 解析过期时间为时间戳
std::tm tm = {};
std::istringstream ss(expiration_time_str);
ss >> std::get_time(&tm, "%Y-%m-%d %H:%M:%S");
std::time_t expiration_time = std::mktime(&tm);

// 获取当前时间并与过期时间进行比较
auto now = std::chrono::system_clock::now();
auto now_time_t = std::chrono::system_clock::to_time_t(now);

```

3. 检查过期状态并返回数据

根据过期检查结果返回状态和数据:

- 如果 `expiration_time` 大于当前时间 `now_time_t`, 说明数据未过期, 将 `actual_value` 传递给 `handle_result` 回调函数, 作为返回值。
- 如果数据已过期, 设置 `Status` 为 `NotFound`, 并提供“Key has expired”信息。

```

// 检查是否过期
    if (expiration_time > now_time_t) {
        // 调用结果处理函数, 返回实际值
        (*handle_result)(arg, block_iter->key(), slice(actual_value));//,
slice(expiration_time_str));
        s = block_iter->status();
        if (! s.ok()) printf("291\n");
    } else {
        // 数据已过期, 处理过期情况
        s = Status::NotFound("Key has expired");
    }

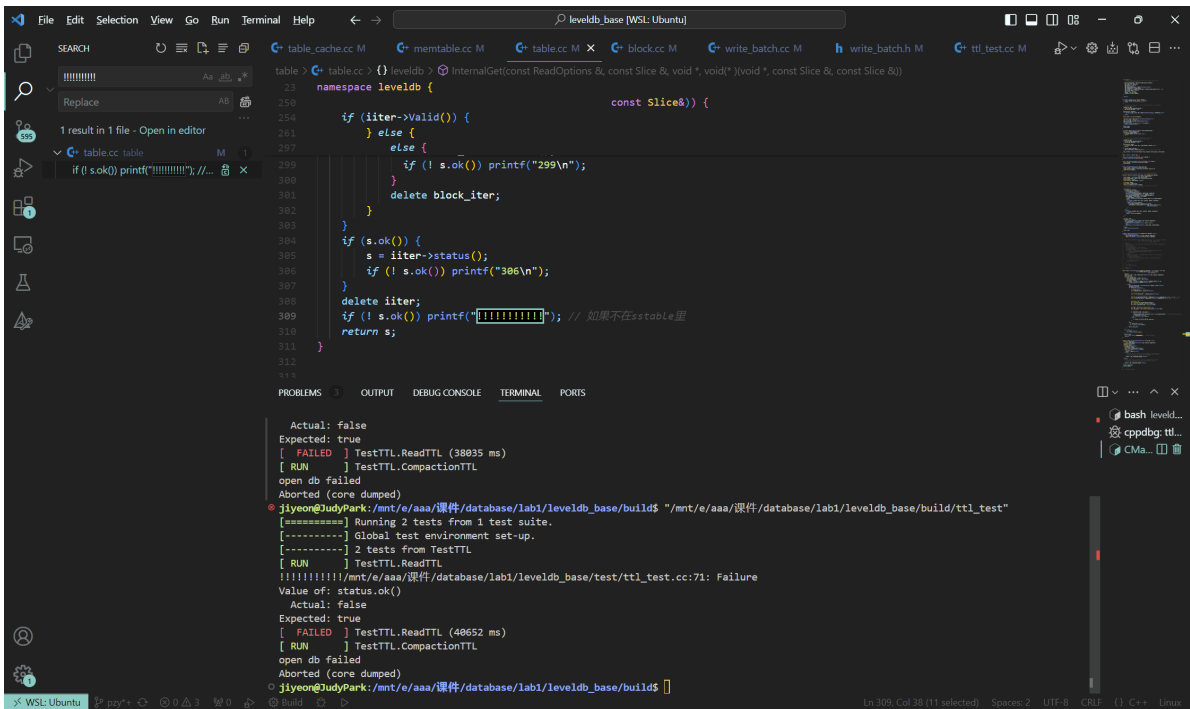
        else {
            s = block_iter->status();
