

浅析下开源时序数据库VictoriaMetrics的存储机制



关注他

20 人赞同了该文章

[VictoriaMetrics](#)是一个快速高效且可扩展的监控解决方案和时序数据库，可以作为Prometheus的长期远端存储，具备的特性有：

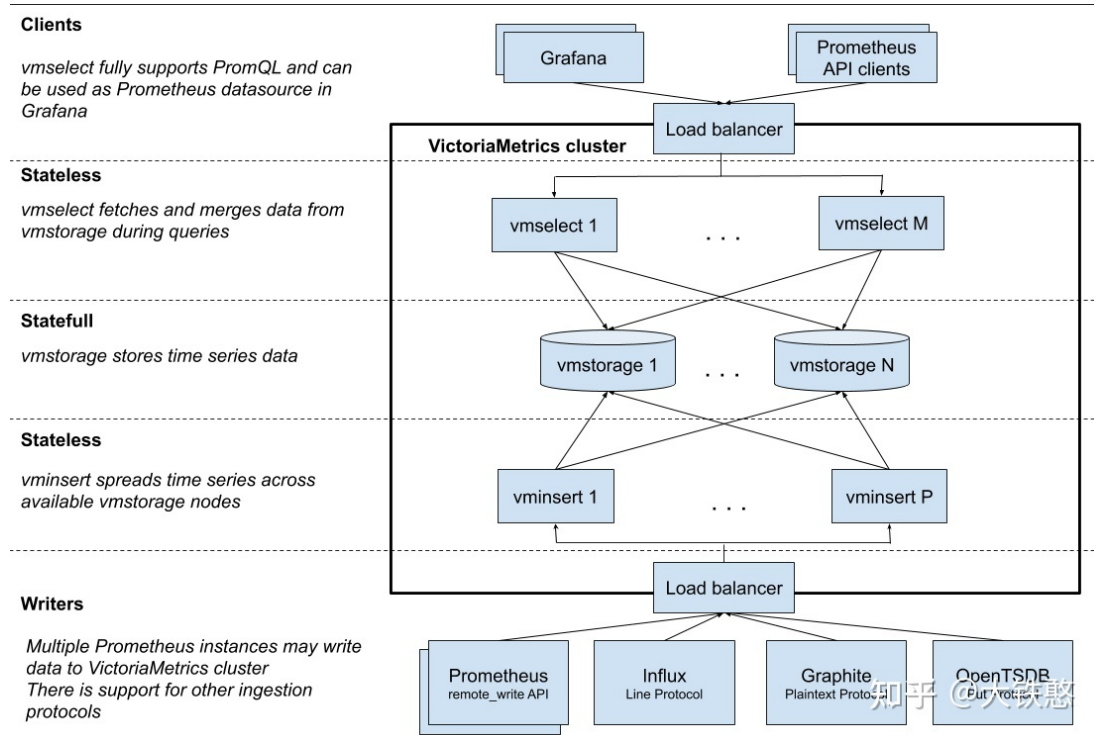
- 支持prometheus查询api，同时实现了一个metricsql 查询语言
- 支持全局查询视图，支持多prometheus 实例写数据到VictoriaMetrics，然后提供一个统一的查询
- 支持集群
- 高性能
- 支持多种协议，包括influxdb line协议，prometheus metrics，graphite，prometheus远端写api，opentsdb http协议等
- 高压缩比

本文主要分析下VictoriaMetrics的存储机制，包括整体架构、数据模型、磁盘目录、文件格式等部分，对应的源码版本号为v1.45.0。

1.整体架构

VictoriaMetrics的集群主要由vmstorage、vminsert、vmselect等三部分组成，其中，vmstorage 负责提供数据存储服务；vminsert 是数据存储 vmstorage 的代理，使用一致性hash算法进行写入分片；vmselect 负责数据查询，根据输入的查询条件从 vmstorage 中查询数据。

