
四 数据库面试题总结

4.1 MySQL

作者：Guide哥。

介绍: Github 70k Star 项目 [JavaGuide](#) (公众号同名) 作者。每周都会在公众号更新一些自己原创干货。公众号后台回复“1”领取Java工程师必备学习资料+面试突击pdf。

4.1.1 精品推荐

书籍推荐

- 《SQL基础教程（第2版）》（入门级）
- 《高性能MySQL：第3版》（进阶）

文字教程推荐

- [SQL Tutorial](#) (SQL语句学习,英文)、[SQL Tutorial](#) (SQL语句学习,中文)、[SQL语句在线练习](#) (非常不错)
- [Github-MySQL入门教程 \(MySQL tutorial book\)](#) (从零开始学习MySQL, 主要是面向MySQL数据库管理系统初学者)
- [官方教程](#)
- [MySQL 教程 \(菜鸟教程\)](#)

相关资源推荐

- [中国5级行政区域mysql库](#)

视频教程推荐

基础入门: [与MySQL的零距离接触-慕课网](#)

MySQL开发技巧: [MySQL开发技巧 \(一\)](#) [MySQL开发技巧 \(二\)](#) [MySQL开发技巧 \(三\)](#)

MySQL5.7新特性及相关优化技巧: [MySQL5.7版本新特性](#) [性能优化之MySQL优化](#)

[MySQL集群 \(PXC\) 入门](#) [MyCAT入门及应用](#)

常见问题总结

4.1.2 什么是MySQL?

MySQL 是一种关系型数据库, 在Java企业级开发中非常常用, 因为 MySQL 是开源免费的, 并且方便扩展。阿里巴巴数据库系统也大量用到了 MySQL, 因此它的稳定性是有保障的。MySQL是开放源代码的, 因此任何人都可以在 GPL(General Public License) 的许可下下载并根据个性化的需要对其进行修改。MySQL的默认端口号是3306。

4.1.3 存储引擎

一些常用命令

查看MySQL提供的所有存储引擎

```
1 | mysql> show engines;
```

MySQL 5.7 Command Line Client

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

```
mysql> show engines;
```

Engine	Support	Comment	Transactions	XA	Savepoints
InnoDB	DEFAULT	Supports transactions, row-level locking, and foreign keys	YES	YES	YES
MRG_MYISAM	YES	Collection of identical MyISAM tables	NO	NO	NO
MEMORY	YES	Hash based, stored in memory, useful for temporary tables	NO	NO	NO
BLACKHOLE	YES	/dev/null storage engine (anything you write to it disappears)	NO	NO	NO
MyISAM	YES	MyISAM storage engine	NO	NO	NO
CSV	YES	CSV storage engine	NO	NO	NO
ARCHIVE	YES	Archive storage engine	NO	NO	NO
PERFORMANCE_SCHEMA	YES	Performance Schema	NO	NO	NO
FEDERATED	NO	Federated MySQL storage engine	NULL	NULL	NULL

从上图我们可以查看出 MySQL 当前默认的存储引擎是InnoDB,并且在5.7版本所有的存储引擎中只有InnoDB 是事务性存储引擎，也就是说只有 InnoDB 支持事务。

查看MySQL当前默认的存储引擎

我们也可以通过下面的命令查看默认的存储引擎。

```
1 | mysql> show variables like '%storage_engine';
```

查看表的存储引擎

```
1 | show table status like "table_name" ;
```

```
mysql> show table status like "tb_base";
```

Name	Engine	Version	Row format	Rows	Avg_row_length	Data_length	Max_data_length	Index_length	Data_free	Auto_increment	Create_time
Update_time	Check_time	Collation	Checksum	Create_options	Comment						
tb_base	InnoDB	10	Dynamic	0	0	16384	0	0	0	2	2019-05-05 15:58:24
NULL	NULL		utf8_general_ci	NULL							

MyISAM和InnoDB区别

MyISAM是MySQL的默认数据库引擎（5.5版之前）。虽然性能极佳，而且提供了大量的特性，包括全文索引、压缩、空间函数等，但MyISAM不支持事务和行级锁，而且最大的缺陷就是崩溃后无法安全恢复。不过，5.5版本之后，MySQL引入了InnoDB（事务性数据库引擎），MySQL 5.5版本后默认的存储引擎为InnoDB。

