

## 分库设计中的主键选择

April 11, 2011

在先前的文章《[又拍网架构中的分库设计](#)》中，我有提到过MySQL分库设计中的主键选择问题。在这篇文章里我想对这个问题进行展开讨论，以此作为对上一篇文章的一个补充。

前面提到[又拍网](#)采用了全局唯一的字段作为主键。比如拿照片表为例，虽然不同用户的照片数据存放在不同的Shard（或者说MySQL节点/实例，请参考《[又拍网架构中的分库设计](#)》）上，但是每一张照片拥有整个站点唯一的ID作为标示。

### 为什么要全局唯一？

我们在对数据库集群作扩容时，为了保证负载的平衡，需要在不同的Shard之间进行数据的移动，如果主键不唯一，我们就没办法这样随意的移动数据。起初，我们考虑采用组合主键来解决这个问题。一般会以user\_id和一个自增的photo\_id来作为主键，这的确能解决移动数据可能带来的主键冲突问题，但是就像在“[又拍网架构中的分库设计](#)”中描述的那样当Shard之间的数据发生关系后，我们需要用更多的字段来组成主键以保证唯一性，因此主键的索引会变的很大，从而影响查询性能，同时也会影响写入性能。

其次，每个Shard由两台MySQL服务器组成，而这两台服务器采用master-master的复制方式，以保证每个Shard一直可写。master-master复制方式必须保证在两台服务器上各自插入的数据有不同的主键，不然当复制到另外一台时就会出现主键重复错误。如果我们保证主键全局唯一，就自然的解决了这个问题。在没有采用数据拆分的设计当中，如果要用自增字段，可以参考[这篇文章](#)里的解决办法。

### 可能的解决方案

#### ■ UUID

或许可以采用UUID作为主键，但是UUID好长的一串，放在URL里好难看啊，有木有？当然这个不是关键所在，更重要的原因还是性能。UUID的生成没有顺序性，所以在写入

时，需要随机更改索引的不同位置，这就需要更多的IO操作，如果索引太大而不能存放在内存中的话就更是如此。而UUID索引时，一个key需要32个字节(当然如果采用二进制形式存储的话可以压缩到16个字节)，因此整个索引也会相对比较大。

## ■ MySQL自增字段

在单个MySQL数据库的应用中一般设置一个自增的字段就可以了，而在水平分库的设计当中，这种方法显然不能保证全局唯一。那么我们可以单独建立一个库用来生成ID，在Shard中的每张表在这个ID库中都有一个对应的表，而这个对应的表只有一个字段，这个字段是自增的。当我们需要插入新的数据，我们首先在ID库中的相应表中插入一条记录，以此得到一个新的ID，然后将这个ID作为插入到Shard中的数据的主键。这个方法的缺点就是需要额外的插入操作，如果ID库变的很大，性能也会随之降低。所以一定要保证ID库的数据集不要太大，一个办法是定期清理前面的记录。

## ■ 引入其它工具

Redis、Memcached等都支持原子性的increment操作，而且因为它们的优秀性能可以减少写入时的额外开销，也许我们可以拿它们当作序列生成器。Memcached的问题在于不持久性，所以我们不会考虑。而Redis也不是实时持久的，当然也可以配置成实时的，但那样怪怪的。当然也有一些持久的工具，比如Kyoto Cabinet、Tokyo Cabinet、MongoDB等等，传说中性能都不错，但是引入其它工具会增加架构的复杂程度，也会增加维护成本。我们的团队很小，精力有限，我们奉行够用就好的原则，也就是没有特别的原因，在可以接受的情况下，尽量用我们熟悉的工具解决问题。所以，我们还是来考虑一下怎么样用MySQL来解决这个问题吧。

## 更好的方案

我们一开始就是采用了上面所描述的MySQL自增字段的方法，后来看到《Ticket Servers: Distributed Unique Primary Keys on the Cheap》这篇文章里所描述的方法，豁然开朗。我经常这样想：如果没有那些开源产品、没有那些无私分享经验的人，光凭我们自己的能力能做到什么程度。很感谢那些人，所以我也尽量多的分享一些自己的经验。