VictoriaMetrics的这个三个组件每个组件都可以单独进行扩展，并运行在大多数合适软件上。vmstorage采用shared-nothing架构，优点是vmstorage的节点相互之间无感知，相互之间无需通信，不共享任何数据，增加了集群的可用性、简化了集群的运维和集群的扩展。



另外，VictoriaMetrics集群支持多个隔离的租户特性，又名命名空间。租户通过accountID (或者accountID:projectID) 来标识，这些ID放在请求URL中。有关VictoriaMetrics租户的一些事实：

- 每个accountID和projectID由[0 .. 2 ^ 32) 范围内的任意32位整数标识。如果缺少projectID，则会将其自动分配为0。预计有关租户的其他信息（例如身份验证令牌，租户名称，限制，计费等）将存储在单独的关系数据库中。此数据库必须由位于VictoriaMetrics群集前面的单独服务（例如vmauth）管理。
- 将第一个数据点写入给定租户时，将自动创建租户。
- 所有租户的数据平均分布在可用的vmstorage节点之间。当不同的租户具有不同的数据量和不同的查询负载时，这可以保证在vmstorage节点之间平均负载。

整体上来说，VictoriaMetrics支持多租户，但租户的信息需要使用额外的关系型数据库来存储，且VictoriaMetrics不支持在单个请求中查询多个租户。

2.数据模型

开门见山，**VictoriaMetrics采用的数据模型是单值模型，且只支持浮

点数指标**。那么到底什么是单值模型呢？目前，常见的时序数据库的数据模型，主要分成单值模型和多值模型。这里简单说明下单值模型和多值模型，整体上，可以认为单值模型是多值模型的一个特例。

单值模型是根据****业务指标数据****建模，按照单个指标的细粒度进行数据使用和逻辑存储，如下图所示，一行数据只有一个指标值，即value列。目前采用单值模型的时序数据库，有OpenTSDB、[KairosDB](#)、Prometheus等。

metric	appName	region	timestamp	value
cpu	tsdb	shanghai	1613706029	0.5
cpu	tsdb	shanghai	1613706039	0.4
io	tsdb	hangzhou	1613706049	0.7
io	tsdb	hangzhou	1613706059	0.8

多值的模型则是针对****数据源****建模，我们每一行数据针对的是一个数据源，它的被测量的多个指标在同一列上，如下图所示，一行数据有多个指标值，即有cpu和io两列。目前采用多值模型的时序数据库，有InfluxDB、TimescaleDB等。

measurement	appName	region	timestamp	cpu	io
monitor	tsdb	shanghai	1613706029	0.5	0.7
monitor	tsdb	shanghai	1613706039	0.5	0.5

3.磁盘目录

VictoriaMetrics的根目录下主要包括4个目录或文件，如下图所示。其中，最主要的是数据目录data和索引目录indexdb，flock.lock文件为文件锁文件，用于VictoriaMetrics进程锁住文件，不允许别的进程进行修改目录或文件。

```
1 |— data
2 |— flock.lock
3 |— indexdb
4 |— snapshots
```

