



虽然 TiDB 具有不同于单机 RDBMS 的数据结构，但顺序的主键值写入，在 TiDB 上也会产生类的效果：TiKV 上一个的 region 被写满，进而分裂出一个新的 region，后续的写入转由新的 reigon 来承载。但甲之蜜糖，乙之砒霜，单机 RDBMS 的最佳实践放到 TiDB 上，会使写入压力总是集中在一个 region 上，这样就构成了持续的写入热点，无法发挥出 TiDB 这种分布式数据库的并行写入能力，降低了业务写入瓶颈，造成了系统资源的浪费。

TiDB v4.0 版本提供了便于迅速识别集群负载的 Dashboard 流量可视化页面（Key Visualizer），下图展示了写入热点的显示效果，中间一条明亮的曲线即标志着存在一张连续写入 Key 值的表。而右上侧的一组线条则显示出一个写入压力较为均匀的负载。Key Visualizer 的具体使用方法请参考[官方文档](#)。

图 2. 写入热点在 Dashboard Key Visualizer 中的显示效果

具体来说，TiDB 的写入热点是由于 TiKV 中 KV 的 Key 值连续写入造成的，根据 [TiDB 的编码规则](#)，在 TiDB v4.0 及更早的版本中，Key 的取值存在以下两种情况：

1. 当表的主键为单一字段，且该字段的类型为整型时，Key 值由该字段构成，Value 为所有字段值的拼接，因此整型主键的表为索引组织表。
2. 其他情况，TiDB 会为表构建一个隐藏列 \_tidb\_rowid，Key 值由该隐藏列构成，Value 为所有字段值的拼接，表的主键（如果有的话）构成一个非聚簇索引，即数据并不以主键来组织。拿具有非整型主键的表来举例，它需要比单 int 型主键的表多写一个索引。

对于第二种情况，为了避免由于隐藏列 \_tidb\_rowid 的顺序赋值而引起写入热点，TiDB 提供一个表属性 SHARD\_ROW\_ID\_BITS 来控制所生成的隐藏列的值分散到足以跳过一个 region 大小（96MB）的多个区间（分片）上，再借助 PD 的热点调度能力，最终将写入压力分摊到整个 TiKV 集群中。由于隐藏列不具有任何业务属性，因此这种打散热点的方法是对用户透明的。一般来说，我们建议用户为所有非单一整型主键的表配置这个表属性，来消除这部分的热点隐患，详细使用方法请参考[官方文档](#)。

在第二章中描述的常见的四种序列号生成方案中，由于自增主键面对的是连续的整型数值的写入，因此它的打散方式比较特殊，请参考[官网文档](#)对自增主键进行打散。

对于其他三种方案而言，它们都具有集成到应用代码的能力，也因此具有一定的灵活性，本文将以 Twitter snowflake 为例，展示如何设计应用逻辑来获得较高的唯一 ID 生成效率。

## 在 TiDB 上高效的运行序列号生成服务

本测试基于两张表进行，在原始表结构中，主键为整型，其中一张表有一个索引，另一张表有两个索引，表结构如下：

```
CREATE TABLE T_TX_GLOBAL_LIST (
  global_tx_id varchar(32) NOT NULL,
  global_tx_no bigint NOT NULL,
  trace_id varchar(18) NOT NULL,
  busi_unique_seq varchar(18) NOT NULL,
  as_code varchar(10) NOT NULL,
  as_version varchar(5) NOT NULL,
  framework_type char(1) NOT NULL,
  tx_stat char(1) NOT NULL,
  code char(2) DEFAULT NULL,
  msg varchar(32) DEFAULT NULL,
  create_time timestamp NOT NULL,
  update_time timestamp NULL DEFAULT NULL ON UPDATE CURRENT_TIMESTAMP,
  PRIMARY KEY (global_tx_no),
  KEY index1 (create_time,tx_stat)
);
CREATE TABLE T_TX_BRANCH_LIST (
  branch_tx_id varchar(32) NOT NULL,
  branch_tx_no bigint NOT NULL,
  global_tx_no bigint NOT NULL,
  trace_id varchar(18) NOT NULL,
  busi_unique_seq varchar(18) NOT NULL,
  ms_code varchar(10) NOT NULL,
  ms_version varchar(5) NOT NULL,
  framework_type char(1) NOT NULL,
  tx_stat char(1) NOT NULL,
  code char(2) DEFAULT NULL,
  msg varchar(32) DEFAULT NULL,
  input_data longtext NOT NULL,
  create_time timestamp NOT NULL,
  update_time timestamp NULL DEFAULT NULL ON UPDATE CURRENT_TIMESTAMP,
  PRIMARY KEY (branch_tx_no),
  KEY index1 (create_time,tx_stat),
  KEY index3 (global_tx_no)
);
```