我先描述一下Flickr那篇文章里所描述的方法，他们使用了REPLACE INTO这个MySQL的扩展功能。REPLACE INTO和INSERT的功能一样，但是当使用REPLACE INTO插入新数据行时，如果新插入的行的主键或唯一键(UNIQUE Key)已有的行重复时，已有的行会先被删除，然后再将新数据行插入。你可以放心，这是原子操作。

建立类似下面的表：

```
CREATE TABLE `tickets64` (  
  `id` bigint(20) unsigned NOT NULL auto_increment,  
  `stub` char(1) NOT NULL default '',  
  PRIMARY KEY (`id`),  
  UNIQUE KEY `stub` (`stub`)  
) ENGINE=MyISAM;
```

当需要获得全局唯一ID时，执行下面的SQL语句：

```
REPLACE INTO `tickets64` (`stub`) VALUES ('a');  
SELECT LAST_INSERT_ID();
```

第一次执行这个语句后，ticket64表将包含以下数据：

```
+-----+-----+  
| id      | stub |  
+-----+-----+  
| 1       | a    |  
+-----+-----+
```

以后再次执行前面的语句，stub字段值为'a'的行已经存在，所以MySQL会先删除这一行，再插入。因此，第二次执行后，ticket64表还是只有一行数据，只是id字段的值为2。这个表将一直只有一行数据。

Flickr为Photo, Group, Account, Task各自建立了一张ticket表以保持各自的ID的连续性。其它业务表的ID都使用同一个ticket表产生。

不错吧，其实还可以更棒。比如，只需要一张ticket表就可以为所有的业务表提供各自连续的ID。下面，来看一下我们的方法。首先来看一下表结构：

```
CREATE TABLE `sequence` (  
  `name` varchar(50) NOT NULL,  
  `id` bigint(20) unsigned NOT NULL DEFAULT '0',  
  PRIMARY KEY (`name`)  
) ENGINE=InnoDB;
```

注意区别，id字段不是自增的，也不是主键。在使用前，我们需要先插入一些初始化数据：

```
INSERT INTO `sequence` (`name`) VALUES  
('users'), ('photos'), ('albums'), ('comments');
```

接下来，我们可以通过执行下面的SQL语句来获得新的照片ID：

```
UPDATE `sequence` SET `id` = LAST_INSERT_ID(`id` + 1) WHERE `name` = 'photos';  
SELECT LAST_INSERT_ID();
```

我们执行了一个更新操作，将id字段增加1，并将增加后的值传递到LAST\_INSERT\_ID函数，从而指定了LAST\_INSERT\_ID的返回值。

实际上，我们不一定需要预先指定序列的名字。如果我们现在需要一种新的序列，我们可以直接执行下面的SQL语句：

```
INSERT INTO `sequence` (`name`) VALUES('new_business') ON DUPLICATE KEY UPDATE  
SELECT LAST_INSERT_ID();
```

这里，我们采用了INSERT ... ON DUPLICATE KEY UPDATE这个MySQL扩展，这个扩展的功能也和INSERT一样插入一行新的记录，但是当新插入的行的主键或唯一键(UNIQUE Key)和已有的行重复时，会对已有行进行UPDATE操作。

需要注意的是，当我们第一次执行上面的语句时，因为还没有name为'new\_business'的字段，所以正常的执行了插入操作，没有执行UPDATE，所以也没有为LAST\_INSERT\_ID传递值。所以之后执行SELECT LAST\_INSERT\_ID()返回的值不可确定，要看当前连接在此之前执行过什么操作，如果没有执行过会影响LAST\_INSERT\_ID值的操作，那么返回值将是0，不然就是该操作产生的值。所以，我们应该尽量避免使用这种方式。