知乎 @大铁憨

3.1 数据目录

数据目录data的具体结构，如下图所示，在图中使用红色文字，对主要目录或文件做了简单说明，其中最主要的是****big****目录和****small****目录，这两个目录的结构是一样的。其中，在VictoriaMetrics中，使用table来表示的数据或者索引的根目录，而实际上VictoriaMetrics中没有实际的表table级别目录。

```
1 |— data # table 目录
2 |   |— big # bigPartition 目录
3 |     |— 2020_11 # Partition 目录，必须和 small 中目录一一对应，big 和 small 中 Partition 目录是一起创建的
4 |       |— tmp
5 |       |— txn
6 |       |— snapshots
7 |   |— flock.lock
8 |   |— small
9 |     |— 2020_11 # Partition 目录，命名模式为 YYYY_MM, ==> partition.go
10 |       |— 354_354_20201103102204.255_20201103102204.255_1643F83394CA24A8 # part 目录, ==> part.go
11 |         |— index.bin
12 |         |— metaindex.bin
13 |         |— timestamps.bin
14 |         |— values.bin
15 |       |— 708_354_20201103102134.255_20201103102149.255_1643F83394CA24A7
16 |         |— index.bin
17 |         |— metaindex.bin
18 |         |— timestamps.bin # 时间戳列数据
19 |         |— values.bin # value 列数据
20 |       |— tmp # 临时目录，merge 或者 flush 文件时的临时目录
21 |       |— txn # 事务目录，在 mergeParts 完成后，将 tmpPartPath 和 dstPartPath 的原子
22 |         |— snapshots # 事务目录，在 mergeParts 完成后，将 tmpPartPath 和 dstPartPath 的原子
23 |   |— flock.lock
24 |   |— indexdb
25 |   |— snapshots
```

命名方式为: rowCount_blocksCount_minTime_maxTime_suffix

事务目录，在 mergeParts 完成后，将 tmpPartPath 和 dstPartPath 的原子写入/txn/%016X 文件中，然后开始执行重命名事务

知乎 @大铁憨

在small目录下，以月为单位不断生成partition目录，比如上图中的2020_11目录，对应的实现在源码lib/storage/partition.go中。partition目录包括part目录、临时目录tmp、事务目录txn等三个目录。

内存中的数据每刷盘一次就会生成一个part目录，如上图中的"708_354_20201103102134.255_20201103102149.255_1643F83394CA24A7"，目录名中的708表示这个目录下包含的数据行数

rowCount, 目录名中的354表示这个目录中包含的数据块数
blocksCount, 20201103102134.255表示目录中包含的数据的最小时间戳, 20201103102149.255表示目录中包含的数据的最大时间戳,
1643F83394CA24A7是生成这个目录时的系统纳秒时间戳的16进制表示, 对应的实现逻辑在源码lib/storage/part.go中;

看到这里, 可能会有一些疑问? 比如为何要分成big和small目录, 或者说big目录和small中的数据关系是什么? 这个需要从VictoriaMetrics的compaction机制讲起。

在VictoriaMetrics中, small目录和big目录都会周期性检查, 是否需要做part的合并。VictoriaMetrics默认每个10ms检查一次partition目录下的part是否需要做merge。如果检查出有满足merge条件的parts, 则这些parts合并为一个part。如果没有满足条件的part进行merge, 则以10ms为基进行指数休眠, 最大休眠时间为10s。

VictoriaMetrics在写数据时, 先写入在small目录下的相应partition目录下面的, small目录下的每个partition最多256个part。VictoriaMetrics在Compaction时, 默认一次最多合并15个part, 且small目录下的part最多包含1000W行数据, 即1000W个数据点。因此, 当一次待合并的parts中包含的行数超过1000W行时, 其合并的输出路径为big目录下的同名partition目录下。

因此, big目录下的数据由small目录下的数据在后台compaction时合并生成的。那么为什么要分成big目录和small目录呢?

这个主要是从磁盘空间占用上来考虑的。时序数据经常读取最近写入的数据, 较少读历史数据。而且, 时序数据的数据量比较大, 通常会使用压缩算法进行压缩。

数据进行压缩后, 读取时需要解压, 采用不同级别的压缩压缩算法其解压速度不一样, 通常压缩级别越高, 其解压速度越慢。在VictoriaMetrics在时序压缩的基础上, 又采用了ZSTD这个通用压缩算法进一步压缩了数据, 以提高压缩率。在small目录中的part数据, 采用的是低级别的ZSTD, 而big目录下的数据, 采用的是高级别的