基于这两张表，我们编写了一个压测程序，压测的场景为批量写入（即 batch insert，形如 `insert into t values(),(),(),(),().....();`），每个事务向 `T_TX_GLOBAL_LIST` 表写入 20 行记录，向 `T_TX_BRANCH_LIST` 表写入 100 行记录。两张表中的 `global_tx_no` 字段和 `branch_tx_no` 字段（高亮）使用 Twitter snowflake 生成。

Twitter snowflake 生成的唯一序列号类型为整型，由于序列号的前面大部分的 bit 位由时间戳和机器号占据，只有最后的几个 bit 位为递增序列值，因此在一个时间段内生成的序列号的前几位数值相同：



```
561632049706827776
561632049706827777
561632049706827778
561632049706827779
561632049706827780
561632049706827781
...  . . .
```

我们将通过以下三个实验来展示如何打散 Twitter snowflake 的写入热点。

1.第一个实验中，我们采用默认的表结构和默认 snowflake 设置，向表写入整型序列号，压测持续了 10h。通过 Key Visualizer 展示的负载可以发现明显的写入热点。写入点有 5 个，对应着两张表和 3 个索引。其中一条写入负载非常明亮，是整张图中写入压力最大的一部分，从左侧的标识可以看到是 T\_TX\_BRANCH\_LIST 的表记录部分的写入。

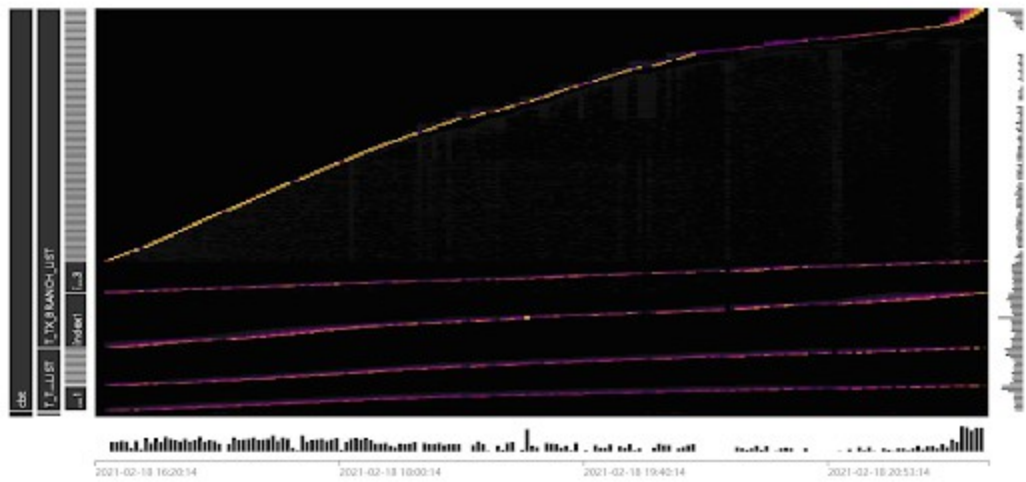


图 3. 默认设置下的写入负载

10h 的测试，两张表写入记录数为：

表名	10h 写入记录数
T_TX_GLOBAL_LIST	76685700
T_TX_BRANCH_LIST	383428500

2.对 Snowflake 生成的序列号进行转换，将最后一位数字移动到左数第二个数字的位置，原左数第二位数字及之后的所有数字向右移动一位。以此来让生成的 ID 跨越 96MB 的 region 容量，落在 10 个不同的 region 中。直接在二进制 id 上做位运算会导致转换后的十进制 id 位数不稳定，因此这个转换需要将整型的序列号先转为字符型，进行文本操作换位之后再转为整型，经测试，这个转换带来 10% 左右的额外消耗，由于这个额外消耗发生在应用程序中，相对于延迟较高的数据库，其带来的额外的影响在整个压测链路中微乎其微。

原始序列号	转换后的序列号
561632371724517376	566163237172451737
561632371728711680	506163237172871168
561632371728711681	516163237172871168
561632371728711682	526163237172871168
561632371732905984	546163237173290598
561632371732905985	556163237173290598
561632371732905986	566163237173290598
561632371732905987	576163237173290598
561632371732905988	586163237173290598
561632371737100288	586163237173710028