大多数时候我们使用的都是 InnoDB 存储引擎，但是在某些情况下使用 MyISAM 也是合适的比如读密集的情况下。（如果你不介意 MyISAM 崩溃恢复问题的话）。

两者的对比：

- 是否支持行级锁**：MyISAM 只有表级锁(table-level locking)，而InnoDB 支持行级锁(row-level locking)和表级锁,默认为行级锁。
- 是否支持事务和崩溃后的安全恢复**：**MyISAM** 强调的是性能，每次查询具有原子性,其执行速度比 InnoDB类型更快，但是不提供事务支持。但是**InnoDB** 提供事务支持事务，外部键等高级数据库功能。具有事务(commit)、回滚(rollback)和崩溃修复能力(crash recovery capabilities)的事务安全(transaction-safe (ACID compliant))型表。
- 是否支持外键**：MyISAM不支持，而InnoDB支持。
- 是否支持MVCC**：仅 InnoDB 支持。应对高并发事务, MVCC比单纯的加锁更高效;MVCC只在 `READ COMMITTED` 和 `REPEATABLE READ` 两个隔离级别下工作;MVCC可以使用 乐观(optimistic)锁 和 悲

观(pessimistic)锁来实现;各数据库中MVCC实现并不统一。推荐阅读: [MySQL-InnoDB-MVCC多版本并发控制](#)

5.

《MySQL高性能》上面有一句话这样写到:

不要轻易相信“MyISAM比InnoDB快”之类的经验之谈, 这个结论往往不是绝对的。在很多我们已知场景中, InnoDB的速度都可以让MyISAM望尘莫及, 尤其是用到了聚簇索引, 或者需要访问的数据都可以放入内存的应用。

一般情况下我们选择 InnoDB 都是没有问题的, 但是某些情况下你并不在乎可扩展能力和并发能力, 也不需要事务支持, 也不在乎崩溃后的安全恢复问题的话, 选择MyISAM也是一个不错的选择。但是一般情况下, 我们都是需要考虑到这些问题的。

4.1.4 字符集及校对规则

字符集指的是一种从二进制编码到某类字符符号的映射。校对规则则是指某种字符集下的排序规则。MySQL中每一种字符集都会对应一系列的校对规则。

MySQL采用的是类似继承的方式指定字符集的默认值, 每个数据库以及每张数据表都有自己的默认值, 他们逐层继承。比如: 某个库中所有表的默认字符集将是该数据库所指定的字符集(这些表在没有指定字符集的情况下, 才会采用默认字符集) PS: 整理自《Java工程师修炼之道》

详细内容可以参考: [MySQL字符集及校对规则的理解](#)

作者: Guide哥。

介绍: Github 70k Star 项目 [JavaGuide](#) (公众号同名) 作者。每周都会在公众号更新一些自己原创干货。公众号后台回复“1”领取Java工程师必备学习资料+面试突击pdf。

4.1.5 索引

MySQL索引使用的数据结构主要有**BTree索引**和**哈希索引**。对于哈希索引来说, 底层的数据结构就是哈希表, 因此在绝大多数需求为单条记录查询的时候, 可以选择哈希索引, 查询性能最快; 其余大部分场景, 建议选择BTree索引。

MySQL的BTree索引使用的是B树中的B+Tree, 但对于主要的两种存储引擎的实现方式是不同的。

- **MyISAM:** B+Tree叶节点的data域存放的是数据记录的地址。在索引检索的时候, 首先按照B+Tree搜索算法搜索索引, 如果指定的Key存在, 则取出其 data 域的值, 然后以 data 域的值作为地址读取相应的数据记录。这被称为“非聚簇索引”。
- **InnoDB:** 其数据文件本身就是索引文件。相比MyISAM, 索引文件和数据文件是分离的, 其表数据文件本身就是按B+Tree组织的一个索引结构, 树的叶节点data域保存了完整的数据记录。这个索引的key是数据表的主键, 因此InnoDB表数据文件本身就是主索引。这被称为“聚簇索引(或聚集索引)”。而其余的索引都作为辅助索引, 辅助索引的data域存储相应记录主键的值而不是地址, 这也是和MyISAM不同的地方。**在根据主索引搜索时, 直接找到key所在的节点即可取出数据; 在根据辅助索引查找时, 则需要先取出主键的值, 再走一遍主索引。因此, 在设计表的时候, 不建议使用过长的字段作为主键, 也不建议使用非单调的字段作为主键, 这样会造成主索引频繁分裂。** PS: 整理自《Java工程师修炼之道》