ZSTD。

因此，VictoriaMetrics分成small目录和big目录，主要是兼顾近期数据的读取和历史数据的压缩率。

3.2 索引目录

索引目录indexdb的具体结构，如下图所示，在图中使用红色文字，对主要目录或文件做了简单说明。与数据目录不同的是，indexdb目录下由多个table目录，每个table目录代表一个完整的自治的索引，每个table目录下，又有多个不同的part目录，part命名方式比较简单，有文件包含的item数itemsCount和block数blocksCount, 以及根据系统纳秒时间戳自增生成的mergeldx的16进制表示。



知乎 @大铁憨

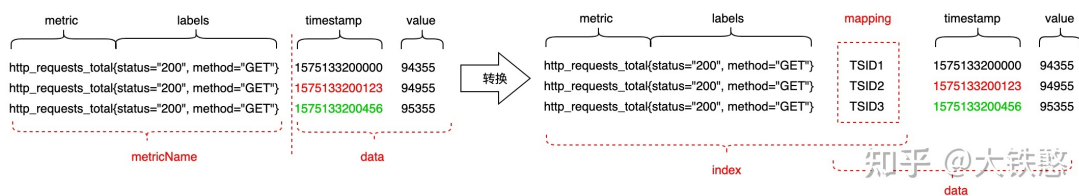
indexdb下面的形如"1643F4F397B53DEE"是怎么生成的，或者什么时候切换新的目录写索引的呢？VictoriaMetrics会根据用户设置的数据保留周期retention来定期滚动索引目录，当前一个索引目录的保留时间到了，就会切换一个新的目录，重新生成索引。

4. 文件格式

在介绍具体的文件格式之前，不得不提下VictoriaMetrics对于写入数据的处理过程。下图是VictoriaMetrics支持的Prometheus协议的一个写入示例。

metric	labels	timestamp	value
http_requests_total	{status="200", method="GET"}	1575133200000	94355
http_requests_total	{status="200", method="GET"}	1575133200123	94955
http_requests_total	{status="200", method="GET"}	1575133200456	95355

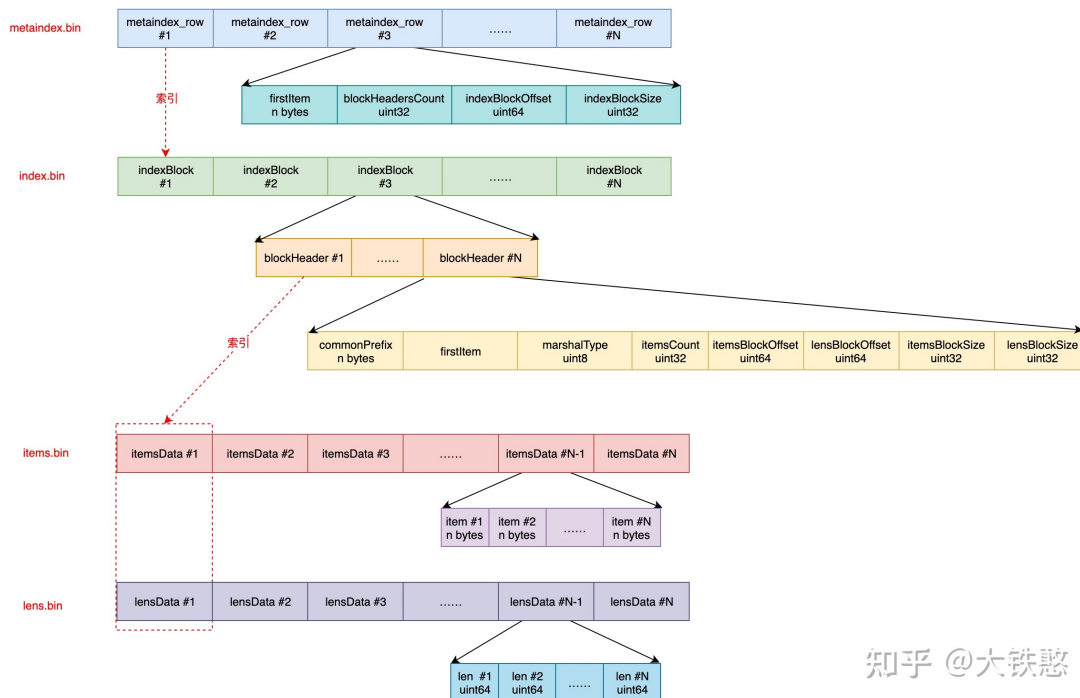
VictoriaMetrics在接受到写入请求时，会对请求中包含的时序数据做转换处理，如下图所示。首先先根据包含metric和labels的MetricName生成一个唯一标识TSID，然后metric + labels + TSID作为索引index，TSID + timestamp + value作为数据data，最后索引index和数据data分别进行存储和检索。