压测持续了 10h。通过 Key Visualizer 展示的负载可以看到，两张表的记录部分已经被各自打散到 10 个写入分片上，三个索引的其中一个由于字段值的转换，也呈现出一种较为分散的负载，负载图的整体亮度比较均衡，没有明显的写入热点。



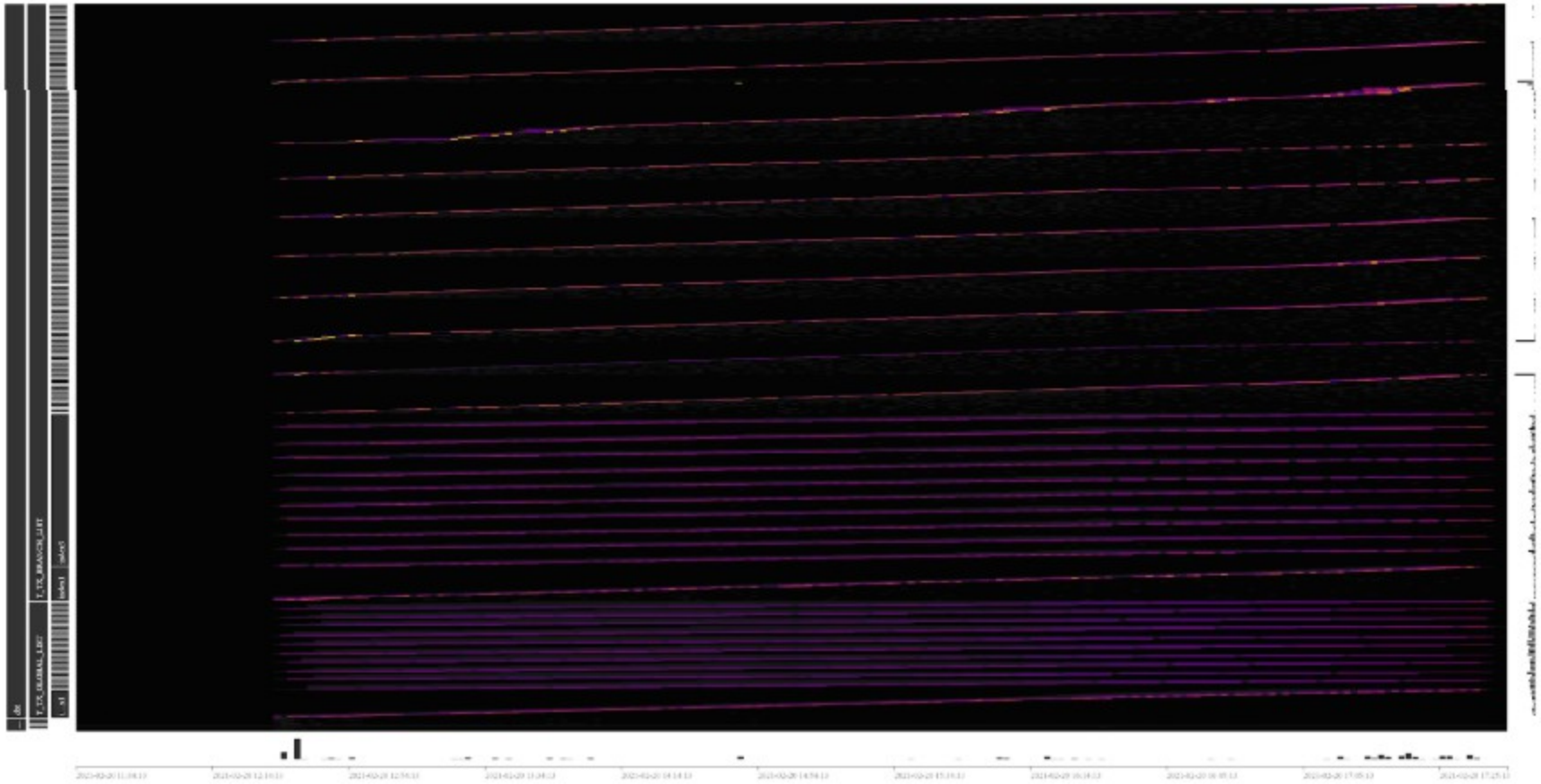


图 4. 序列号换位后的写入负载

10h 的测试，两张表写入记录数为：

表名	10h 写入记录数
T_TX_GLOBAL_LIST	150778840
T_TX_BRANCH_LIST	753894200

3.将两张表中的 global\_tx\_no 字段和 branch\_tx\_no 字段改为字符型，这样两张表从单一整型主键的索引组织表变为了按隐藏列组织的表。对两张表增加 shard\_row\_id\_bits=4 pre\_split\_regions=4 参数，以分散写入压力。由于主键类型发生了变化，还需要再程序中对 snowflake 生成的序列号类型做整型到字符型的转换。