更多关于索引的内容可以查看文档首页MySQL目录下关于索引的详细总结。

4.1.6 查询缓存的使用

执行查询语句的时候, 会先查询缓存。不过, MySQL 8.0 版本后移除, 因为这个功能不太实用

my.cnf加入以下配置, 重启MySQL开启查询缓存

```
1 query_cache_type=1
2 query_cache_size=600000
```

MySQL执行以下命令也可以开启查询缓存

```
1 set global query_cache_type=1;
2 set global query_cache_size=600000;
```

如上，开启查询缓存后在同样的查询条件以及数据情况下，会直接在缓存中返回结果。这里的查询条件包括查询本身、当前要查询的数据库、客户端协议版本号等一些可能影响结果的信息。因此任何两个查询在任何字符上的不同都会导致缓存不命中。此外，如果查询中包含任何用户自定义函数、存储函数、用户变量、临时表、MySQL库中的系统表，其查询结果也不会被缓存。

缓存建立之后，MySQL的查询缓存系统会跟踪查询中涉及的每张表，如果这些表（数据或结构）发生变化，那么和这张表相关的所有缓存数据都将失效。

缓存虽然能够提升数据库的查询性能，但是缓存同时也带来了额外的开销，每次查询后都要做一次缓存操作，失效后还要销毁。因此，开启缓存查询要谨慎，尤其对于写密集的应用来说更是如此。如果开启，要注意合理控制缓存空间大小，一般来说其大小设置为几十MB比较合适。此外，还可以通过 `sql_cache` 和 `sql_no_cache` 来控制某个查询语句是否需要缓存：

```
1 select sql_no_cache count(*) from usr;
```

4.1.7 什么是事务？

事务是逻辑上的一组操作，要么都执行，要么都不执行。

事务最经典也经常被拿出来例子就是转账了。假如小明要给小红转账1000元，这个转账会涉及到两个关键操作就是：将小明的余额减少1000元，将小红的余额增加1000元。万一在这两个操作之间突然出现错误比如银行系统崩溃，导致小明余额减少而小红的余额没有增加，这样就不对了。事务就是保证这两个关键操作要么都成功，要么都要失败。

4.1.8 事物的四大特性(ACID)



1. **原子性 (Atomicity)**：事务是最小的执行单位，不允许分割。事务的原子性确保动作要么全部完成，要么完全不起作用；
2. **一致性 (Consistency)**：执行事务前后，数据保持一致，多个事务对同一个数据读取的结果是相同的；
3. **隔离性 (Isolation)**：并发访问数据库时，一个用户的事务不被其他事务所干扰，各并发事务之间数据库是独立的；
4. **持久性 (Durability)**：一个事务被提交之后，它对数据库中数据的改变是持久的，即使数据库发生故障也不应该对其有任何影响。

4.1.9 并发事务带来哪些问题？

在典型的应用程序中，多个事务并发运行，经常会操作相同的数据来完成各自的任务（多个用户对同一数据进行操作）。并发虽然是必须的，但可能会导致以下的问题。

- **脏读 (Dirty read)**：当一个事务正在访问数据并且对数据进行了修改，而这种修改还没有提交到数据库中，这时另外一个事务也访问了这个数据，然后使用了这个数据。因为这个数据是还没有提交的数据，那么另外一个事务读到的这个数据是“脏数据”，依据“脏数据”所做的操作可能是不正确的。
- **丢失修改 (Lost to modify)**：指在一个事务读取一个数据时，另外一个事务也访问了该数据，那么在第一个事务中修改了这个数据后，第二个事务也修改了这个数据。这样第一个事务内的修改结果就被丢失，因此称为丢失修改。例如：事务1读取某表中的数据A=20，事务2也读取A=20，事务1修改A=A-1，事务2也修改A=A-1，最终结果A=19，事务1的修改被丢失。
- **不可重复读 (Unrepeatableread)**：指在一个事务内多次读同一数据。在这个事务还没有结束时，另一个事务也访问该数据。那么，在第一个事务中的两次读数据之间，由于第二个事务的修改导致第一个事务两次读取的数据可能不太一样。这就发生了在一个事务内两次读到的数据是不一样的情况，因此称为不可重复读。
- **幻读 (Phantom read)**：幻读与不可重复读类似。它发生在一个事务 (T1) 读取了几行数据，接着另一个并发事务 (T2) 插入了一些数据时。在随后的查询中，第一个事务 (T1) 就会发现多了一些原本不存在的记录，就好像发生了幻觉一样，所以称为幻读。

