



ClickHouse 索引与优化

易顺（宙游）阿里云数据库研发

本次分享尝试回答的问题

- 1.ClickHouse 数据在物理文件中怎么存的，不同文件之间的联系是什么？
- 2.针对不同业务，我们应该怎么选择 Order by key、Primary key？
- 3.Skipping Index 是什么？怎么选择？
- 4.怎么用 explain 语句查看 SQL 的执行计划？
- 5.ClickHouse 的一般优化策略。

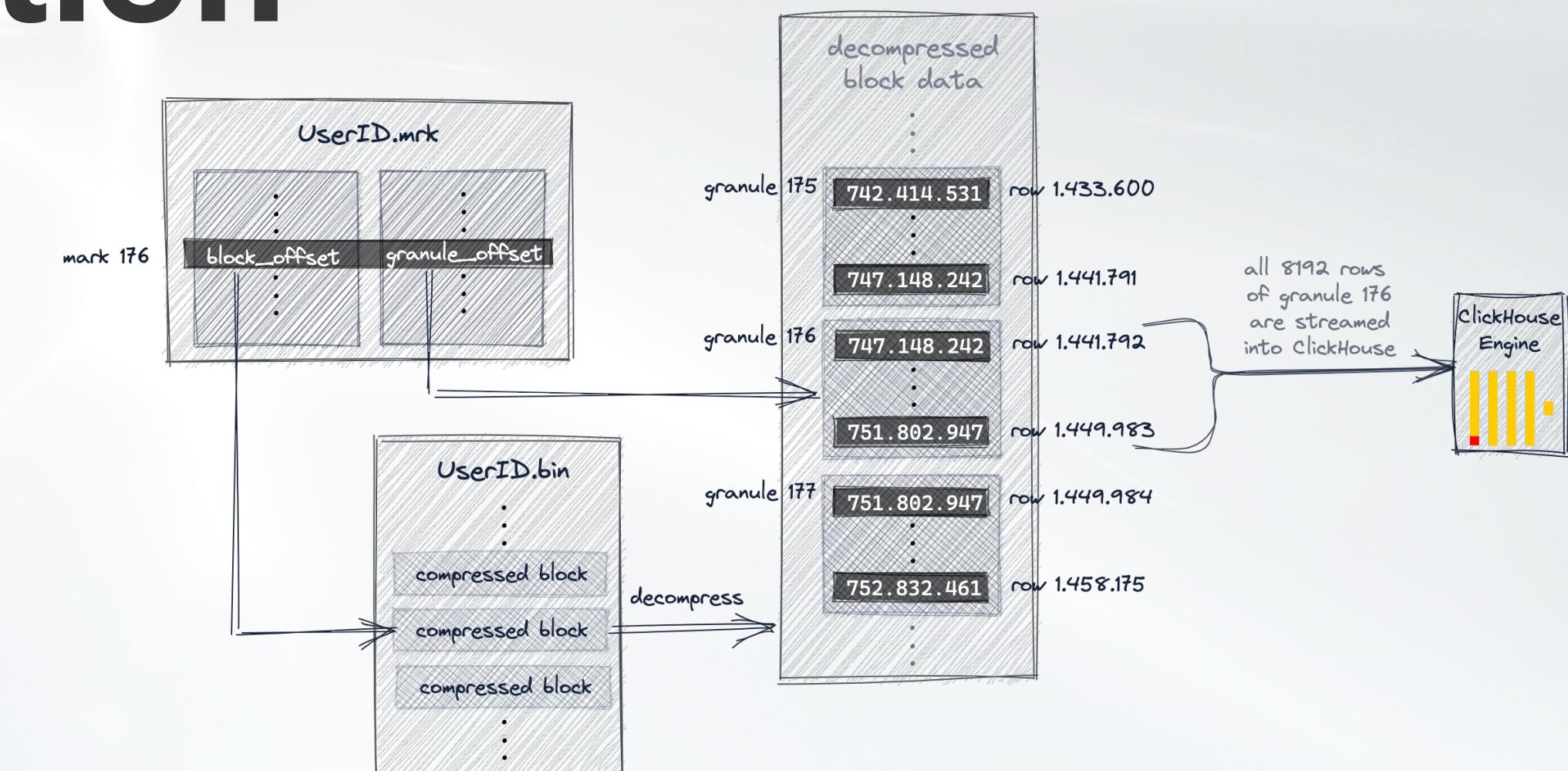
01 ClickHouse Files and Keys

02 ClickHouse Skipping Indexes

03 ClickHouse Explain

04 ClickHouse Optimization

05 Q&A





01 ClickHouse Files and Keys

1. ClickHouse Files and Keys

```

1 CREATE TABLE hits_UserID_URL
2 (
3     `UserID` UInt32,
4     `URL` String,
5     `EventTime` DateTime
6 )
7 ENGINE = MergeTree
8 PRIMARY KEY (UserID, URL)
9 ORDER BY (UserID, URL, EventTime)
10 SETTINGS index_granularity = 8192
11
12 Query id: f772b7c2-26fc-4b02-8ca8-45e281298bca
13
14 Ok.
15
16 0 rows in set. Elapsed: 0.007 sec.

```

```

1 INSERT INTO hits_UserID_URL SELECT
2     intHash32(UserID) AS UserID,
3     URL,
4     EventTime
5 FROM url('https://datasets.clickhouse.com/hits/tsv/hits_v1.tsv.xz', 'TSV', 'WatchID UInt64, JavaEnable UInt8,
6 WHERE URL != ''
7
8 Query id: a4eeba29-0f09-41f7-86d2-2a6494a4e035
9
10 Ok.
11
12 0 rows in set. Elapsed: 92.086 sec. Processed 8.87 million rows, 897.42 MB (96.37 thousand rows/s., 9.75 MB/s.)
13 Peak memory usage: 418.22 MiB.

```

```

1 v part_type: Wide
2 path: /var/lib/clickhouse/store/805/8058f914-1a43-4228-a48d-c99aae564a16/all_1_9_2/
3 rows: 8.87 million
4 data_uncompressed_bytes: 733.28 MiB
5 data_compressed_bytes: 206.92 MiB
6 primary_key_bytes_in_memory: 96.93 KiB
7 marks: 1083
8 bytes_on_disk: 206.98 MiB

```

Note:

1. 使用 Merge Tree 引擎;
2. Sorting Key 选择 (UserID, URL, EventTime) 3 列;
3. Primary Key 选择 (UserID, URL);
4. Primary Key 索引常驻内存，是 Sorting Key 的前缀；
5. 如果只指定了 sorting key，那么 primary key 跟 sorting key 相同；
6. Part 数据以 Wide 格式存储，每一列单独一个文件。

1. ClickHouse Files and Keys

```

1 [root@iZ2zecl1pe0cgqf3lyev5lZ all_1_9_2]# ll -lh
2 total 208M
3 -rw-r----- 1 clickhouse clickhouse 468 Aug 20 23:02 checksums.txt
4 -rw-r----- 1 clickhouse clickhouse 87 Aug 20 23:02 columns.txt
5 -rw-r----- 1 clickhouse clickhouse 7 Aug 20 23:02 count.txt
6 -rw-r----- 1 clickhouse clickhouse 10 Aug 20 23:02 default_compression_codec.txt
7 -rw-r----- 1 clickhouse clickhouse 33M Aug 20 23:02 EventTime.bin
8 -rw-r----- 1 clickhouse clickhouse 2.8K Aug 20 23:02 EventTime.cmrk
9 -rw-r----- 1 clickhouse clickhouse 1 Aug 20 23:02 metadata_version.txt
10 -rw-r----- 1 clickhouse clickhouse 48K Aug 20 23:02 primary.cidx
11 -rw-r----- 1 clickhouse clickhouse 238 Aug 20 23:02 serialization.json
12 -rw-r----- 1 clickhouse clickhouse 175M Aug 20 23:02 URL.bin
13 -rw-r----- 1 clickhouse clickhouse 3.8K Aug 20 23:02 URL.cmrk
14 -rw-r----- 1 clickhouse clickhouse 797K Aug 20 23:02 UserID.bin
15 -rw-r----- 1 clickhouse clickhouse 2.2K Aug 20 23:02 UserID.cmrk

```

Note:

- 最重要的两个文件 bin 文件和 mrk 文件。

```

1 [root@iZ2zecl1pe0cgqf3lyev5lZ all_1_9_2]# cat columns.txt
2 columns format version: 1
3 columns:
4 `UserID` UInt32
5 `URL` String
6 `EventTime` DateTime
7
8 [root@iZ2zecl1pe0cgqf3lyev5lZ all_1_9_2]# cat count.txt
9 8867681
10
11 [root@iZ2zecl1pe0cgqf3lyev5lZ all_1_9_2]# cat default_compression_codec.txt
12 CODEC(LZ4)
13
14 [root@iZ2zecl1pe0cgqf3lyev5lZ all_1_9_2]# cat serialization.json
15 {"columns": [{"kind": "Default", "name": "EventTime", "num_defaults": 0, "num_rows": 8867681},
   {"kind": "Default", "name": "URL", "num_defaults": 0, "num_rows": 8867681},
   {"kind": "Default", "name": "UserID", "num_defaults": 0, "num_rows": 8867681}], "version": 0}