压测持续了 10h。通过 Key Visualizer 展示的负载可以看到，存在 5 条明亮的线条，这是由于两张表的主键变为了非聚簇索引，导致需要单独 region 来存放主键，索引的数目因此变为了 5 个。而数据部分由于增加了打散参数，各自呈现出 16 个分片的均匀写入负载。

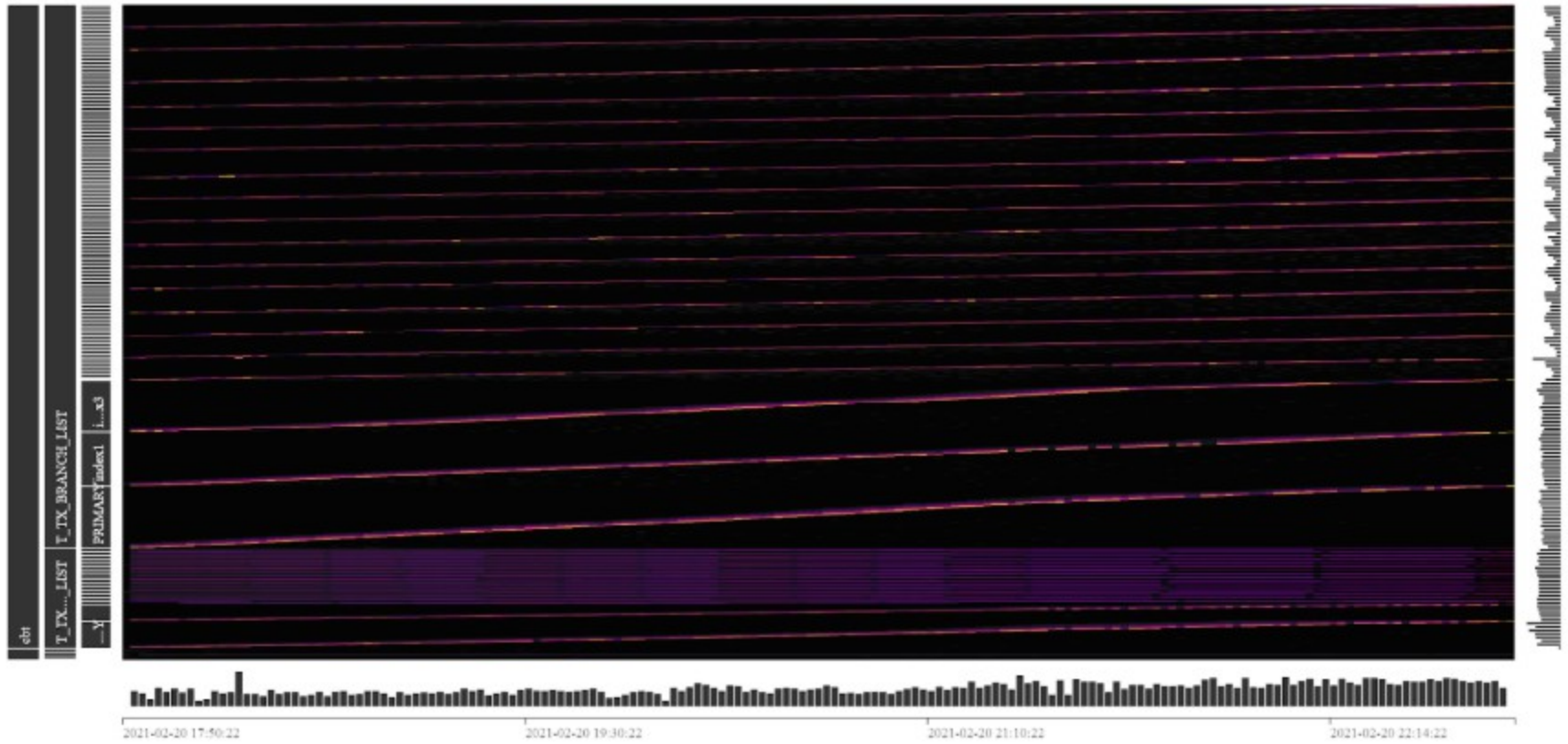


图 5. 将主键转为字符型并打散后的写入负载

10h 的测试，两张表写入记录数为：

表名	10h 写入记录数
T_TX_GLOBAL_LIST	136921680
T_TX_BRANCH_LIST	684608400

## 测试结论及示例代码

### 测试结论

- a. 从下面的测试成绩表可以看出，默认表结构配合 snowflake 默认配置生成的序列号，由于存在严重的写入热点，其写入性能较另外两个测试有较大的差距。
- b. 整型主键配合序列号换位，获得了本次测试中的最佳性能。我们还另外进行了末尾 2 位数字与末尾 3 位数字的换位测试，但过多的写入分片（2 位数字 100 个分片，3 位数字 1000 个分片）反而拖慢了写入性能，一般来讲使分片数量接近集群 tikv 实例个数可以充分的发挥集群的性能，用户需要根据自身的集群规模来制订换位策略。