            if (! s.ok()) printf("299\n");
        }
        delete block_iter;
    }
}
if (s.ok()) {
    s = iiter->status();
    if (! s.ok()) printf("306\n");
}
delete iiter;
if (! s.ok()) printf("!!!!!!!!!!!!"); // 如果不在sstable里
return s;

```

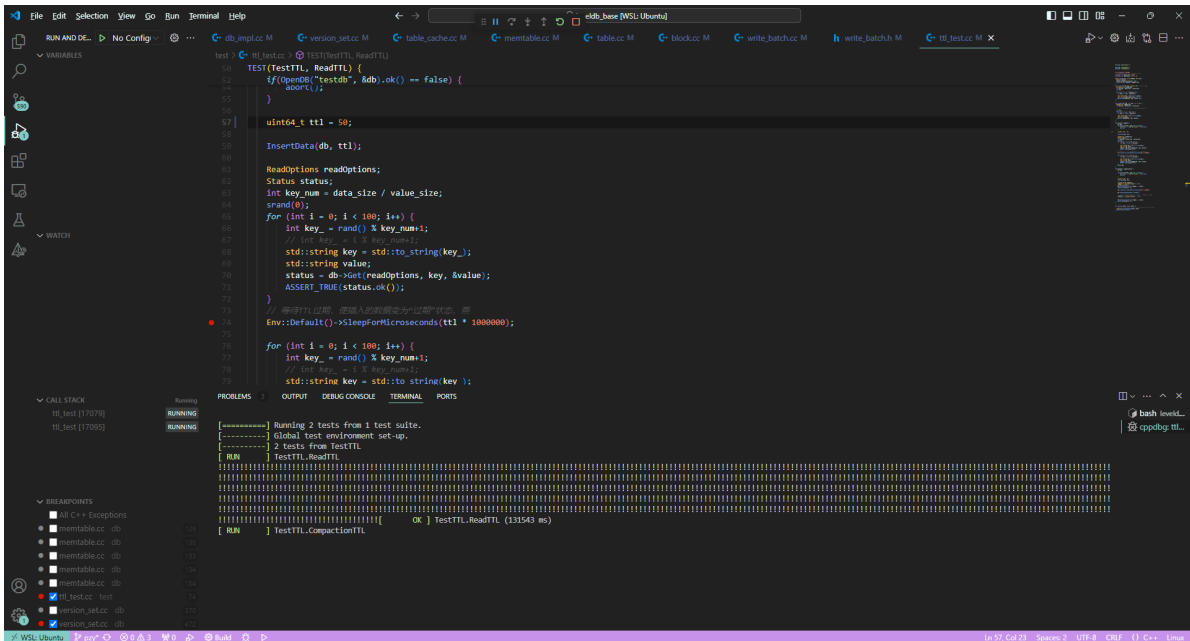
4. 删除迭代器并返回状态

在完成数据查找和过期检查后, 释放迭代器资源 `iiter` 和 `block_iter`。

代码逻辑没有问题, 但是TEST跑不通?



打断点，添加输出查找是哪一行的问题，一开始愣是没找出来。最后发现是 `(*handle_result)(arg, block_iter->key(), slice(actual_value))`; 这里的参数问题，第二个参数没有传对导致的。修正后，运行成功：



那为什么没有在其他调用层添加？

理论上应该都是可以添加的，但是因为其他层要么就是数据结构不对应，要么就是会在调用回来的时候出现问题，所以没有在其他层添加。在最底层seek中应该也可以实现，但是由于时间上优化应该没有那么大，而且逻辑上容易乱，所以还是选择table/table.cc的 `InternalGet` 来实现TTL功能。

4.3 修改合并流程--Compaction

在阅读合并部分的代码时，我们注意到不同种类的合并都需要调用 `docompactionword` 函数，且这个函数对要合并的sstable中的每个键值对都进行了遍历。除此之外更为重要的是，这个函数中的 `drop` 变量表示键值对是否保留，若 `drop` 被设为 `false`，当前键值对就不会被插入到合并后的新sstable中去，这与我们删除过期键值对的需求不谋而合。

因此，我们在给drop变量定值的条件语句中新增了一个判断条件，即对value进行解析，若时间戳超过当前时间，就将drop设为false。