```

1. ClickHouse Files and Keys – bin file

```

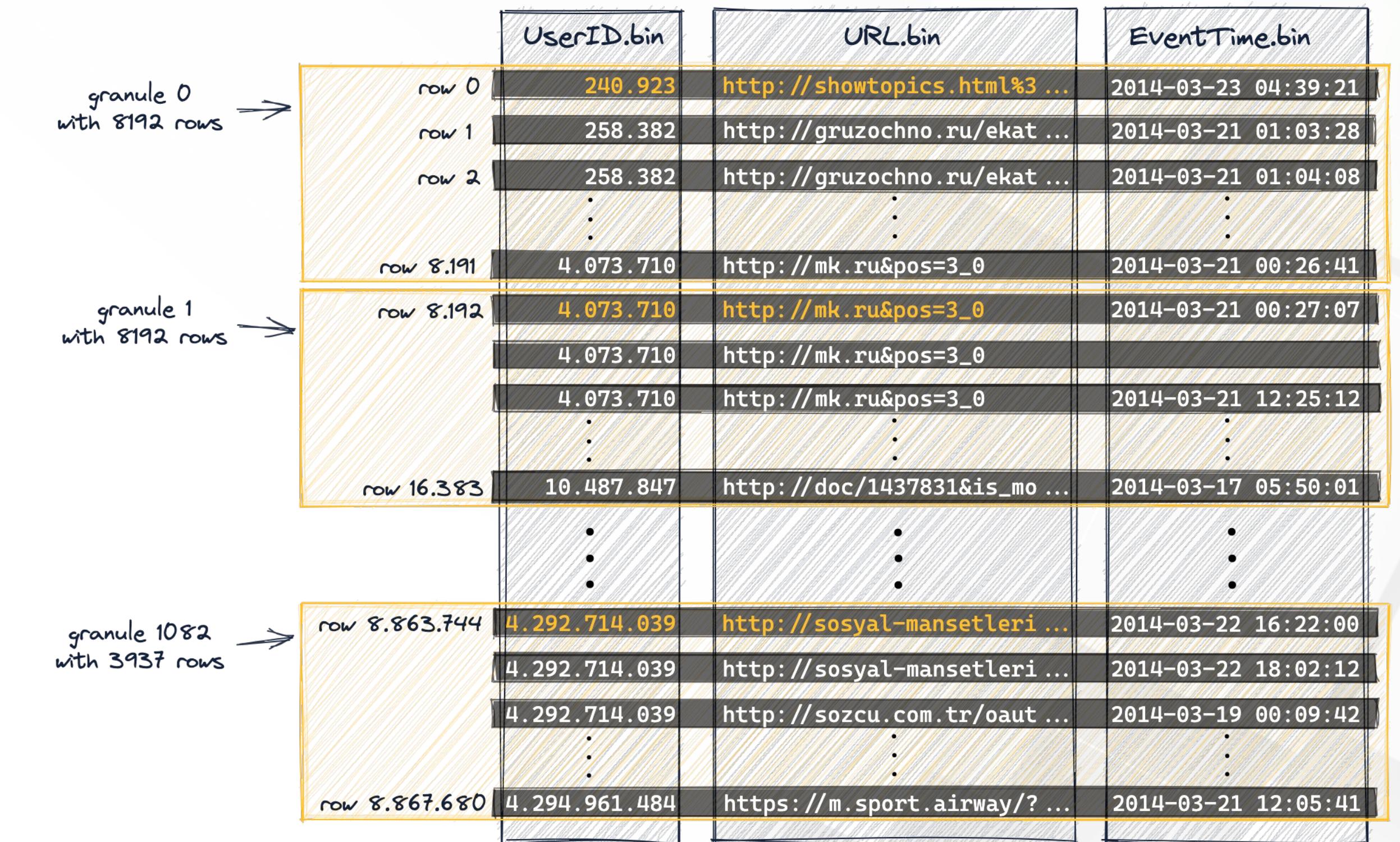
1 CREATE TABLE hits_UserID_URL
2 (
3   `UserID` UInt32,
4   `URL` String,
5   `EventTime` DateTime
6 )
7 ENGINE = MergeTree
8 PRIMARY KEY (UserID, URL)
9 ORDER BY (UserID, URL, EventTime)
10 SETTINGS index_granularity = 8192

```

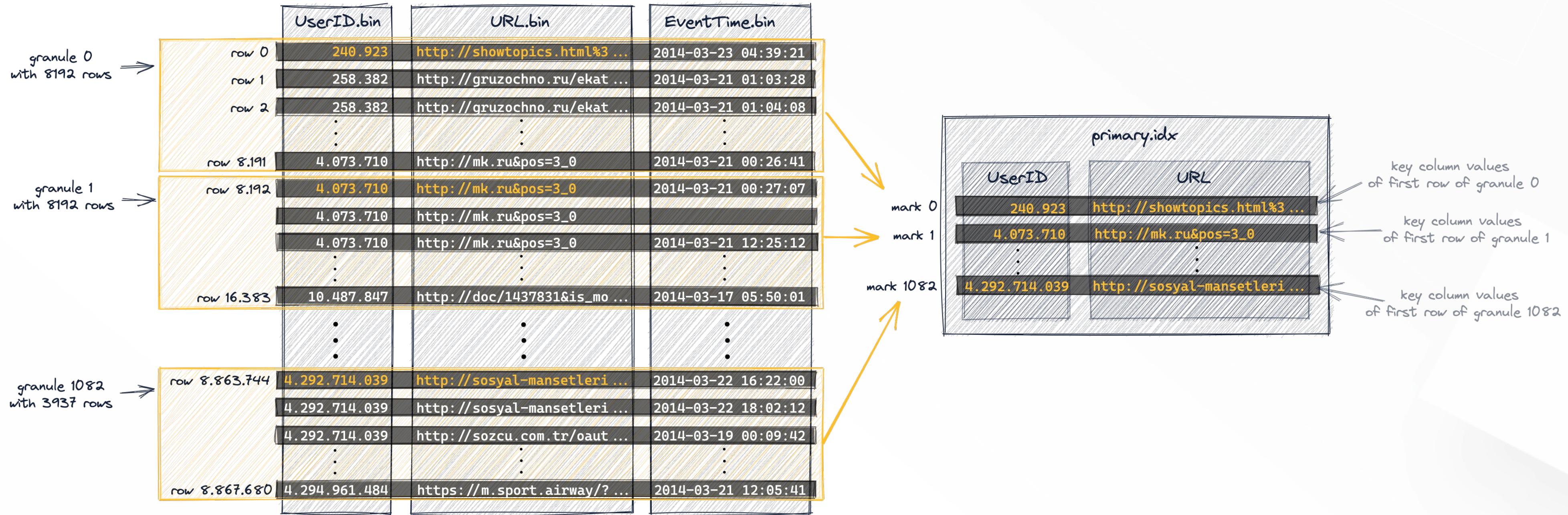
	UserID.bin	URL.bin	EventTime.bin
row 0	240.923	http://showtopics.html%3 ...	2014-03-23 04:39:21
row 1	258.382	http://gruzochno.ru/ekat ...	2014-03-21 01:03:28
row 2	258.382	http://gruzochno.ru/ekat ...	2014-03-21 01:04:08
⋮	⋮	⋮	⋮
row 8.191	4.073.710	http://mk.ru&pos=3_0	2014-03-21 00:26:41
row 8.192	4.073.710	http://mk.ru&pos=3_0	2014-03-21 00:27:07
	4.073.710	http://mk.ru&pos=3_0	
	4.073.710	http://mk.ru&pos=3_0	2014-03-21 12:25:12
⋮	⋮	⋮	⋮
row 16.383	10.487.847	http://doc/1437831&is_mo ...	2014-03-17 05:50:01
⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮
row 8.863.744	4.292.714.039	http://sosyal-mansetleri ...	2014-03-22 16:22:00
	4.292.714.039	http://sosyal-mansetleri ...	2014-03-22 18:02:12
	4.292.714.039	http://sozcu.com.tr/oaut ...	2014-03-19 00:09:42
⋮	⋮	⋮	⋮
row 8.867.680	4.294.961.484	https://m.sport.airway/? ...	2014-03-21 12:05:41

Note:

- 每一列在逻辑上被划分为多个 granule;
- granule 是 ClickHouse 数据处理的最小单位。

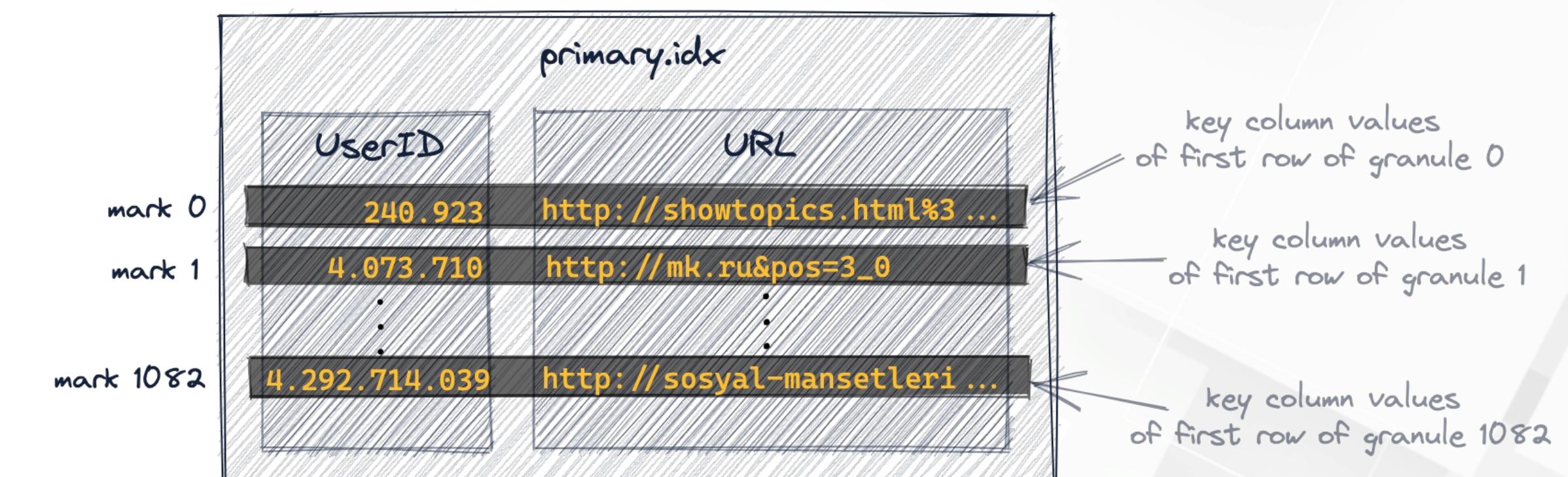


1. ClickHouse Files and Keys – idx file

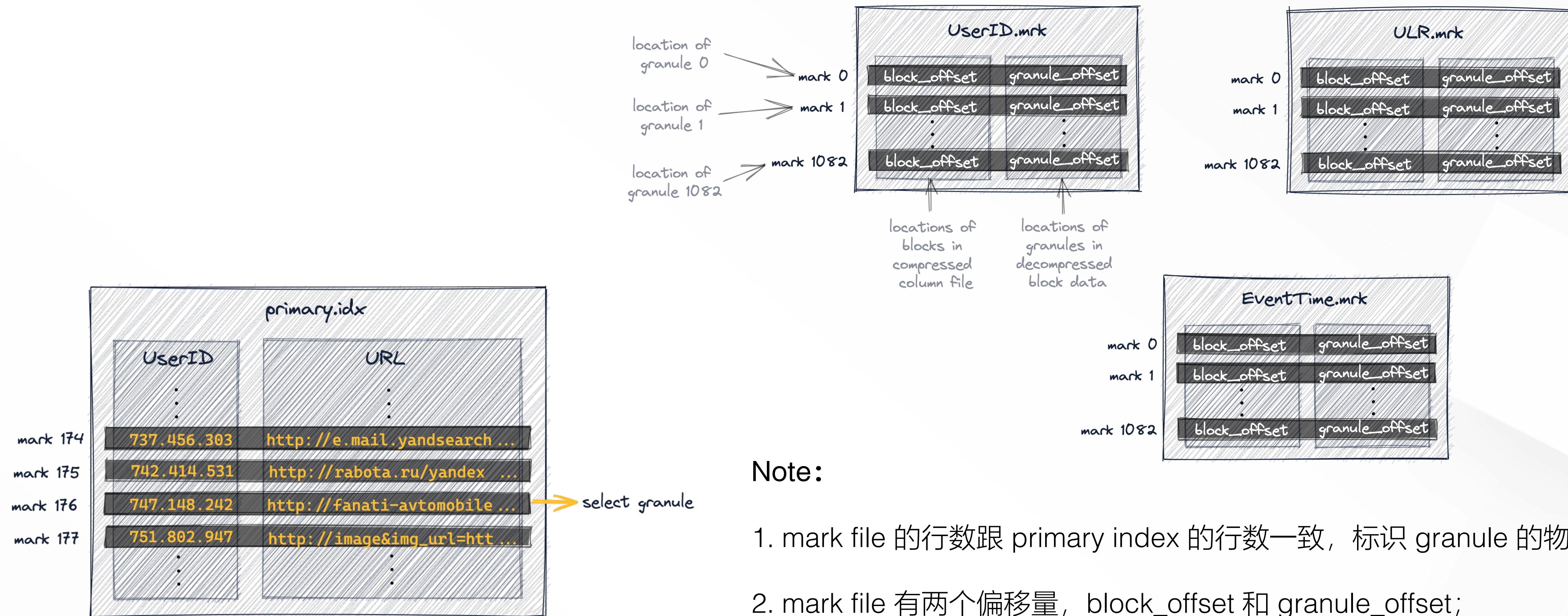


Note:

- 只有每个 granule 的第一行会被写到 primary.idx 中；
- 总共 887万行，分为 1083个 granule，主键就有 1083 行；
- primary.idx 完全加载到内存中，bytes: 96.93 KiB。



1. ClickHouse Files and Keys – mrk file



Note:

1. mark file 的行数跟 primary index 的行数一致，标识 granule 的物理位置；
2. mark file 有两个偏移量，block_offset 和 granule_offset；
3. bin 文件中，多个 granule 会被压缩成 1 个 block，当前数据的最大 block_number: 9；
4. block offset 是 bin 文件中 block 的偏移量； granule offset 是 block 中 granule 的偏移量；
5. mark file 不是常驻内存的，按需 load 到内存，marks_bytes: 8.61 KiB。

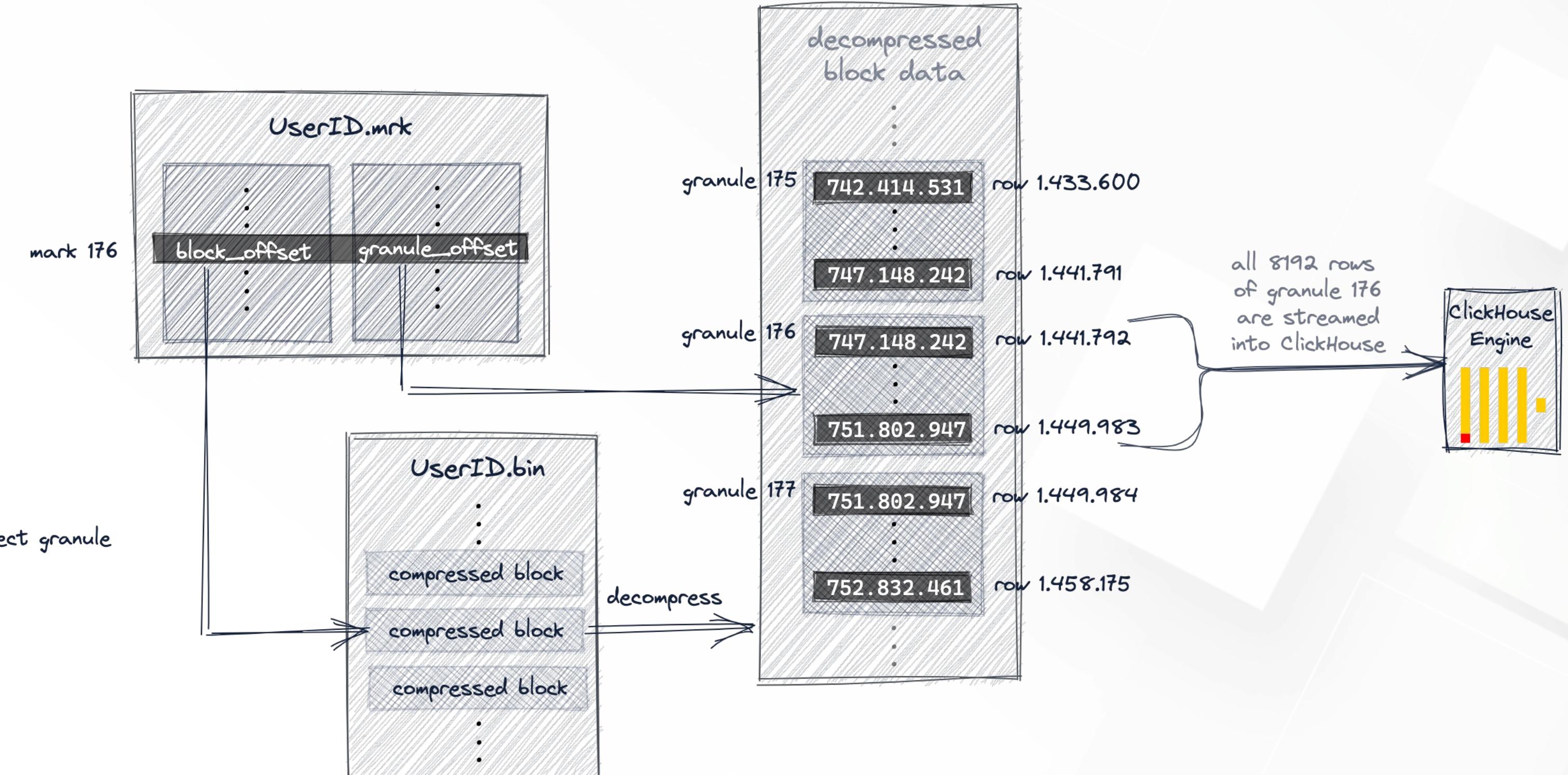
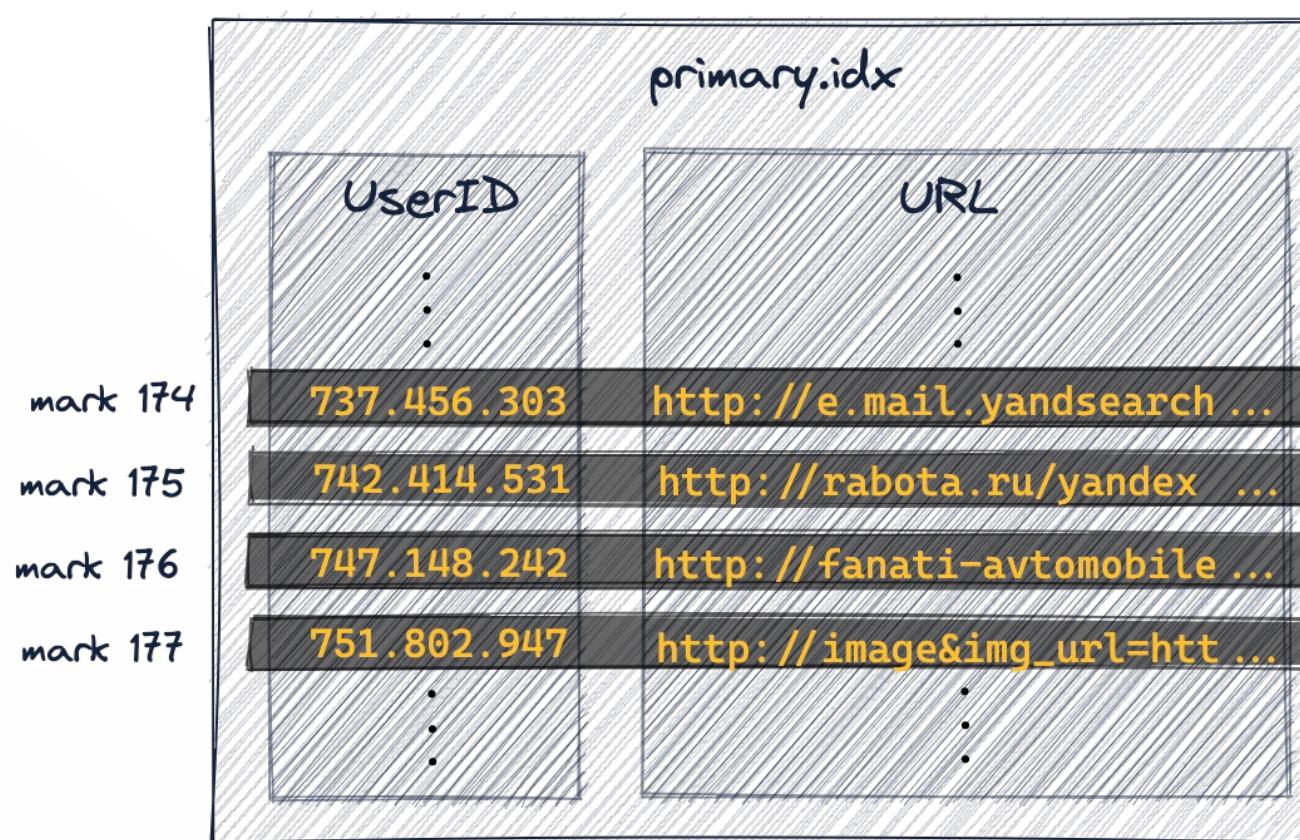
1. ClickHouse Files and Keys – 查询示例

```
1  SELECT
2      URL,
3      count(URL) AS Count
4  FROM hits_UserID_URL
5  WHERE UserID = 749927693
6  GROUP BY URL
7  ORDER BY Count DESC
8  LIMIT 10
9
10 Query id: 85706ec8-14cb-40b8-b851-85ef0de3f5eb
11
12   URL-----Count
13 1. http://auto.ru/chatay-barana... 170
14 2. http://auto.ru/chatay-id=371... 52
15 3. http://public_search           45
16 4. http://kovrik-medvedevushku... 36
17 5. http://forumal               33
18 6. http://korablitz.ru/L_10FFER... 14
19 7. http://auto.ru/chatay-id=371... 14
20 8. http://auto.ru/chatay-john-D... 13
21 9. http://auto.ru/chatay-john-D... 10
22 10. http://wot/html?page/23600_m... 9
23
24
25 10 rows in set. Elapsed: 0.007 sec. Processed 8.19 thousand rows, 740.18 KB (1.10 million rows/s., 99.05 MB/s.)
26 Peak memory usage: 243.55 KiB.
```