不可重复读和幻读区别：

不可重复读的重点是修改比如多次读取一条记录发现其中某些列的值被修改，幻读的重点在于新增或者删除比如多次读取一条记录发现记录增多或减少了。

4.1.10 事务隔离级别有哪些？MySQL的默认隔离级别是？

SQL 标准定义了四个隔离级别：

- **READ-UNCOMMITTED(读取未提交)**：最低的隔离级别，允许读取尚未提交的数据变更，**可能会导致脏读、幻读或不可重复读。**
- **READ-COMMITTED(读取已提交)**：允许读取并发事务已经提交的数据，**可以阻止脏读，但是幻读或不可重复读仍有可能发生。**
- **REPEATABLE-READ(可重复读)**：对同一字段的多次读取结果都是一致的，除非数据是被本身事务自己所修改，**可以阻止脏读和不可重复读，但幻读仍有可能发生。**
- **SERIALIZABLE(可串行化)**：最高的隔离级别，完全服从ACID的隔离级别。所有的事务依次逐个执行，这样事务之间就完全不可能产生干扰，也就是说，**该级别可以防止脏读、不可重复读以及幻读。**

隔离级别	脏读	不可重复读	幻影读
READ-UNCOMMITTED	√	√	√
READ-COMMITTED	×	√	√
REPEATABLE-READ	×	×	√
SERIALIZABLE	×	×	×

MySQL InnoDB 存储引擎的默认支持的隔离级别是 **REPEATABLE-READ (可重读)**。我们可以通过 `SELECT @@tx_isolation;` 命令来查看


```

1  mysql> SELECT @@tx_isolation;
2  +-----+
3  | @@tx_isolation |
4  +-----+
5  | REPEATABLE-READ |
6  +-----+

```

这里需要注意的是：与 SQL 标准不同的地方在于 InnoDB 存储引擎在 **REPEATABLE-READ (可重读)** 事务隔离级别下使用的是 Next-Key Lock 锁算法，因此可以避免幻读的产生，这与其他数据库系统(如 SQL Server)

是不同的。所以说 InnoDB 存储引擎的默认支持的隔离级别是 **REPEATABLE-READ (可重读)** 已经可以完全保证事务的隔离性要求，即达到了

SQL 标准的 **SERIALIZABLE(可串行化)** 隔离级别。因为隔离级别越低，事务请求的锁越少，所以大部分数据库系统的隔离级别都是 **READ-COMMITTED(读取提交内容)**，但是你要知道的是 InnoDB 存储引擎默认使用 **REPEATABLE-READ (可重读)** 并不会有任何性能损失。

InnoDB 存储引擎在 **分布式事务** 的情况下一般会用到 **SERIALIZABLE(可串行化)** 隔离级别。

4.1.11 锁机制与InnoDB锁算法

MyISAM和InnoDB存储引擎使用的锁：

- MyISAM采用表级锁(table-level locking)。
- InnoDB支持行级锁(row-level locking)和表级锁,默认为行级锁

表级锁和行级锁对比：

- **表级锁：**MySQL中锁定 **粒度最大** 的一种锁，对当前操作的整张表加锁，实现简单，资源消耗也比较少，加锁快，不会出现死锁。其锁定粒度最大，触发锁冲突的概率最高，并发度最低，MyISAM和 InnoDB引擎都支持表级锁。
- **行级锁：**MySQL中锁定 **粒度最小** 的一种锁，只针对当前操作的行进行加锁。行级锁能大大减少数据库操作的冲突。其加锁粒度最小，并发度高，但加锁的开销也最大，加锁慢，会出现死锁。

详细内容可以参考：MySQL锁机制简单了解一下：

https://blog.csdn.net/qg_34337272/article/details/80611486

InnoDB存储引擎的锁的算法有三种：

- Record lock：单个行记录上的锁
- Gap lock：间隙锁，锁定一个范围，不包括记录本身
- Next-key lock：record+gap 锁定一个范围，包含记录本身