因此，VictoriaMetrics的数据整体上分成索引和数据两个部分，因此文件格式整体上会有两个部分。其中，索引部分主要是用于支持按照label或者tag进行多维检索。与大多数时序数据库的数据组织方式一样，比如InfluxDB、Prometheus、OpenTSDB等，VictoriaMetrics也是按时间线来组织数据的，即数据存储时，先将数据按TSID进行分组，然后每个TSID的包含的数据点各自使用列式压缩存储。

4.1 索引文件

VictoriaMetrics每次内存Flush或者后台Merge时生成的索引part，主要包含metaindex.bin、index.bin、lens.bin、items.bin等4个文件。这四个文件的关系如下图所示，metaindex.bin文件通过metaindex_row索引index.bin文件，index.bin文件通过indexBlock中的blockHeader同时索引lens.bin文件和items.bin文件。



知乎 @大铁憨

metaindex.bin文件中，包含一系列的metaindex_row，每个metaindex_row中包含最小项firstItem、索引块包含的块头部数blockHeadersCount、索引块偏移indexBlockOffset、索引块大小indexBlockSize。

- metaindex_row在文件中的位置按照firstItem的大小的字典序排序存储，以支持二分检索；
- metaindex.bin文件使用ZSTD进行压缩；
- metaindex.bin文件中的内容在part打开时，会全部读出加载至内存中，以加速查询过滤；
- metaindex_row包含的firstItem为其索引的IndexBlock中所有blockHeader中的字典序最小的firstItem；
- 查找时根据firstItem进行二分检索；

index.bin文件中，包含一系列的indexBlock, 每个indexBlock又包含一系列blockHeader，每个blockHeader的包含item的公共前缀commonPrefix、最小项firstItem、itemsData的序列化类型marshalType、itemsData包含的item数、item块的偏移itemsBlockOffset等内容。

- 每个indexBlock使用ZSTD压缩算法进行压缩；
- 在indexBlock中查找时，根据firstItem进行二分检索blockHeader；

items.bin文件中，包含一系列的itemsData, 每个itemsData又包含一系列的item。

- itemsData会根据情况是否使用ZSTD压缩，当item个数小于2时，或者itemsData的长度小于64字节时，不压缩；当itemsData使用ZSTD压缩后的大小，大于原始itemsData的0.9倍时，则不压缩，否则使用ZSTD算法进行压缩。
- 每个item在存储时，去掉了blockHeader中的公共前缀commonPrefix以提高压缩率。

lens.bin文件中，包含一系列的lensData, 每个lensData又包含一系列8字节的长度len，长度len标识items.bin文件中对应item的长度。在读取或者需要解析itemsData中的item时，先要读取对应的lensData中对应的长度len。当itemsData进行压缩时，lensData会先使用异或算法进行压缩，然后再使用ZSTD算法进一步压缩。