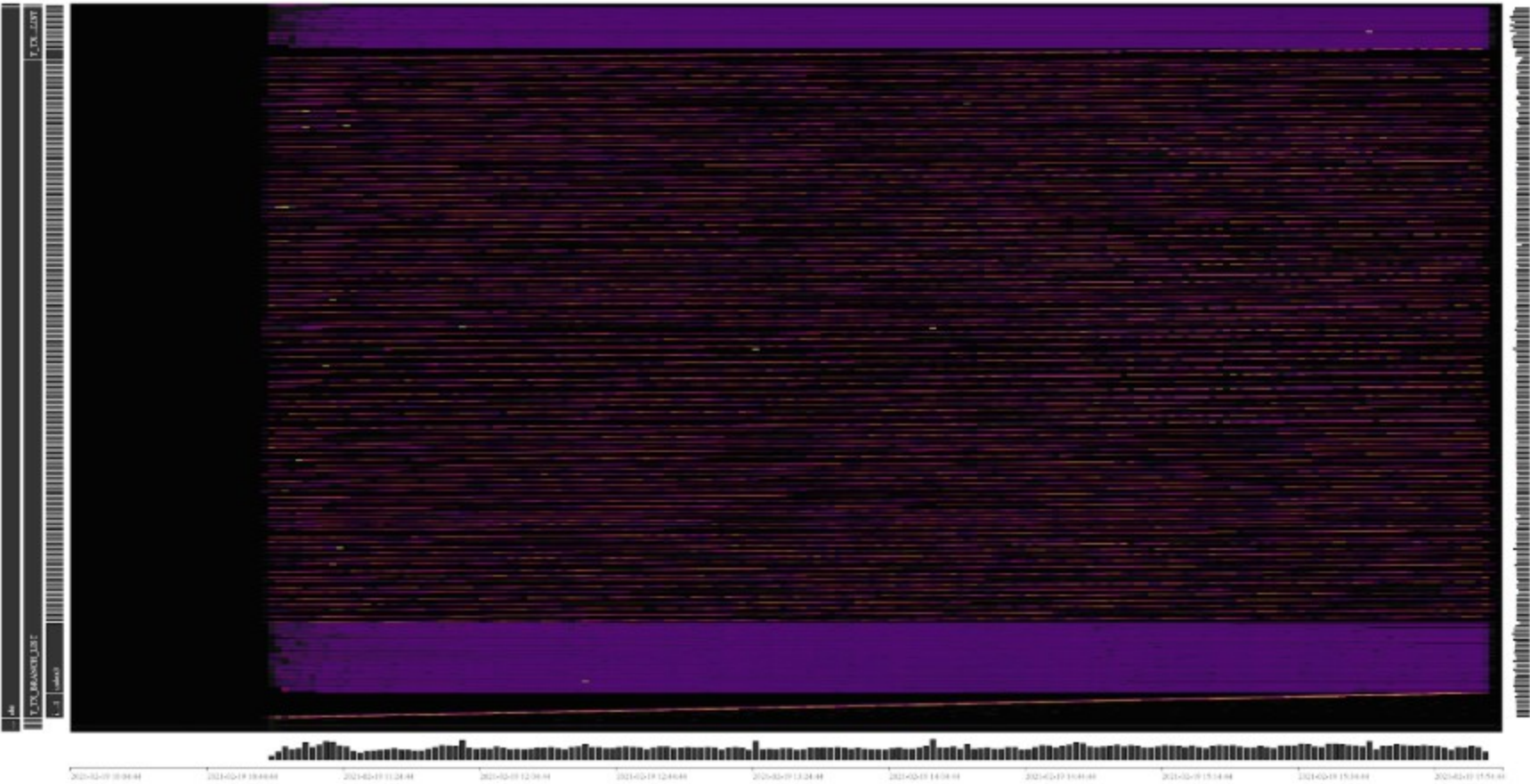


图 6. 整型主键换位（2 位置换），100 个分片



图 7. 整型主键换位（3 位置换），1000 个分片

c. 字符型主键 shard 也获得了不错的写入性能，但由于额外的热点索引写入，其性能略低于序列号换位方案。易用性是它的优势，用户可以通过简单的表结构变更来获取优异的写入性能。

测试轮次	T_TX_GLOBAL_LIST 表记录数（行）	T_TX_BRANCH_LIST 表记录数（行）
测试一，整型主键默认配置	76685700	383428500
测试二，整型主键换位（1 位）	150778840	753894200
测试二，整型主键换位（2 位）	123447080	617235400
测试二，整型主键换位（3 位）	101330340	506651700
测试三，字符型主键 shard	136921680	684608400

示例代码

```
###
public synchronized long nextId() {
    long currStmp = getNewstmp();
    if (currStmp < lastStmp) {
        throw new RuntimeException("Clock moved backwards. Refusing to generate id");
    }
    if (currStmp == lastStmp) {
        //相同毫秒内, 序列号自增
        sequence = (sequence + 1) & MAX_SEQUENCE;
        //同一毫秒的序列数已经达到最大
        if (sequence == 0L) {
            currStmp = getNextMill();
        }
    } else {
        //不同毫秒内, 序列号置为0
        sequence = 0L;
    }
    lastStmp = currStmp;
}
/**
 * XX.....XX XX000000 00000000 00000000    时间差 XX
 *           XXXXX0 00000000 00000000    数据中心ID XX
 *           X XXXX0000 00000000    机器ID XX
 *           XXXX XXXXXXXX    序列号 XX
 * 三部分进行|位或运算: 如果相对位都是0, 则结果为0, 否则为1
 */
long id = (currStmp - START_STMP) << TIMESTMP_LEFT //时间戳部分