相关知识点：

1. innodb对于行的查询使用next-key lock
2. Next-locking keying为了解决Phantom Problem幻读问题
3. 当查询的索引含有唯一属性时，将next-key lock降级为record key
4. Gap锁设计的目的是为了阻止多个事务将记录插入到同一范围内，而这会导致幻读问题的产生
5. 有两种方式显式关闭gap锁：（除了外键约束和唯一性检查外，其余情况仅使用record lock）A. 将事务隔离级别设置为RC B. 将参数innodb_locks_unsafe_for_binlog设置为1

4.1.12 大表优化

当MySQL单表记录数过大时，数据库的CRUD性能会明显下降，一些常见的优化措施如下：

限定数据的范围

务必禁止不带任何限制数据范围条件的查询语句。比如：我们当用户在查询订单历史的时候，我们可以控制在一个月的范围内；

读/写分离

经典的数据库拆分方案，主库负责写，从库负责读；

垂直分区

根据数据库里面数据表的相关性进行拆分。例如，用户表中既有用户的登录信息又有用户的基本信息，可以将用户表拆分成两个单独的表，甚至放到单独的库做分库。

简单来说垂直拆分是指数据表列的拆分，把一张列比较多的表拆分为多张表。如下图所示，这样来说大家应该就更容易理解了。

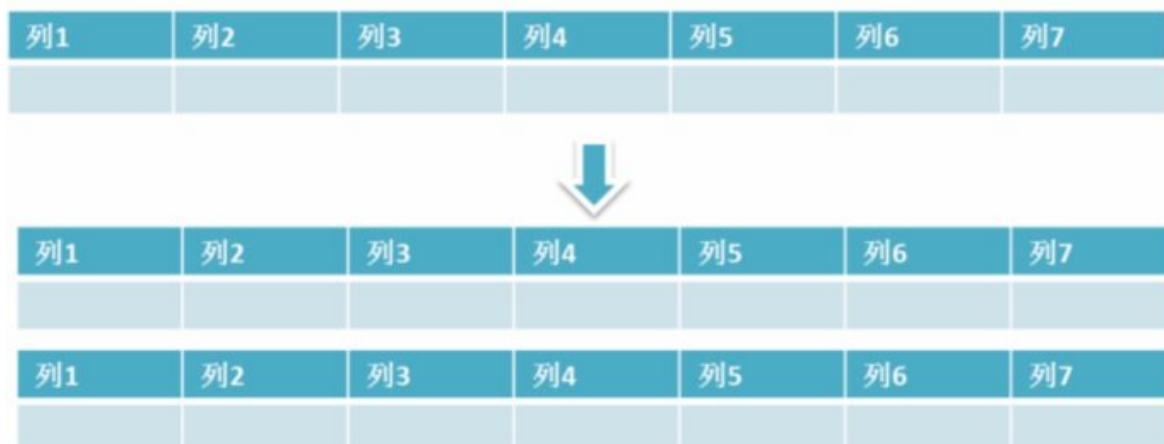


- **垂直拆分的优点：**可以使得列数据变小，在查询时减少读取的Block数，减少I/O次数。此外，垂直分区可以简化表的结构，易于维护。
- **垂直拆分的缺点：**主键会出现冗余，需要管理冗余列，并会引起Join操作，可以通过在应用层进行Join来解决。此外，垂直分区会让事务变得更加复杂；

水平分区

保持数据表结构不变，通过某种策略存储数据分片。这样每一片数据分散到不同的表或者库中，达到了分布式的目的。水平拆分可以支撑非常大的数据量。

水平拆分是指数据表行的拆分，表的行数超过200万行时，就会变慢，这时可以把一张的表的数据拆成多张表来存放。举个例子：我们可以将用户信息表拆分成多个用户信息表，这样就可以避免单一表数据量过大对性能造成影响。



水平拆分可以支持非常大的数据量。需要注意的一点是：分表仅仅是解决了单一表数据过大的问题，但由于表的数据还是在同一台机器上，其实对于提升MySQL并发能力没有什么意义，所以 **水平拆分最好分库**。