1. ClickHouse Files and Keys – 查询示例

```

1 SELECT
2   URL,
3   count(URL) AS Count
4 FROM hits_UserID_URL
5 WHERE UserID = 749927693
6 GROUP BY URL
7 ORDER BY Count DESC
8 LIMIT 10
  
```



```

6 8-a48d-c99aae564a16) (SelectExecutor): Key condition: (column 0 in [749927693, 749927693])
7 8-a48d-c99aae564a16) (SelectExecutor): Filtering marks by primary and secondary keys
8 8-a48d-c99aae564a16) (SelectExecutor): Running binary search on index range for part all_1_9_2 (1083 marks)
9 8-a48d-c99aae564a16) (SelectExecutor): Found (LEFT) boundary mark: 176
10 8-a48d-c99aae564a16) (SelectExecutor): Found (RIGHT) boundary mark: 177
11 8-a48d-c99aae564a16) (SelectExecutor): Found continuous range in 19 steps
12 8-a48d-c99aae564a16) (SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1/1083 marks by primary key, 1 marks to read from 1 ranges
13 8-a48d-c99aae564a16) (SelectExecutor): Spreading mark ranges among streams (default reading)
14 8-a48d-c99aae564a16) (SelectExecutor): Reading 1 ranges in order from part all_1_9_2, approx. 8192 rows starting from 1441792
  
```

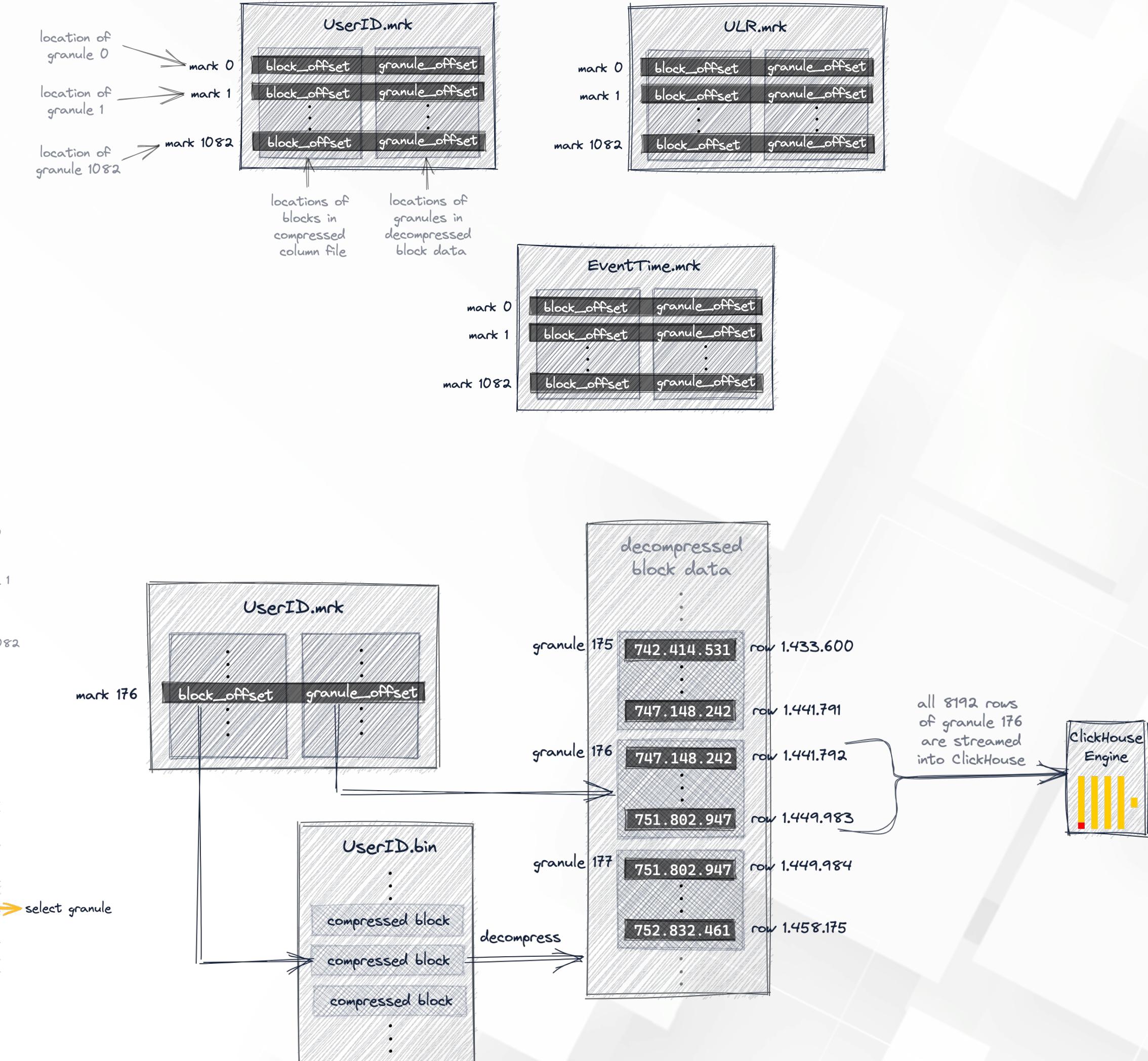
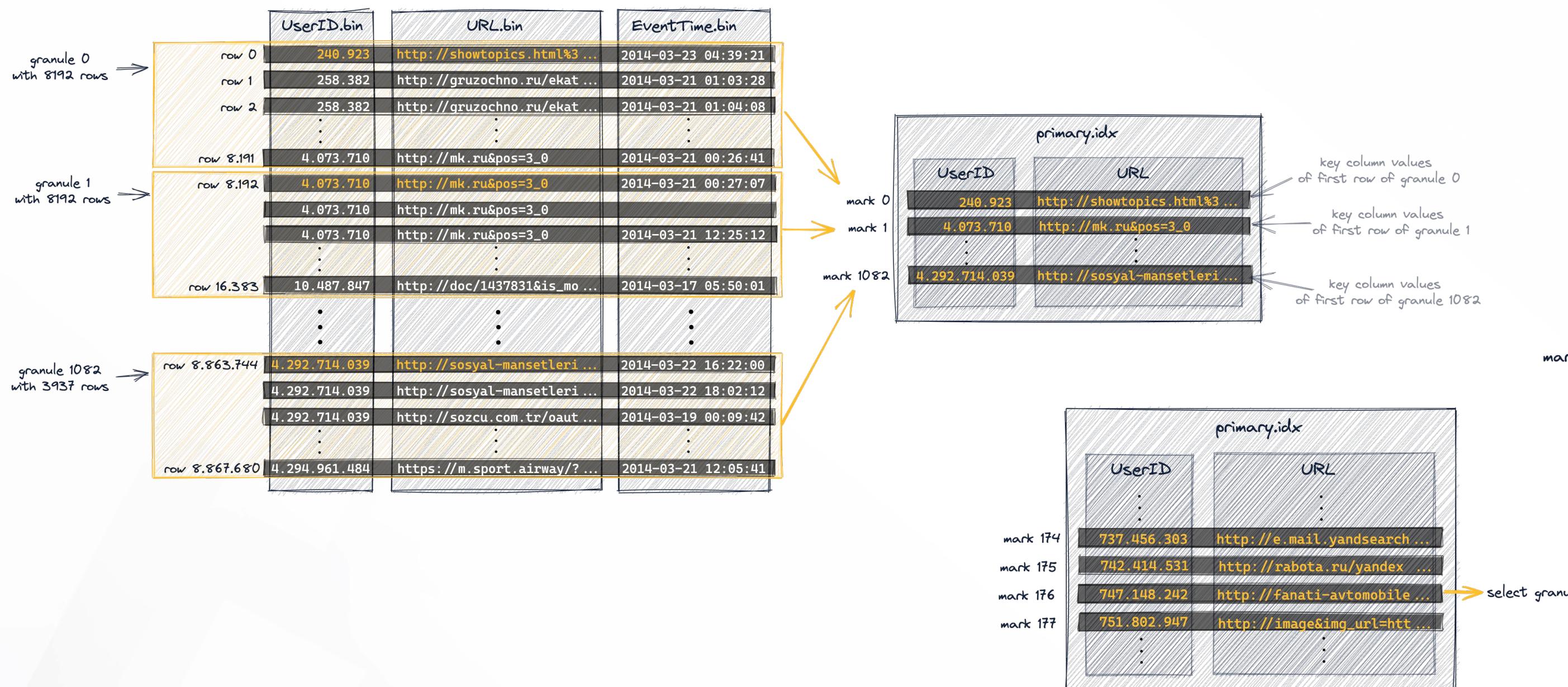
1. ClickHouse Files and Keys 小结



```

1 [root@iZ2zecl1pe0cgqf3lyev5lZ all_1_9_2]# ll -lh
2 total 208M
3 -rw-r----- 1 clickhouse clickhouse 468 Aug 20 23:02 checksums.txt
4 -rw-r----- 1 clickhouse clickhouse 87 Aug 20 23:02 columns.txt
5 -rw-r----- 1 clickhouse clickhouse 7 Aug 20 23:02 count.txt
6 -rw-r----- 1 clickhouse clickhouse 10 Aug 20 23:02 default_compression_codec.txt
7 -rw-r----- 1 clickhouse clickhouse 33M Aug 20 23:02 EventTime.bin
8 -rw-r----- 1 clickhouse clickhouse 2.8K Aug 20 23:02 EventTime.cmrk
9 -rw-r----- 1 clickhouse clickhouse 1 Aug 20 23:02 metadata_version.txt
10 -rw-r----- 1 clickhouse clickhouse 48K Aug 20 23:02 primary.cidx
11 -rw-r----- 1 clickhouse clickhouse 238 Aug 20 23:02 serialization.json
12 -rw-r----- 1 clickhouse clickhouse 175M Aug 20 23:02 URL.bin
13 -rw-r----- 1 clickhouse clickhouse 3.8K Aug 20 23:02 URL.cmrk
14 -rw-r----- 1 clickhouse clickhouse 797K Aug 20 23:02 UserID.bin
15 -rw-r----- 1 clickhouse clickhouse 2.2K Aug 20 23:02 UserID.cmrk