VictoriaMetrics索引文件都是围绕着item来组织的，那么item的结构是什么样子的？或者item的种类有哪些？在VictoriaMetrics中item的整体上是一个KeyValue结构的字节数组，共计有7种类型，每种类型的item通过固定前缀来区分，前缀类型如下图所示。

```

29  const (
30      // Prefix for MetricName->TSID entries.
31      nsPrefixMetricNameToTSID = 0
32
33      // Prefix for Tag->MetricID entries.
34      nsPrefixTagToMetricIDs = 1
35
36      // Prefix for MetricID->TSID entries.
37      nsPrefixMetricIDToTSID = 2
38
39      // Prefix for MetricID->MetricName entries.
40      nsPrefixMetricIDToMetricName = 3
41
42      // Prefix for deleted MetricID entries.
43      nsPrefixDeletedMetricID = 4
44
45      // Prefix for Date->MetricID entries.
46      nsPrefixDateToMetricID = 5
47
48      // Prefix for (Date,Tag)->MetricID entries.
49      nsPrefixDateTagToMetricIDs = 6
50  )

```

知乎 @大铁憨

VictoriaMetrics是怎么基于item支持tag多维检索的呢？这里就不得不提下VictoriaMetrics的TSID的生成过程。

VictoriaMetrics的MetricName的结构如下所示，包含MetricGroup字节数组和Tag数组，其中，MetricGroup是可选的，每个Tag由Key和Value等字节数组构成。

```

1 // MetricName represents a metric name.
2 type MetricName struct {
3     MetricGroup []byte
4
5     // Tags are optional. They must be sorted by tag Key for canonical view.
6     // Use sortTags method.
7     Tags []Tag
8 }

```

```

1 // Tag represents a (key, value) tag for metric.
2 type Tag struct {
3     Key []byte
4     Value []byte
5 }

```

知乎 @大铁憨

VictoriaMetrics的TSID的结构如下所示，包含MetricGroupID, JobID, InstanceID, MetricID等四个字段，其中除了MetricID外，其他三个字段都是可选的。这个几个ID的生成方法如下：

- metricGroupID根据MetricName中的MetricGroup使用xxhash的sum64算法生成。
- JobID和InstanceID分别由MetricName中的tag的第一个tag和第二个tag使用xxhash的sum64算法生成。为什么使用第一个tag和第二个tag？这是因为VictoriaMetrics在写入时，将写入请求中的JobID和InstanceID放在了Tag数组的第一个和第二个位置。
- MetricID，使用VictoriaMetrics进程启动时的系统纳秒时间戳自增生

```
1 // TSID is unique id for a time series.
2 // Time series blocks are sorted by TSID.
3 // All the fields except MetricID are optional.
4 type TSID struct {
5     // MetricGroupID is the id of metric group.
6     MetricGroupID uint64
7
8     // JobID is the id of an individual job (aka service).
9     // JobID must be unique.
10    JobID uint32
11
12    // InstanceID is the id of an instance (aka process).
13    // InstanceID must be unique.
14    InstanceID uint32
15
16    // MetricID is the unique id of the metric (time series).
17    // All the other TSID fields may be obtained by MetricID.
18    MetricID uint64
19 }
```

知乎 @大铁憨

因为TSID中除了MetricID外，其他字段都是可选的，因此TSID中可以始终作为有效信息的只有metricID，因此VictoriaMetrics的在构建tag到TSID的字典过程中，是直接存储的tag到metricID的字典。

以写入http_requests_total{status="200", method="GET"}为例，则MetricName为http_requests_total{status="200", method="GET"}，假设生成的TSID为{metricGroupID=0, jobID=0, instanceID=0, metricID=51106185174286}，则VictoriaMetrics在写入时就构建了如下几种类型的索引item，其他类型的索引item是在后台或者查询时构建的。