水平拆分能够支持非常大的数据量存储，应用端改造也少，但分片事务难以解决，跨节点Join性能较差，逻辑复杂。《Java工程师修炼之道》的作者推荐尽量不要对数据进行分片，因为拆分会带来逻辑、部署、运维的各种复杂度，一般的数据表在优化得当的情况下支撑千万以下的数据量是没有太大问题的。如果实在要分片，尽量选择客户端分片架构，这样可以减少一次和中间件的网络I/O。

下面补充一下数据库分片的两种常见方案：

- **客户端代理：**分片逻辑在应用端，封装在jar包中，通过修改或者封装JDBC层来实现。当当网的Sharding-JDBC、阿里的TDDL是两种比较常用的实现。
- **中间件代理：**在应用和数据中间加了一个代理层。分片逻辑统一维护在中间件服务中。我们现在谈的Mycat、360的Atlas、网易的DDB等等都是这种架构的实现。

详细内容可以参考：MySQL大表优化方案：<https://segmentfault.com/a/1190000006158186>

4.1.13 解释一下什么是池化设计思想。什么是数据库连接池？为什么需要数据库连接池？

池化设计应该不是一个新名词。我们常见的如java线程池、jdbc连接池、redis连接池等就是这类设计的代表实现。这种设计会初始预设资源，解决的问题就是抵消每次获取资源的消耗，如创建线程的开销，获取远程连接的开销等。就好比你去食堂打饭，打饭的大妈会先把饭盛好几份放那里，你来了就直接拿着饭盒加菜即可，不用再临时又盛饭又打菜，效率就高了。除了初始化资源，池化设计还包括如下这些特征：池子的初始值、池子的活跃值、池子的最大值等，这些特征可以直接映射到java线程池和数据库连接池的成员属性中。这篇文章对[池化设计思想](#)介绍的还不错，直接复制过来，避免重复造轮子了。

数据库连接本质就是一个 socket 的连接。数据库服务端还要维护一些缓存和用户权限信息之类的 所以占用了一些内存。我们可以把数据库连接池是看做是维护的数据库连接的缓存，以便将来需要对数据库的请求时可以重用这些连接。为每个用户打开和维护数据库连接，尤其是对动态数据库驱动的网站应用程序的请求，既昂贵又浪费资源。**在连接池中，创建连接后，将其放置在池中，并再次使用它，因此不必建立新的连接。如果使用了所有连接，则会建立一个新连接并将其添加到池中。**连接池还减少了用户必须等待建立与数据库的连接的时间。

4.1.14 分库分表之后,id 主键如何处理？

因为要是分成多个表之后，每个表都是从 1 开始累加，这样是不对的，我们需要一个全局唯一的 id 来支持。

生成全局 id 有下面这几种方式：

- **UUID：**不适合作为主键，因为太长了，并且无序不可读，查询效率低。比较适合用于生成唯一的名字的标示比如文件的名字。
- **数据库自增 id：**两台数据库分别设置不同步长，生成不重复ID的策略来实现高可用。这种方式生成的 id 有序，但是需要独立部署数据库实例，成本高，还会有性能瓶颈。
- **利用 redis 生成 id：**性能比较好，灵活方便，不依赖于数据库。但是，引入了新的组件造成系统更加复杂，可用性降低，编码更加复杂，增加了系统成本。
- **Twitter的snowflake算法：**Github 地址：<https://github.com/twitter-archive/snowflake>。
- **美团的Leaf分布式ID生成系统：**Leaf 是美团开源的分布式ID生成器，能保证全局唯一性、趋势递增、单调递增、信息安全，里面也提到了几种分布式方案的对比，但也需要依赖关系数据库、Zookeeper等中间件。感觉还不错。美团技术团队的一篇文章：<https://tech.meituan.com/2017/04/21/mt-leaf.html>。
-

4.1.15 一条SQL语句在MySQL中如何执行的

[一条SQL语句在MySQL中如何执行的](#)

4.1.16 MySQL高性能优化规范建议

[MySQL高性能优化规范建议](#)

4.1.17 一条SQL语句执行得很慢的原因有哪些？

[腾讯面试：一条SQL语句执行得很慢的原因有哪些？ ---不看后悔系列](#)

4.1.19 后端程序员必备：书写高质量SQL的30条建议

[后端程序员必备：书写高质量SQL的30条建议](#)