        | datacenterId << DATACENTER_LEFT           //数据中心部分
        | machineId << MACHINE_LEFT                 //机器标识部分
        | sequence;                                  //序列号部分

String strid = String.valueOf(id);
int length = strid.length();
String lastnum = strid.substring(length - 1, length);
String str = strid.substring(1, length - 1);
String head = strid.substring(0, 1);
StringBuilder sb = new StringBuilder(head).append(lastnum).append(str);
return Long.valueOf(sb.toString());
}
###
```

- 必须使用时间差作为首段做位运算，因为其它段调整为首段做位运算会出现生成序列号位数不一致的问题。
- 使用字符串拼接方式效率虽然降低，但是从一次交易总体时间上看是可以忽略不记的。

## 结语

当前版本（v4.0）的易用性还有待加强，TiDB v5.0 版本将正式推出聚簇索引功能，新版本中的聚簇索引将支持任意类型的索引字段，而具有整型主键的表也可以被设置为非主键组织表，这代表采用整型主键的表可以很便捷的通过表属性 SHARD\_ROW\_ID\_BITS 来分散写入热点，大家敬请期待！

TiDB

### 关于我们

公司概况

发展历程

新闻中心

市场活动

加入我们

隐私政策

Cookie 政策

安全合规

版本支持策略

漏洞管理策略

网站使用条款

### 资源中心

社区

TiDB 文档

TiDB 6.x in Action

快速上手指南

社区问答-AskTUG

博客

GitHub

PingCAP Education

商标下载及使用

### 联系我们

商务咨询

400-6790-886

010-58400041

info@pingcap.com

前台总机

010-53326356

媒体合作

pr@pingcap.com

### PingCAP 公司

PingCAP 是业界领先的企业级开源分布式数据库企业，提供包括开源分布式数据库产品、解决方案与咨询、技术支持与培训认证服务，致力于为全球行业用户提供稳定高效、安全可靠、开放兼容的新型数据服务平台，解放企业生产力，加速企业数字化转型升级。



联系我们

友情链接