UPDATE: 这个方法更容易解决单点问题，也不局限于两个服务器，只要对不同的服务器设置不同的初始值（但必须是连续的），然后将增量变为服务器数就行了。

## 总结一下

我还是那句话，够用就好。当然，也不是说就不要去了解其它产品、方案了。又拍网也在一些新兴的产品，比如Redis（在10年3月就开始在正式环境下使用了，算是比较早的使用者），因为它的引入的确能够更好、更方便、更高效的解决我们的某些问题。关键还是需要在购买前对其进行足够的了解。我会在后面的文章中介绍一下Redis的使用情况。



加入讨论...



chinalu • 4年前

Good

1 ^ | ▾ • 回复 • 分享 >



DolphinBoy • 2年前

UPDATE: 这个方法更容易解决单点问题，也不局限于两个服务器，只要对不同的服务器设置不同的初始值（但必须是连续的），然后将增量变为服务器数就行了。

这一句不是很明白，想不通这样怎么可以实现跨库情况？求解答！

^ | ▾ • 回复 • 分享 >



zolazhou 管理员 ➔ DolphinBoy • 2年前

比如有3台机器，那么我们只要将初始值分别设置为1，2，3. 然后执行下面的语句即可：

```
UPDATE `sequence` SET `id` = LAST_INSERT_ID(`id` + 3) WHERE `name` = 'photos';
```

```
SELECT LAST_INSERT_ID();
```

^ | ▾ • 回复 • 分享 >



DolphinBoy ➔ zolazhou • 2年前

哦，大概明白了，意思就是有一个主键数据库，往不同数据库写入主键的时候采取的步长不一样？

那我有两个问题:1. 怎么知道每次要写入的数据库是哪个数据库, 2. 这个步骤是在那个时候处理的, 是在mybatis中处理, 还是用触发器呢  
还有, 请教一下主键用int 还是 char 更好 呢？

^ | ▾ • 回复 • 分享 >



zolazhou 管理员 ➔ DolphinBoy • 2年前

数据对应应该写到哪个库？是指sharding的规则吗？写这个文章的时候，我们的设计是根据用户ID来拆分数据的，这个是在我们的数据访问框架这层完成。

如果是顺序的，就是说新的ID总比老的ID大或小，那么我觉得int和char没有太大差别。

^ | ▾ • 回复 • 分享 >



zhanghb • 2年前

敲入localhost结果进的不是自己机器的服务器, 跑到这里来了 --

^ | ▾ • 回复 • 分享 >



swoky • 2年前

恩，不错，回去试试，BTW，你这的代码高亮用的什么

^ | v • 回复 • 分享 >



zolazhou 管理员 → swoky • 2年前

<http://pygments.org>

^ | v • 回复 • 分享 >



yejianfeng • 3年前

学习了，还可以这样使用

^ | v • 回复 • 分享 >



Xiwang2116 • 4年前

请问：那个insert 之后是不是要马上提交事务？

^ | v • 回复 • 分享 >



zolazhou 管理员 → Xiwang2116 • 4年前

如果使用InnoDB, 并且关闭了autocommit, 在当前事务提交之前，其它连接将无法马上获取新的ID（因为行被锁定，会等待释放）。

所以，不建议将ID库和业务库放在一起，或者说，不要将ID获取包含在业务逻辑/事务中。

^ | v • 回复 • 分享 >



BeerBubble • 4年前

```
UPDATE `sequence` SET `id` = LAST_INSERT_ID(`id` + 1) WHERE `name` = 'photos';
```

```
SELECT LAST_INSERT_ID();
```

多次执行后，再执行

```
INSERT INTO `sequence` (`name`) VALUES('new_business') ON DUPLICATE KEY
```

```
UPDATE `id` = LAST_INSERT_ID(`id` + 1);
```

```
SELECT LAST_INSERT_ID();
```