```



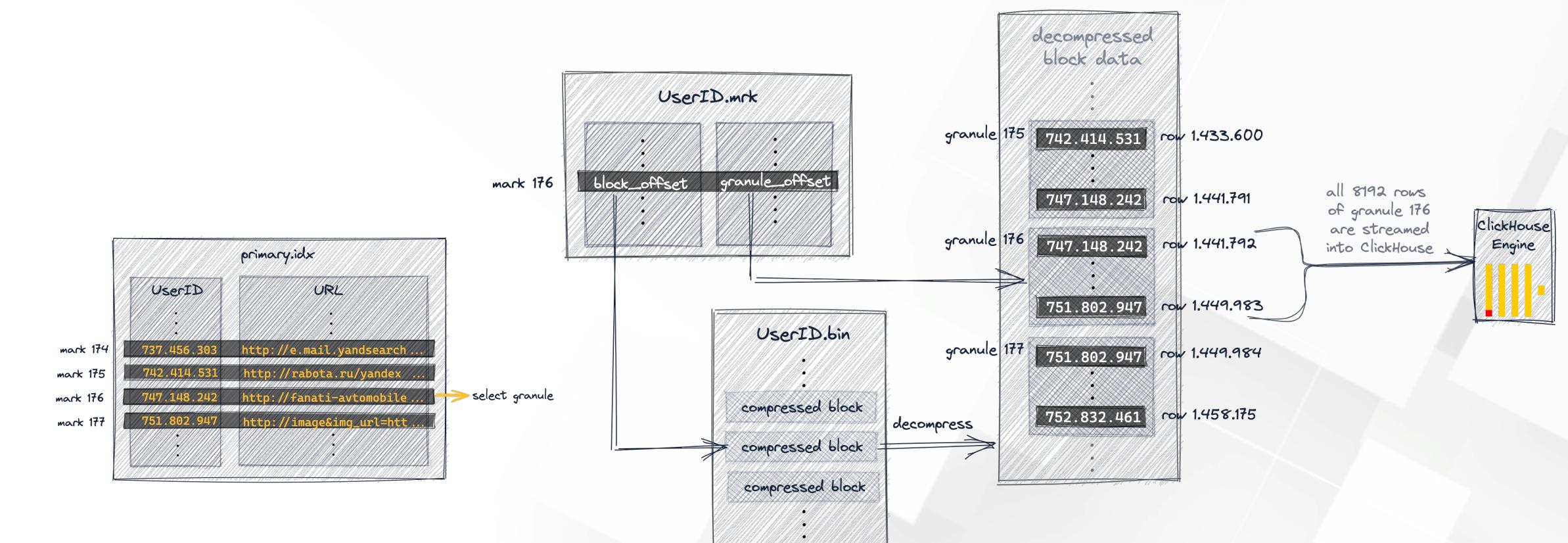
1. ClickHouse Files and Keys 小结

排序建、主键、分区键的一般选择方式：

1. Sorting Key 是 ClickHouse 最主要依赖的扫描加速技术，Part 内的数据按 Sorting Key 严格有序，查询使用二分查找；
2. 按照业务场景，将查询最频繁 SQL 命中的那些列，按照基数升序放在 Primary Key 中，取得相对较好的压缩率和性能；
3. Primary Key 是 Sorting key 的前缀，为了提高内存使用效率，Primary Key 只包含查询要过滤的列；
4. 因为相同的数据放在一起能够得到更好的压缩率，所以 Sorting Key 中，除了 Primary Key 还能放基数低的列，降低磁盘使用率；
5. Partition Key 定义了不同分区，不同分区的数据是在物理上隔离开的，每个 Part 只属于一个分区；
6. 由于 Partition Key 会对 Part 起到物理隔离，所以数据分区过细的情况下会导致 Part 数量膨胀，文件IO一多，就可能影响查询性能。

Question:

Primary Key 查找是很快，但是如果有些 SQL 无法命中 Primary Key 怎么办？





02 ClickHouse Skipping Indexes

2. Skipping Indexes

```

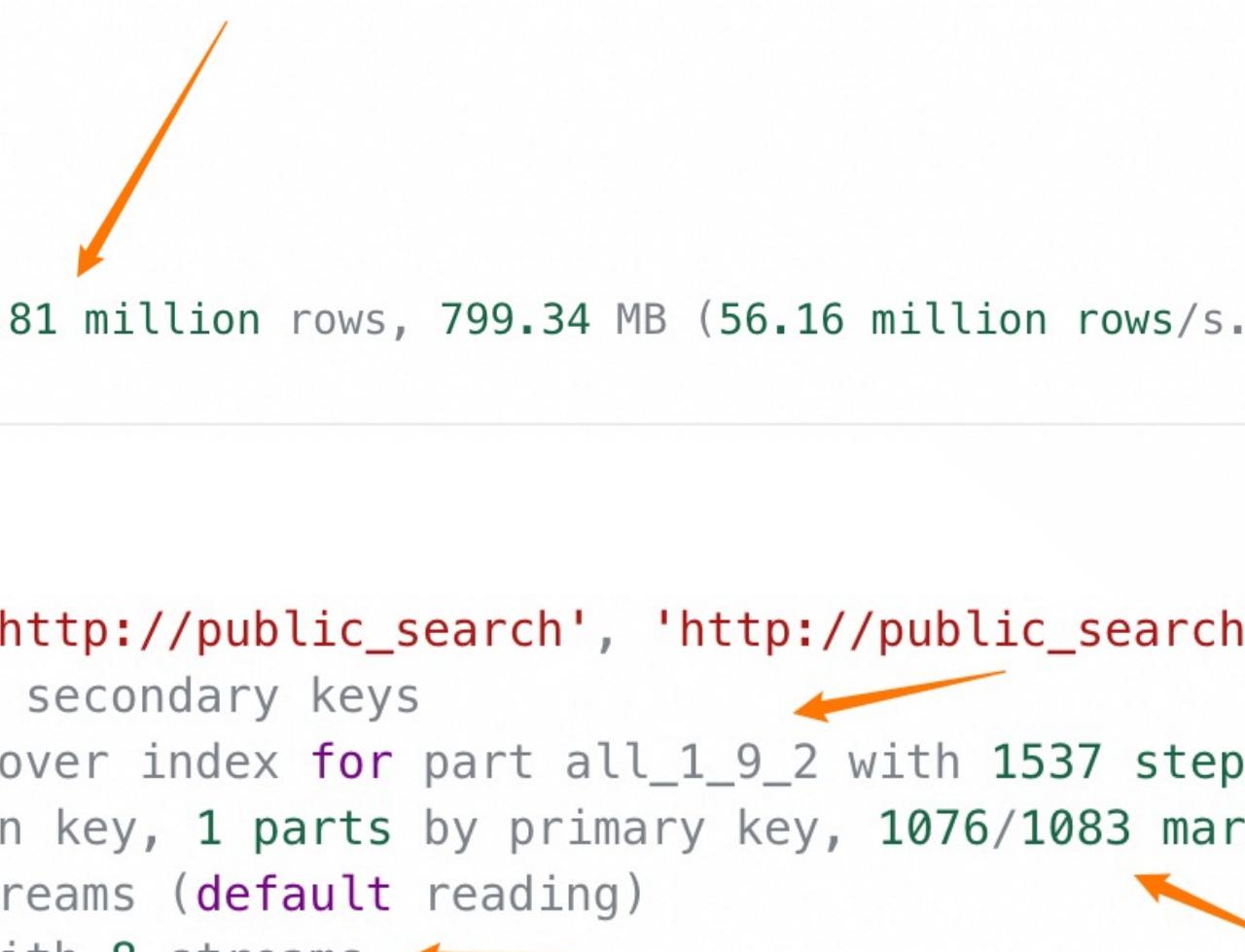
1  SELECT
2    UserID,
3    count(UserID) AS Count
4  FROM hits_UserID_URL
5  WHERE URL = 'http://public_search'
6  GROUP BY UserID
7  ORDER BY Count DESC
8  LIMIT 10
9
10 Query id: 4aed4fe5-4d36-47c0-8ee0-e757d2db87af
11
12    +-----+-----+
13    | UserID | Count |
14    +-----+-----+
15    | 2459550954 | 3741 |
16    | 1084649151 | 2484 |
17    | 723361875 | 729 |
18    | 3087145896 | 695 |
19    | 2754931092 | 672 |
20    | 1509037307 | 582 |
21    | 3085460200 | 573 |
22    | 2454360090 | 556 |
23    | 3884990840 | 539 |
24    | 765730816 | 536 |
25  10 rows in set. Elapsed: 0.157 sec. Processed 8.81 million rows, 799.34 MB (56.16 million rows/s., 5.10 GB/s.)
26  Peak memory usage: 25.27 MiB.

```

```

(SelectExecutor): Key condition: (column 1 in ['http://public_search', 'http://public_search'])
(SelectExecutor): Filtering marks by primary and secondary keys
(SelectExecutor): Used generic exclusion search over index for part all_1_9_2 with 1537 steps
(SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1076/1083 marks by primary key, 1076 marks to read from 5 ranges
(SelectExecutor): Spreading mark ranges among streams (default reading)
(SelectExecutor): Reading approx. 8810337 rows with 8 streams

```



2. Skipping Indexes

```

1  SELECT
2    UserID,
3    count(UserID) AS Count
4  FROM hits_UserID_URL
5  WHERE URL = 'http://public_search'
6  GROUP BY UserID
7  ORDER BY Count DESC
8  LIMIT 10
9
10 Query id: 4aed4fe5-4d36-47c0-8ee0-e757d2db87af
11
12   +-----+-----+
13   | UserID | Count |
14   +-----+-----+
15   | 2459550954 | 3741 |
16   | 1084649151 | 2484 |
17   | 723361875 | 729 |
18   | 3087145896 | 695 |
19   | 2754931092 | 672 |
20   | 1509037307 | 582 |
21   | 3085460200 | 573 |
22   | 2454360090 | 556 |
23   | 3884990840 | 539 |
24   | 765730816 | 536 |
25
26 10 rows in set. Elapsed: 0.157 sec. Processed 8.81 million rows, 799.34 MB (56.16 million rows/s., 5.10 GB/s.)
27 Peak memory usage: 25.27 MiB.

```

Question:

1. 从日志来看，为什么 SQL 基本上执行全表扫？
2. 被过滤掉的那几个 marks 是什么情况？
3. 怎么加速这种查询？

```

(SelectExecutor): Key condition: (column 1 in ['http://public_search', 'http://public_search'])
(SelectExecutor): Filtering marks by primary and secondary keys
(SelectExecutor): Used generic exclusion search over index for part all_1_9_2 with 1537 steps
(SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1076/1083 marks by primary key, 1076 marks to read from 5 ranges
(SelectExecutor): Spreading mark ranges among streams (default reading)
(SelectExecutor): Reading approx. 8810337 rows with 8 streams

```

2. Skipping Indexes



```
1 SELECT
2     UserID,
3     count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE URL = 'http://public_search'
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
```

	UserID.bin	URL.bin	EventTime.bin
granule 0 with 8192 rows	row 0 240.923	http://showtopics.html%3 ...	2014-03-23 04:39:21
	row 1 258.382	http://gruzochno.ru/ekat ...	2014-03-21 01:03:28
	row 2 258.382	http://gruzochno.ru/ekat ...	2014-03-21 01:04:08
	⋮	⋮	⋮
	row 8.191 4.073.710	http://mk.ru&pos=3_0	2014-03-21 00:26:41
granule 1 with 8192 rows	row 8.192 4.073.710	http://mk.ru&pos=3_0	2014-03-21 00:27:07
	4.073.710	http://mk.ru&pos=3_0	
	4.073.710	http://mk.ru&pos=3_0	2014-03-21 12:25:12
	⋮	⋮	⋮
	row 16.383 10.487.847	http://doc/1437831&is_mo ...	2014-03-17 05:50:01
	⋮	⋮	⋮
granule 1082 with 3937 rows	row 8.863.744 4.292.714.039	http://sosyal-mansetleri ...	2014-03-22 16:22:00
	4.292.714.039	http://sosyal-mansetleri ...	2014-03-22 18:02:12
	4.292.714.039	http://sozcu.com.tr/oaut ...	2014-03-19 00:09:42
	⋮	⋮	⋮
	row 8.867.680 4.294.961.484	https://m.sport.airway/? ...	2014-03-21 12:05:41

Question:

Q: 从日志来看, 为什么 SQL 基本上执行全表扫?
A: 因为 URL 这一列没有排序。

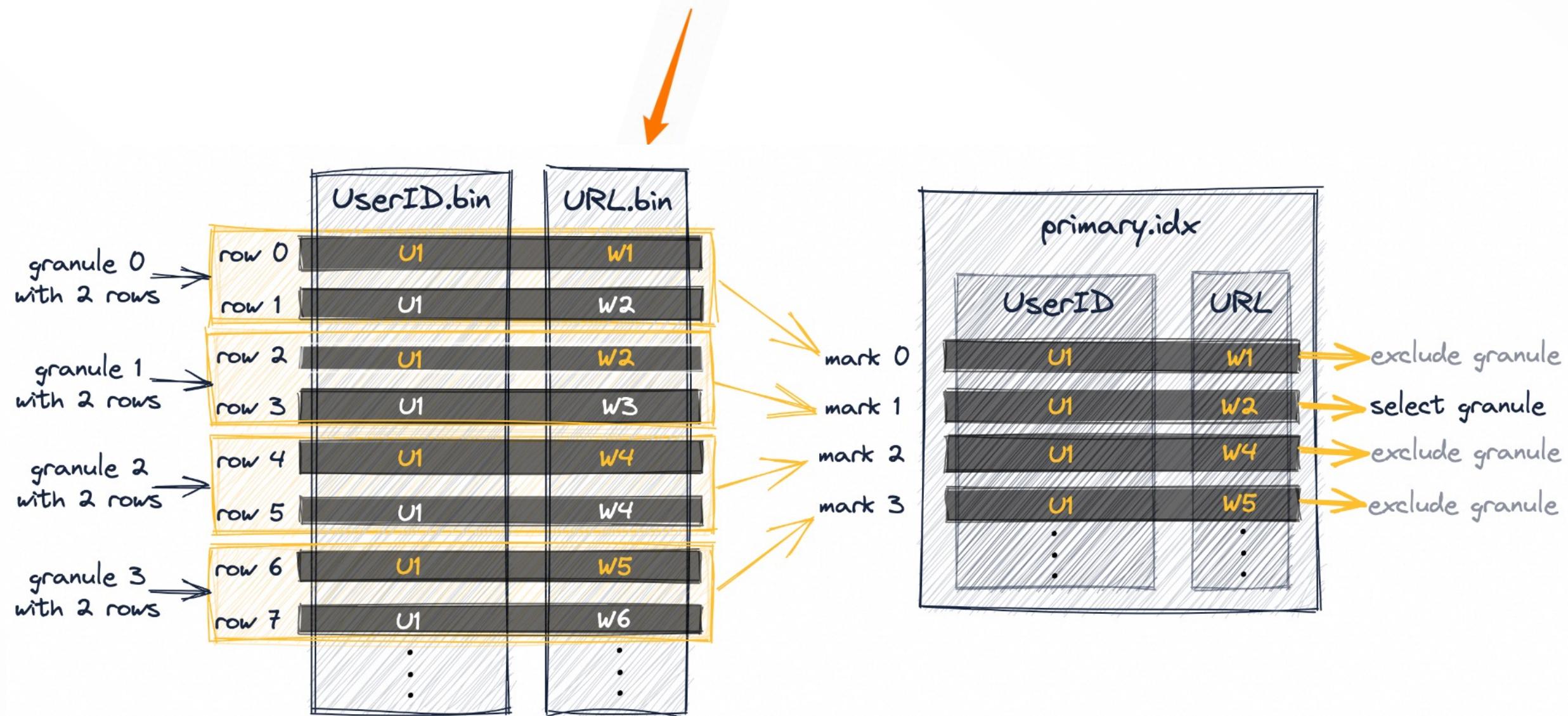
Question:

- 1. 从日志来看，为什么 SQL 基本上执行全表扫？
 - 2. 被过滤掉的那几个 marks 是什么情况？
 - 3. 怎么加速这种查询？

2. Skipping Indexes

```

1 SELECT
2   UserID,
3   count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE URL = 'http://public_search'
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
  
```



Question:

Q: 被过滤掉的那几个 marks 是什么情况?