```
// Handle key/value, add to state, etc.
bool drop = false;
if (!ParseInternalKey(key, &ikey)) {
    // Do not hide error keys
    current_user_key.clear();
    has_current_user_key = false;
    last_sequence_for_key = kMaxSequenceNumber;
} else {
    if (!has_current_user_key ||
        user_comparator()->Compare(ikey.user_key, slice(current_user_key)) !=
            0) {
        // First occurrence of this user key
        current_user_key.assign(ikey.user_key.data(), ikey.user_key.size());
        has_current_user_key = true;
        last_sequence_for_key = kMaxSequenceNumber;
    }

    if (last_sequence_for_key <= compact->smallest_snapshot) {
        // Hidden by a newer entry for same user key
        drop = true; // (A)
    } else if (ikey.type == kTypeDeletion &&
                ikey.sequence <= compact->smallest_snapshot &&
                compact->compaction->IsBaseLevelForKey(ikey.user_key)) {
        // For this user key:
        // (1) there is no data in higher levels
        // (2) data in lower levels will have larger sequence numbers
        // (3) data in layers that are being compacted here and have
        //     smaller sequence numbers will be dropped in the next
        //     few iterations of this loop (by rule (A) above).
        // Therefore this deletion marker is obsolete and can be dropped.
        drop = true;
    }
    else{
        slice value = input->value();
        std::string combined_str = value.ToString();
        // 从右往左提取固定长度的字符串
        std::string time_part = combined_str.substr(combined_str.length() - 19,
19);

        // 解析时间字符串
        std::tm parsed_tm = {};
        const char* result = strptime(time_part.c_str(), "%Y-%m-%d %H:%M:%S",
&parsed_tm);

        // 将解析出的时间转为 time_t
        std::time_t parsed_time_t = std::mktime(&parsed_tm);
        // 获取当前时间
        auto now = std::chrono::system_clock::now();

        // 转换为 time_t
        std::time_t now_time_t = std::chrono::system_clock::to_time_t(now);
```