插入新数据时，返回的id是错的~

^ | v • 回复 • 分享 >



zolazhou 管理员 → BeerBubble • 4年前

```
INSERT INTO `sequence` (`name`) VALUES('new_business') ON  
DUPLICATE KEY UPDATE `id` = LAST_INSERT_ID(`id` + 1);
```

在第一次执行时，不会更新LAST\_INSERT\_ID值，所以会返回前面的photos的ID值。

所以建议还是事先初始化好要生成的序列名称和初始值。

^ | v • 回复 • 分享 >



国恺 韩 → zolazhou • 4年前

这个bug最好原文更新一下，因为第一次生成的时候

LAST\_INSERT\_ID() 返回了上次别的类型的最后id，而不是从0开始的。而之后这个类型的id生成到第一个错误的 id 的时候就会导致生成的主键不唯一，最终出现主键重复错误！

我今天就遇到了这个错误，调试了半天才明白是这里的问题。所以想过来指出，原来前几天就有人指出了，可惜没来看评论，呵呵。

^ | v · 回复 · 分享 >



**zolazhou** 管理员 → 国恺 韩 · 4年前

好的，已更新

^ | v · 回复 · 分享 >



**tianshengpan alex** · 4年前

```
UPDATE `sequence` SET `id` = LAST_INSERT_ID(`id` + 1) WHERE `name` = 'photos';  
SELECT LAST_INSERT_ID();
```

这是一个原子操作？

^ | v · 回复 · 分享 >



**zolazhou** 管理员 → tianshengpan alex · 4年前

UPDATE 是原子操作，后面的SELECT加一起就不是了，不过没有关系，在同一个连接里，只要中间不执行会导致LAST\_INSERT\_ID值变更的操作，不管是什么时候SELECT LAST\_INSERT\_ID() 返回的都是刚才的值。

^ | v · 回复 · 分享 >



**tianshengpan alex** → zolazhou · 4年前

在多个连接中，并发的时候，是不是就会导致SELECT LAST\_INSERT\_ID();返回的是同一个值了

^ | v · 回复 · 分享 >



**zolazhou** 管理员 → tianshengpan alex · 4年前

不会，其它连接不会影响当前连接的LAST\_INSERT\_ID

^ | v · 回复 · 分享 >



**tianshengpan alex** → zolazhou · 4年前

SELECT LAST\_INSERT\_ID(); 只是跟当前的连接有关系？

^ | v · 回复 · 分享 >



**zolazhou** 管理员 → tianshengpan alex · 4年前

是的

^ | v · 回复 · 分享 >



**Abioy Sun** → zolazhou · 4年前

但这样你就要保证每次NEW一个ID就要新起一个连接？另外如  
用不加其他保护的连接，具有单点风险的，所以左避免单点的方

来个加其他保护的话，是有单点风险的。MICKI有避免单点的方法，也可以一起参考。用ttserver之类的incr操作还是挺方便的，这个维护成长作为长远的设计来讲应该是值得。

^ | v • 回复 • 分享 >

**zolazhou** 管理员 ➔ Abioy Sun • 4年前

一般情况下，在同步应用（非异步的，或者说基于线程/进程模型的应用），处理一个业务请求时不管新建连接（如果是PHP，建议新建连接）或者重用连接，都不会有问题。

至于单点问题，我描述的方法更容易解决，只要在不同的服务器上变一下初始值和增量就行了。

引入其它工具在管理上带来的麻烦是不这样做的一方面原因，另一个原因是我们需要将复杂的ID逻辑封装到我们开发框架里，我们不希望这个框架依赖其它工具。

^ | v • 回复 • 分享 >

**bruce** ➔ zolazhou • 6个月前

您好

我想请问一下，步长的设定

比如设定步长为2，一个步长为4

那么在步长为2执行第二次的时候，会取到步长为4的第一次id不就冲突了吗？

^ | v • 回复 • 分享 >

**zolazhou** 管理员 ➔ bruce • 6个月前

步长一定要设定为一样的值，不同在于初始值

^ | v • 回复 • 分享 >

**bruce** ➔ zolazhou • 6个月前

非常感谢，不过还是有点蒙。。

能麻烦详细点说一下吗？让您费心了

比如我有服务器A和B，分别有表sequence  
name | id 两个字段

如果理解没错的话，初始值指的是 sequence里面id字段么

服务器A初始值:

name | id

photo | 1

服务器B初始值:

name | id

photo | 2

如果只有图片需求:

在服务器A中执行



在服务器A中执行

```
UPDATE `sequence` SET `id` = LAST_INSERT_ID(`id` + 1)
WHERE `name` = 'photos';
```

在服务器B中执行

```
UPDATE `sequence` SET `id` = LAST_INSERT_ID(`id` + 1)
WHERE `name` = 'photos';
```

如果A中执行3次，B执行1次，他们得到的自增ID不是一样的吗？

^ | v • 回复 • 分享 >

**zolazhou** 管理员 ➔ bruce • 6个月前

步长设为2的话，就应该是：

```
UPDATE `sequence` SET `id` = LAST_INSERT_ID(`id` + 2)
WHERE `name` = 'photos';
```

因此服务器A得到的ID始终是奇数，B始终是偶数。

^ | v • 回复 • 分享 >

**bruce** ➔ zolazhou • 6个月前

嗯，谢谢啊，flicker那个倒是看明白了，就是服务器设定增量都是2，但是步长1|2 就能实现奇偶组合。

您这个我看的不太懂怎么实现不止两个服务器。。

“UPDATE: 这个方法更容易解决单点问题，也不局限于两个服务器，只要对不同的服务器设置不同的初始值（但必须是连续的），然后将增量变为服务器数就行了。”

:(

^ | v • 回复 • 分享 >

**Ramd Wang** ➔ zolazhou • 4年前

你好，zolazhou：

这个问题在php这种线程模型中确实没有问题，但是在java环境下，由于是多线程的，同一个连接可能有多个并发访问操作，会不会造成一些隐患？例如 SELECT LAST\_INSERT\_ID(); 得到的并不是上次执行update操作的id

^ | v • 回复 • 分享 >

**zolazhou** 管理员 ➔ Ramd Wang • 4年前

不需要担心这个问题，我在例子里面用的都是SQL，分成两个语句执行，所以会有这样的顾虑。实际上当我们使用客户端实现时，这个ID会在执行INSERT/UPDATE语句的返回包中就包含，客户端实现通常解析这个包后就将这个值保持在连接对象里。也就是说当你在调用mysql\_insert\_id()时（C客户端的方法，PHP基于它），是不需要再请求服务器的。Java的客户端也一样。

^ | v . 回复 . 分享 >



**Ramd Wang** → zolazhou • 4年前

我还不是很明白。mysql\_insert\_id()如果说是把last\_insert\_id放在了返回值里，那么如果我在java里执行这两句sql又是什么情况呢？

因为这个last\_insert\_id毕竟是上一句sql赋值才能得到，所以必须使用这两句sql才能完成操作呀。

^ | v . 回复 . 分享 >



**zolazhou** 管理员 → Ramd Wang • 4年前

不，一个MySQL Command就可以了，当执行了INSERT 或UPDATE语句，MySQL的返回包里就包含了这个ID。请参考<http://forge.mysql.com/wiki/My...>

^ | v . 回复 . 分享 >



**Liukaixuan** • 4年前

lilo?

^ | v . 回复 . 分享 >



**anhong wang** • 4年前

这个方法很巧，

^ | v . 回复 . 分享 >



**lamxjb** • 4年前

学习

^ | v . 回复 . 分享 >



订阅



在您的网站上使用Disqus



隐私

**DISQUS**

Copyright © 2010 Zola Zhou. All rights reserved.

本站采用HTML5和CSS3构建，所以请使用支持这些特性的浏览器浏览，推荐：Safari, Chrome, Firefox