- metricName -> TSID, 即http_requests_total{status="200",

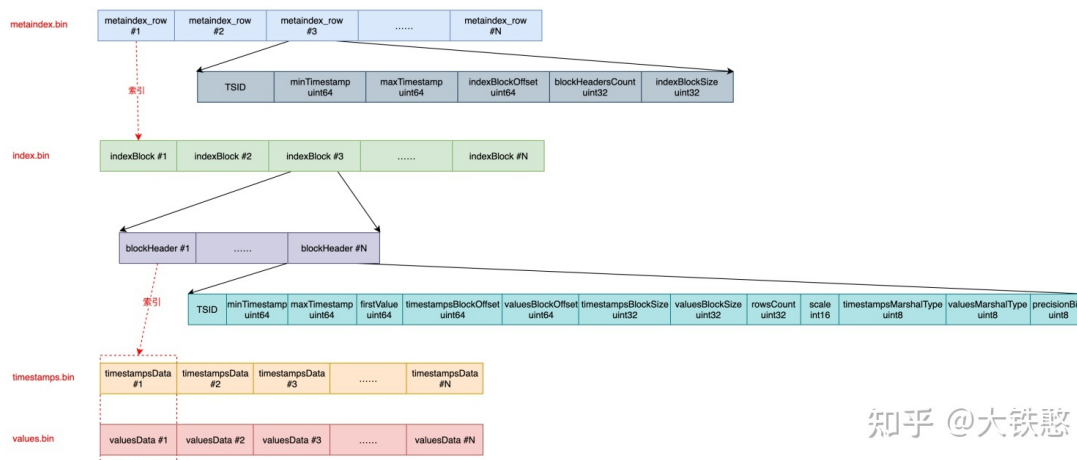
```
method="GET"} -> {metricGroupID=0, jobID=0, instanceID=0,
metricID=51106185174286}
- metricID -> metricName, 即51106185174286 ->
http_requests_total{status="200", method="GET"}
- metricID -> TSID, 即51106185174286 -> {metricGroupID=0, jobID=0,
instanceID=0, metricID=51106185174286}
- tag -> metricID, 即 status="200" -> 51106185174286, method="GET"
-> 51106185174286, "" = http_requests_total -> 51106185174286
```

有了这些索引的item后，就可以支持基于tag的多维检索了，在当给定查询条件http_requests_total{status="200"}时，VictoriaMetrics先根据给定的tag条件，找出每个tag的metricID列表，然后求所有tag的metricID列表的交集，然后根据交集集中的metricID，再到索引文件中检索出TSID，根据TSID就可以到数据文件中查询数据了，在返回结果之前，再根据TSID中的metricID，到索引文件中检索出对应的写入时的原始MetricName。

但是由于VictoriaMetrics的tag到metricID的字典，没有将相同tag的所有metricID放在一起存储，在检索时，一个tag可能需要查询多次才能得到完整的metricID列表。另外查询出metricID后，还要再到索引文件中去检索TSID才能去数据文件查询数据，又增加了一次IO开销。这样来看的话，VictoriaMetrics的索引文件在检索时，如果命中的时间线比较多的情况下，其IO开销会比较大，查询延迟也会比较高。

4.2 数据文件

VictoriaMetrics每次内存Flush或者后台Merge时生成的数据part，包含metaindex.bin、index.bin、timestamps.bin、values.bin等4个文件。这四个文件的关系如下所示。metaindex.bin文件索引index.bin文件，index.bin文件同时索引timestamps.bin和values.bin文件。



知乎 @大铁憨

metaindex.bin文件中，包含一系列的`metaindex_row`，每个`metaindex_row`中包含时间线标识`TSID`、最小时间戳`minTimestamp`、最大时间戳`maxTimestamp`、索引块偏移`indexBlockOffset`、索引块大小`indexBlockSize`、索引块包含的块头部数`blockHeadersCount`。