```

    if (parsed_time_t <= now_time_t) drop = true; // 如果过期了就丢弃
    // 心
}

```

增加上述代码后，我们本以为已经完成了合并时删除过期键值对的逻辑，但多次尝试仍无法通过测试点。

按照测试逻辑，在手动触发全数据库的合并后，数据库中应该不存在有效的sstable文件，因此两次调用 `versions->ApproximateOffsetOf()` 函数返回的值都应为0；通过打断点观察在第二次调用 `db->GetApproximateSizes(ranges, 1, sizes)` 时 `versions->ApproximateOffsetOf(v, 'A')` 中 `result` 增加的时机，我们发现结果具有高度重复性，即第二层中总存在一个文件没有被删除，在扫描到第二层第0个文件时，`icmp_.Compare(files[i]->largest, ikey) <= 0` 恒成立，因此返回该文件大小。

通过检查 `DBImpl::CompactRange()` 函数执行的逻辑，我们发现程序首先通过执行 `leveldb::Version::OverlapInLevel()` 找到与目标合并范围有重叠的最大层，记为 `max_level_with_files`，然后依次对小于最大重叠层的层进行合并。

我们的数据在sstable中最高存储在第二层，在对第一层的各个sstable进行合并时，会找到与第二层的数据范围重叠的文件，但部分第二层的文件可能并不会与第一层的文件有范围重叠，因此在合并阶段其中的键值对没有机会被扫描到，从而导致即使进行了全局的手动合并，也无法对所有过期数据进行删除。

鉴于以上分析，我们认为只需将合并进行到 `max_level_with_files` 层就可以解决这一问题，经实验验证后均通过测试，下面是修改后的代码：

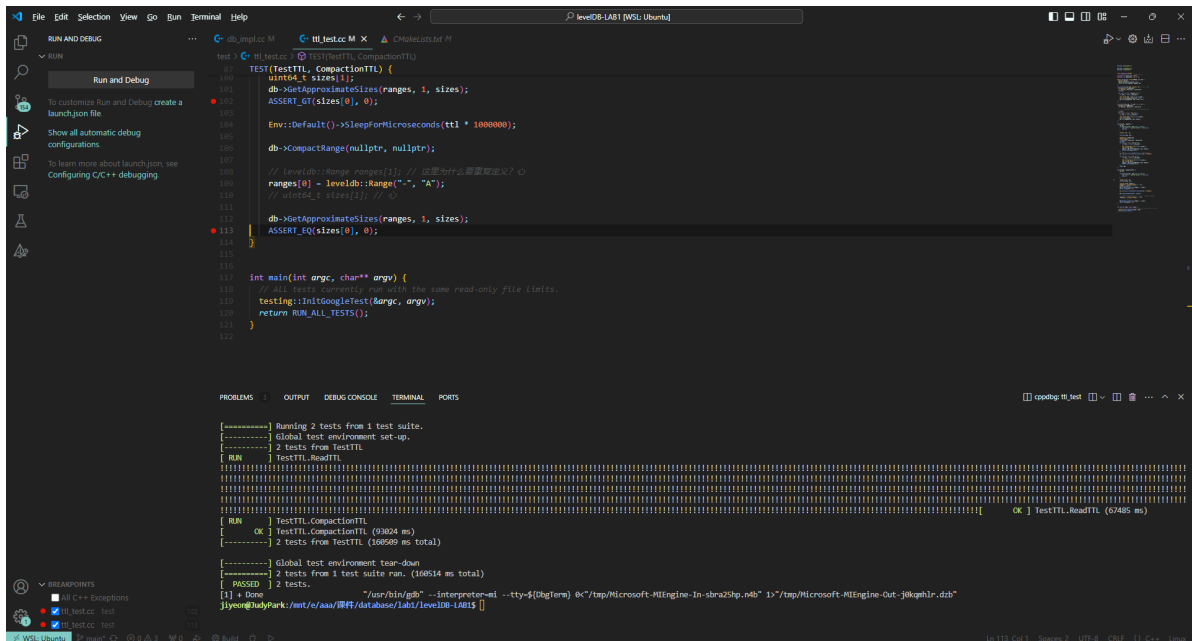
```

void DBImpl::CompactRange(const Slice* begin, const Slice* end) {
    int max_level_with_files = 1;
    {
        MutexLock l(&mutex_);
        Version* base = versions_->current();
        for (int level = 1; level < config::kNumLevels; level++) {
            if (base->OverlapInLevel(level, begin, end)) {
                max_level_with_files = level;
            }
        }
    }
    TEST_CompactMemTable(); // TODO(sanjay): skip if memtable does not overlap

    for (int level = 0; level <= max_level_with_files; level++) { // 此处将小于号修改
        TEST_CompactRange(level, begin, end);
    }
}

```

最后GET和Compaction的逻辑合并后运行，成功执行：



五、测试用例

初始时ttl的值为20，组员在增加了一些读和合并时处理ttl的操作后，发现无法通过TEST(TestTTL, ReadTTL)中第一部分的断言：