A: 因为在 UserID 相同的时候, URL 局部有序。

Question:

1. 从日志来看, 为什么 SQL 基本上执行全表扫?

2. 被过滤掉的那几个 marks 是什么情况?

3. 怎么加速这种查询?



2. Skipping Indexes

```

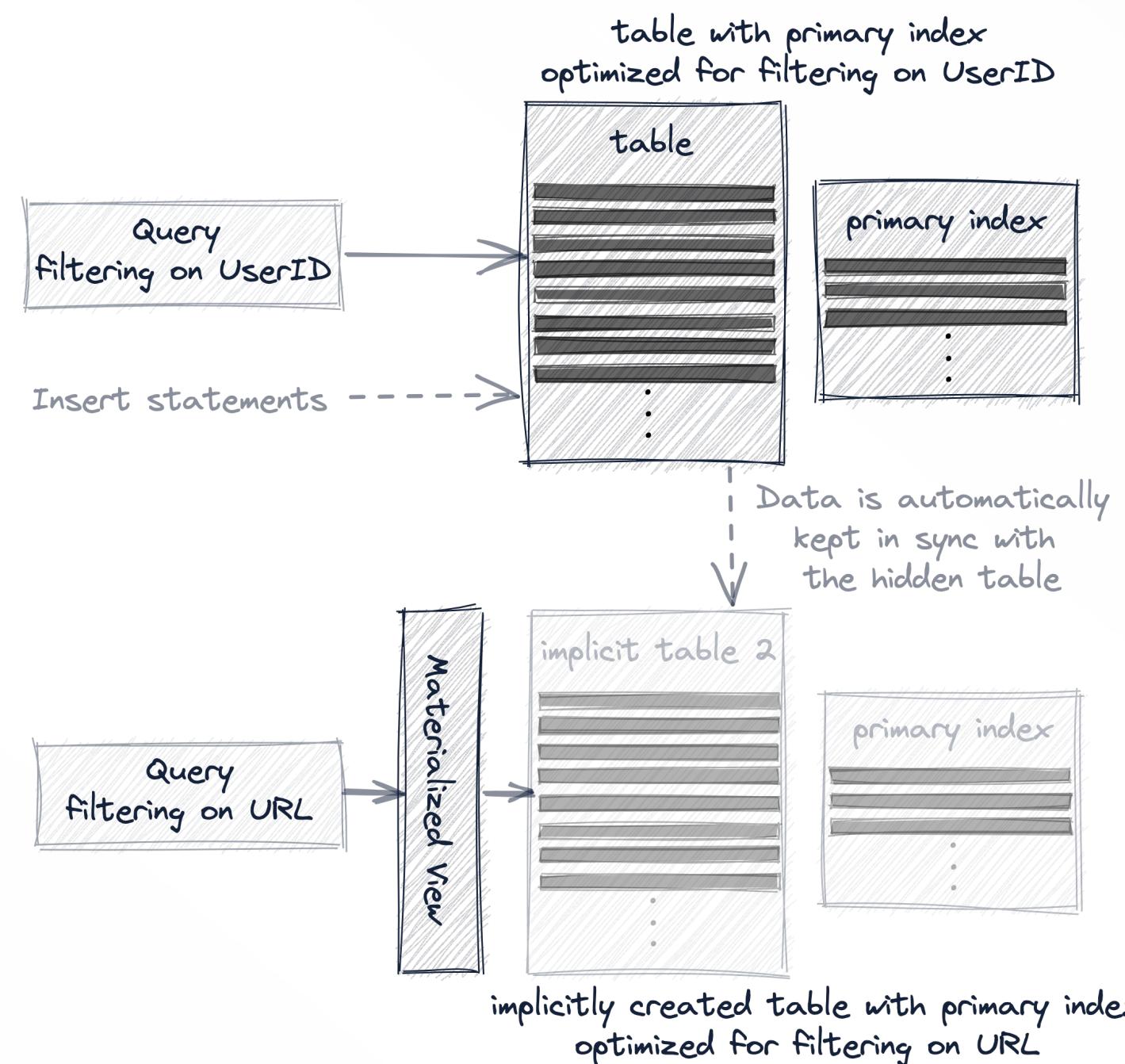
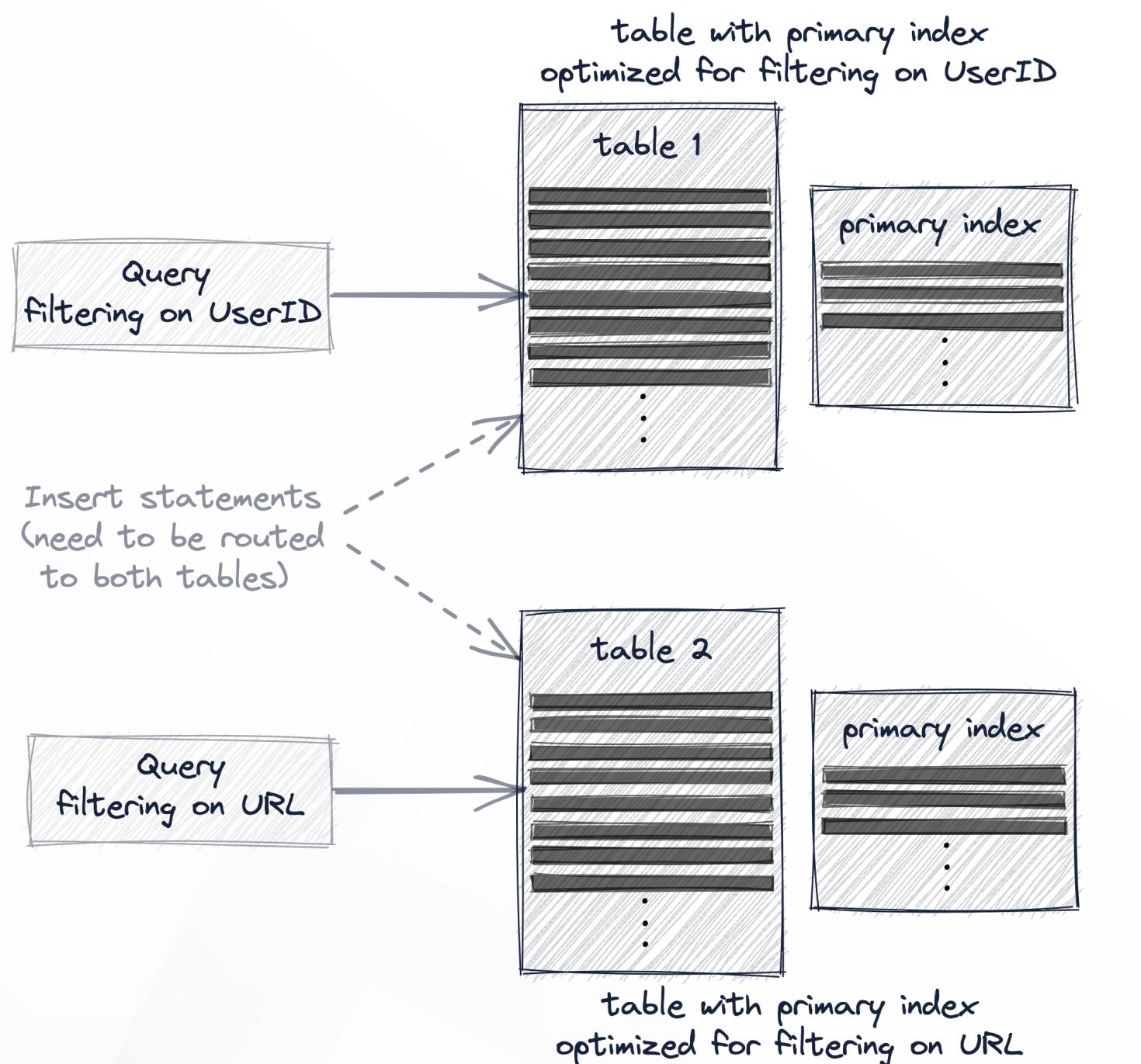
1 SELECT
2   UserID,
3   count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE URL = 'http://public_search'
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
  
```

Question:

Q: 怎么加速这种查询?