- `metaindex_row`在文件中的位置按照`TSID`的大小的字典序排序存储；
- `metaindex.bin`文件使用ZSTD进行压缩；
- `metaindex.bin`文件中的内容在`part`打开时，会全部读出加载至内存中，以加速查询过滤；
- `metaindex_row`包含时间线标识`TSID`为其索引的`IndexBlock`中所有`blockHeader`中的最小时间标识`TSID`；
- `metaindex_row`包含最小时间戳`minTimestamp`为其索引的`IndexBlock`中所有`blockHeader`中的最小时间戳`minTimestamp`；
- `metaindex_row`包含最大时间戳`maxTimestamp`为其索引的`IndexBlock`中所有`blockHeader`中的最大时间戳`maxTimestamp`；
- 查找时根据`TSID`进行二分检索；

index.bin文件中，包含一系列的`indexBlock`，每个`indexBlock`又包含一系列`blockHeader`，每个`blockHeader`的包含时间线标识`TSID`、最小时间戳`minTimestamp`、最大时间戳`maxTimestamp`、第一个指标值`firstValue`、时间戳数据块偏移`timestampsBlockOffset`、指标值数据块偏移`valuesBlockOffset`等内容。

- 每个`indexBlock`使用ZSTD压缩算法进行压缩；
- 查找时，线性遍历`blockHeader`查找`TSID`；

timestamps.bin文件中，包含一系列时间线的时间戳压缩块
timestampsData; values.bin文件中，包含的一系列时间线的指标值压缩块valuesData。其中，timestampsData和values.data会根据时序数据特征进行压缩，整体上的压缩思路是：先做时序压缩，然后在做通用压缩。比如，先做delta-of-delta计算或者异或计算，然后根据情况做zig-zag，最后再根据情况做一次ZSTD压缩，VictoriaMetrics支持的压缩算法或者类型主要有6种，如下图所示，压缩编码源码在lib/encoding/encoding.go文件中。

```
20  const (
21      // MarshalTypeZSTDNearestDelta2 is used for marshaling counter
22      // timeseries.
23      MarshalTypeZSTDNearestDelta2 = MarshalType(1)
24
25      // MarshalTypeDeltaConst is used for marshaling constantly changed
26      // time series with constant delta.
27      MarshalTypeDeltaConst = MarshalType(2)
28
29      // MarshalTypeConst is used for marshaling time series containing only
30      // a single constant.
31      MarshalTypeConst = MarshalType(3)
32
33      // MarshalTypeZSTDNearestDelta is used for marshaling gauge timeseries.
34      MarshalTypeZSTDNearestDelta = MarshalType(4)
35
36      // MarshalTypeNearestDelta2 is used instead of MarshalTypeZSTDNearestDelta2
37      // if compression doesn't help.
38      MarshalTypeNearestDelta2 = MarshalType(5)
39
40      // MarshalTypeNearestDelta is used instead of MarshalTypeZSTDNearestDelta
41      // if compression doesn't help.
42      MarshalTypeNearestDelta = MarshalType(6)
43  )
```

知乎 @大铁憨

[VictoriaMetrics文档](#)中提及在生产环境中，每个数据点（8字节时间戳 + 8字节value共计16字节）压缩后小于1个字节，最高可达 0.4字节，如下所示。

VictoriaMetrics

Based on production data from our customers, sample size is 0.4 byte That means one metric with 10 seconds resolution will need 6307200 points * 0.4 bytes/point = 2522880 bytes or 2.4 megabytes. Calculation for number of metrics can be stored in 2 tb disk: 219902325552 (disk size) / 2522880 (one metric for 2 year) = 871632 metrics So in 2tb we can store 871 632 metrics

5. VictoriaMetrics总结

VictoriaMetrics的整体的存储设计还是不错的，比如数据时间分区、数据压缩率高、丰富的生态协议等。但VictoriaMetrics的标签索引、数据

可靠性、支持的数据类型等方面还存在一些不足。比如，标签索引查询多次IO，可能在时间线数量非常多的场景下，其检索效率会比较低，且没有WAL，在写入时可能会存在数据丢失的风险。目前只支持浮点数类型，不支持布尔、字符串、字节数组等类型。

编辑于 2021-11-13 20:42