```

for (int i = 0; i < 100; i++) {
    int key_ = rand() % key_num+1;
    // int key_ = i % key_num+1;
    std::string key = std::to_string(key_);
    std::string value;
    status = db->Get(readOptions, key, &value);
    ASSERT_TRUE(status.ok());
}

```

考虑到这一步的检查与数据过期无关，且在修改代码前这一部分可以正确通过，因此我们推测是检查ttl是否到期的操作使得该测试循环无法在20秒内完成，从而违背了测试的逻辑（假设所有数据都没有过期，都可以被读到）。

我们尝试了使用较大的ttl(ttl=30,50,60),发现ttl较小（20、30）时，通过测试的概率不能保证，但取较大值（50、60）时，多次测试均能通过。

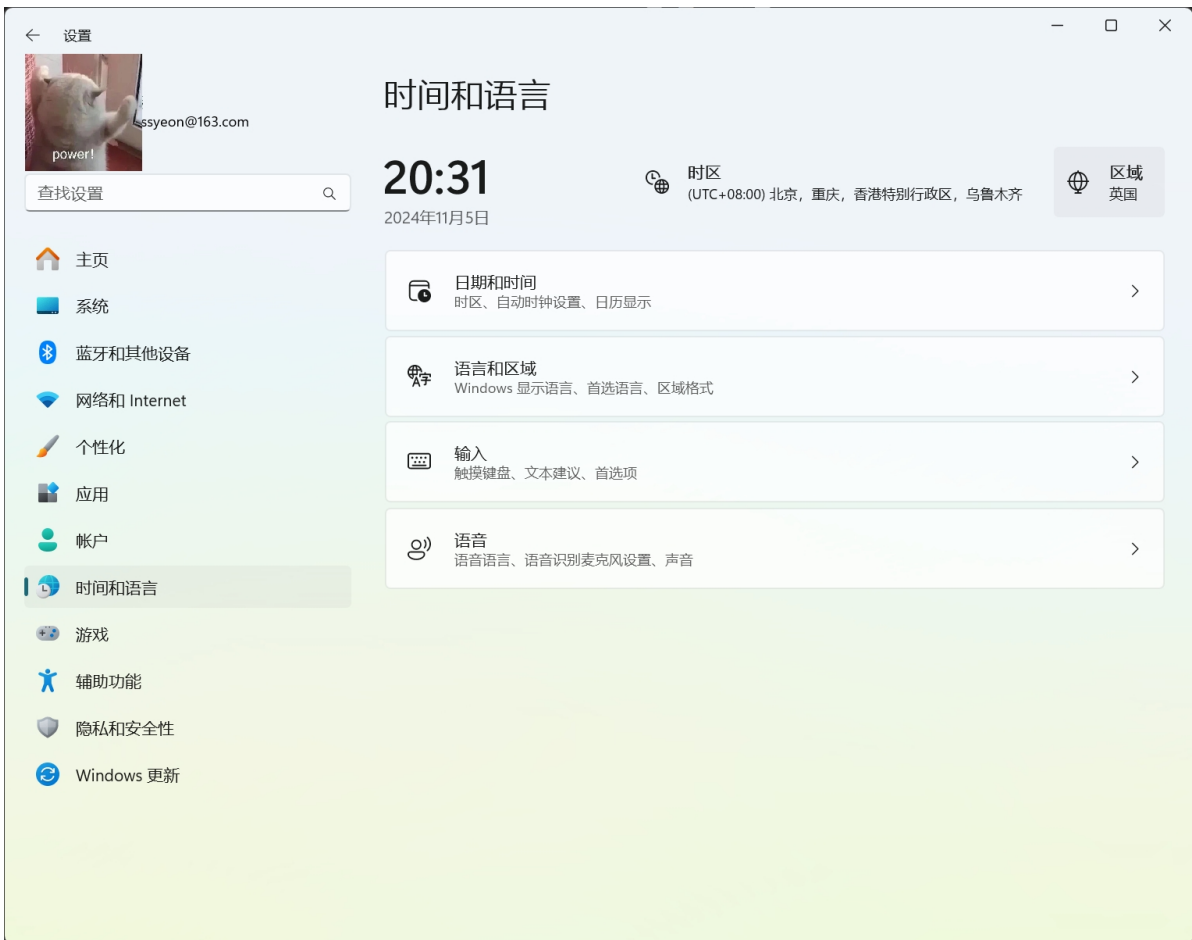
六、特点/出现的问题及解决

6.1 系统时区和地区

一位组员在进行测试时，出现已经超过ttl的键值对没有被正确删除。经过打断点调试发现获取的当前时间是正确的，但是存入的过期时间总是会比正确时间晚八小时（得益于我们将std::tm类型转化为可读的格式化字符串类型）

经过一番排查后该组员想起来曾经将电脑系统中的时区设置为英国时区，英国时区正好比北京时间晚大约八小时。修改系统时区后程序果然按预期运行了。

```
namespace leveldb {
    const Slice& {
        if (iiter->Valid()) {
        } else {
            if (block_iter->Valid()) {
                // 假设过期时间是字符串的最后19个字符
                std::string expiration_time_str = combined_str.substr(combined_str.size() - 19, 19); // 获取过期时间字符串
                std::string actual_value = combined_str.substr(0, combined_str.size() - 19); // 获取实际值
                // 解析过期时间为时间戳
                std::tm tm = {};
                std::stringstream ss(expiration_time_str);
                ss >> std::get_time(&tm, "%Y-%m-%d %H:%M:%S");
                std::time_t expiration_time = std::mktime(&tm);
                // 获取当前时间与过期时间进行比较
                auto now = std::chrono::system_clock::now(); std::chrono::sys_time = { 1730816595706549965n }; [2024-11-05 14:23:15]
                auto now_time_t = std::chrono::system_clock::to_time_t(now);
                // 检查是否过期
                if (expiration_time > now_time_t) {
                    // 调用结果处理函数, 返回实际值
                    (*handle_result)(arg, block_iter->key(), slice(actual_value)); // slice(expiration_time_str);
                }
            }
        }
    }
};
```



从这个小插曲中，我们意识到获取的时间戳在拓展到不同地区、不同系统上可能出现潜在的问题。因此，我们希望获得一种可以不受系统时间影响的、统一的时间戳。

c++中获取时间的方式有：

```
auto now = std::chrono::system_clock::now();
// 系统时钟 (system_clock)：表示当前系统范围的实时日历时钟，通常与系统的钟同步。

auto now = std::chrono::steady_clock::now();
// 稳定时钟 (steady_clock)：表示一个单调递增的时钟，不受系统时间的调整或修改影响。它可以
用于测量时间间隔，以及实现定时器和延时等功能。

std::time_t now = std::time(0);
// 其值表示从CUT (Coordinated Universal Time) 时间1970年1月1日00:00:00 (称为UNIX系
统的Epoch时间) 到当前时刻的秒数
```

虽然上述方法可以返回不依赖于系统设置的时间，但是为了将时间数据格式化为可读的格式，需要使用 `strftime(buffer, sizeof(buffer), "%Y-%m-%d %H:%M:%S", timeptr)` 函数，该函数接受的最后一个参数为 `std::tm` 类型，而将 `std::time_t` 类型的时间转化为 `std::tm` 类型则可以使用依赖于系统时间的 `std::localtime` 函数或使用UTC时间的 `std::gmtime` 函数。

因此，解决方案有使用 `std::gmtime()` 配合不受系统影响的秒数获取函数或 直接使用时间获取函数并记录秒数的位数。

6.2 GET时选择调用层

由于GET是一层层调用，一层层返回参数和状态的，所以具体选择在哪个函数的时候我们小组有进行一些讨论。再考虑性能的同时考虑代码可实现性，如果在最底层 `seek` 里实现，那就是等找到数据后，过期数据就直接认为 `NotFound`，一层层返回；如果是在上层其他层实现，`seek` 满心欢喜认为找到了某个数据，但在返回出来时如果过期突然被一棒子打死会不会导致性能不高？我们的最终讨论结果是：影响不大。因为就算是在 `seek` 里实现，它还是要去找数据，归根结底“找”这个过程是不可避免的。所以，就考虑代码可实现性，就选择在 `table/table.cc` 的 `InternalGet` 来实现TTL功能。

七、总结

在本次实验中，我们通过在LevelDB数据库中加入TTL (Time To Live) 功能，使其具备对过期数据的自动管理能力。这一改进在数据管理上带来了更高的效率，使数据库能够在存储和访问上更加轻量化。设计过程中，我们基于现有的类结构进行改动，尽量不引入额外的成员变量或复杂的函数修改，以保持代码结构的简洁性。同时，为实现TTL，我们是在 `MemTable` 和 `SSTable` 的读写操作中增加了过期检查，并在 `Version::Get` 函数中引入了TTL支持，实现了在不更改函数参数和函数体的情况下对 `seek` 操作中的数据进行有效过滤，确保过期数据不被错误地返回。

通过阶段性测试，我们验证了TTL功能在不同存储层次中能够有效生效，并最终解决了数据在 `MemTable` 和 `SSTable` 之间可能存在的失效问题。总体来说，这次实验不仅提升了我们对数据库底层原理的理解，还在解决复杂功能的代码实现和性能优化上获得了宝贵经验。