A: Secondary Table、MV、Projection、Skipping Index。

给另外的一份数据重新设置 Sorting Key，让 SQL 命中排好序的列。



Secondary Table

Materialized View

Question:

1. 从日志来看，为什么 SQL 基本上执行全表扫？

2. 被过滤掉的那几个 marks 是什么情况？

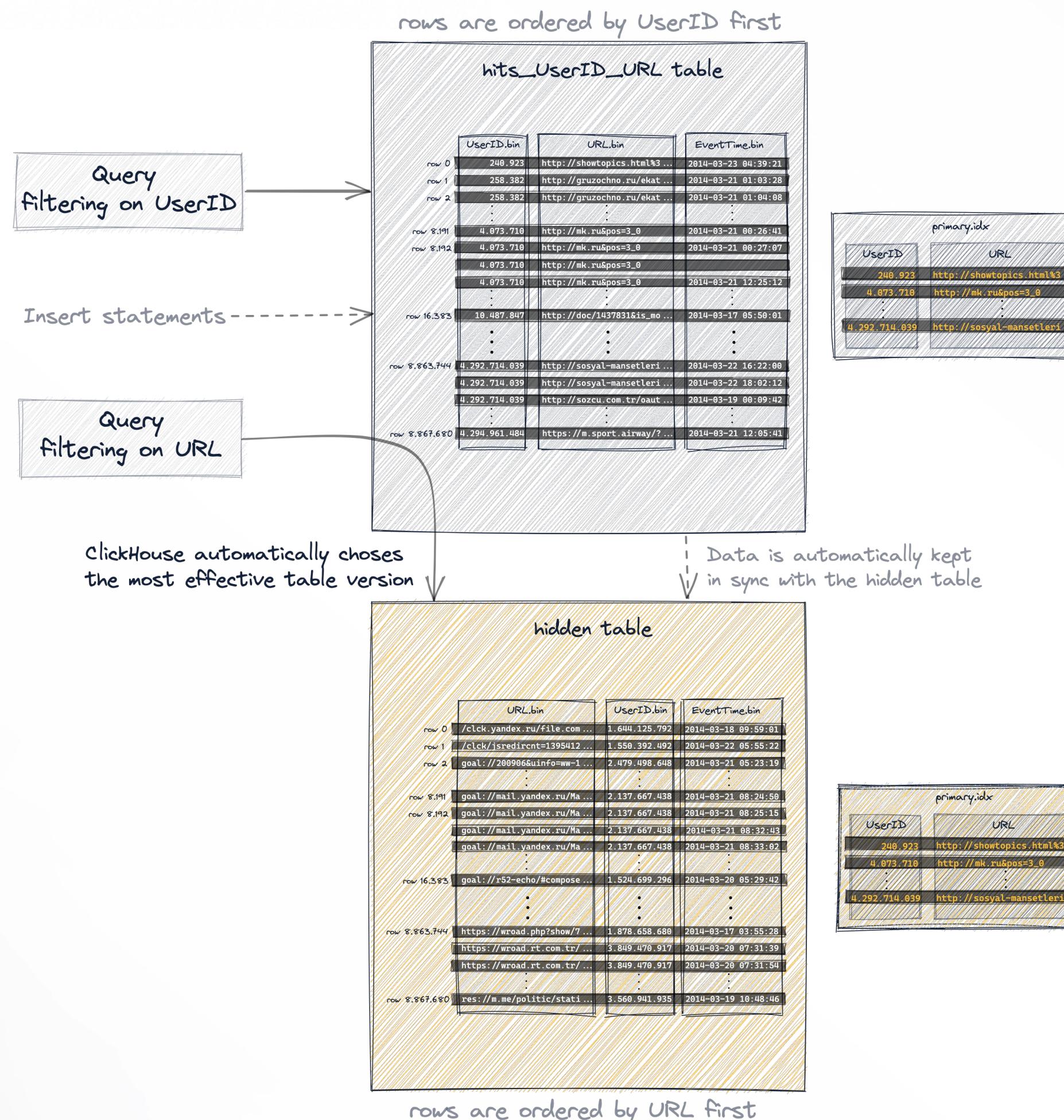
3. 怎么加速这种查询？



2. Skipping Indexes

```

1 SELECT
2   UserID,
3   count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE URL = 'http://public_search'
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
  
```



Projection

Question:

Q: 怎么加速这种查询?

A: Secondary Table、MV、Projection、Skipping Index。

给另外的一份数据重新设置 Sorting Key，让 SQL 命中排好序的列。

Question:

1. 从日志来看，为什么 SQL 基本上执行全表扫？

2. 被过滤掉的那几个 marks 是什么情况？

3. 怎么加速这种查询？

2. Skipping Indexes



```
1 SELECT
2     UserID,
3     count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE URL = 'http://public_search'
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
```

Question

Q: 怎么加速这种查询？

A: Secondary Table、MV、Projection、**Skipping Index**。

数据不冗余的情况下，可能可行的解法。

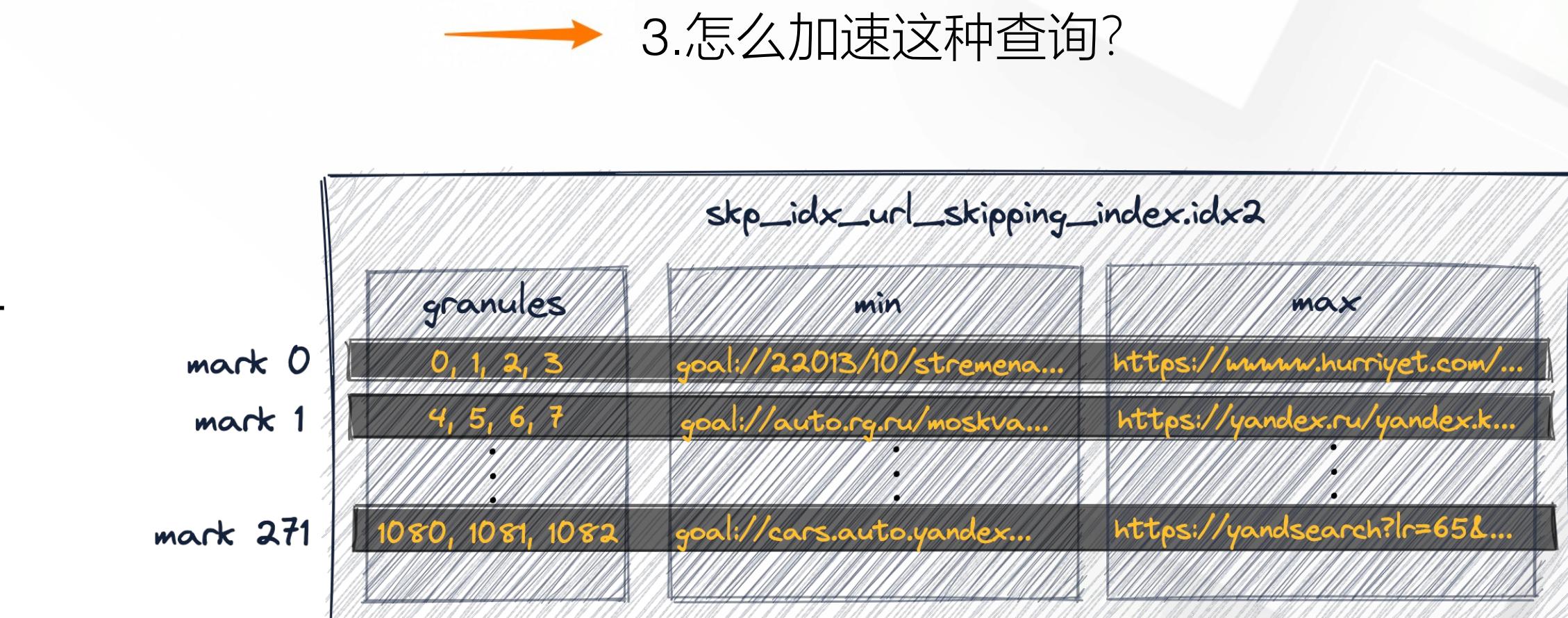
```
1 ALTER TABLE hits_UserID_URL ADD INDEX url_skipping_index URL TYPE minmax GRANULARITY 4  
2 ALTER TABLE hits_UserID_URL MATERIALIZE INDEX url_skipping_index;
```

Question:

1.从日志来看，为什么 SQL 基本上执行全表扫？

2. 被过滤掉的那几个 marks 是什么情况？

3. 怎么加速这种查询？



2. Skipping Indexes

```

1  SELECT
2      UserID,
3      count(UserID) AS Count
4  FROM hits_UserID_URL
5  WHERE URL = 'http://public_search'
6  GROUP BY UserID
7  ORDER BY Count DESC
8  LIMIT 10

```

Question:

Q: 怎么加速这种查询?

A: Secondary Table、MV、Projection、**Skipping Index**。

数据不冗余的情况下，可能可行的解法。

```

1  ALTER TABLE hits_UserID_URL ADD INDEX url_skipping_index URL TYPE minmax GRANULARITY 4;
2  ALTER TABLE hits_UserID_URL MATERIALIZE INDEX url_skipping_index;

```

```

1  v SELECT
2      UserID,
3      count(UserID) AS Count
4  FROM hits_UserID_URL
5  WHERE URL = 'http://public_search'
6  GROUP BY UserID
7  ORDER BY Count DESC
8  LIMIT 10
9
10 Query id: 49fc94ff-8876-4b4d-b67b-bcbbc3a8f89b
11
12 1. UserID Count
13 1. 2459550954 3741
14 2. 1084649151 2484
15 3. 723361875 729
16 4. 3087145896 695
17 5. 2754931092 672
18 6. 1509037307 582
19 7. 3085460200 573
20 8. 2454360090 556
21 9. 3884990840 539
22 10. 765730816 536
23
24
25 10 rows in set. Elapsed: 0.179 sec. Processed 8.81 million rows, 799.34 MB (49.21 million rows/s., 4.46 GB/s.)
26 Peak memory usage: 35.30 MiB.

```

Note:

- 选择 URL 值可能介于 skipping index 为每组 granules 存储的最小值和最大值之间，导致这个 skipping index 没有起到过滤效果，甚至在数据量大的情况下，扫描 Skipping index 带来的开销，使得 SQL 变得“更慢”了。
- 一般来说，当主键和该列存在强相关，且查找的值在很多 granule 中都不会出现的时候，用 Skipping Index 就能得到较好的过滤效果。

Question:

1. 从日志来看，为什么 SQL 基本上执行全表扫?

2. 被过滤掉的那几个 marks 是什么情况?

3. 怎么加速这种查询?



2. Skipping Indexes - types



```
1 INDEX index_name expr TYPE type(...) [GRANULARITY granularity_value]
```

```
1 ✓ MinMax — 存储 granules 内的最大最小值  
2  
3 Syntax: minmax
```

```
1 ✓ Set — 存储指定表达式的唯一值(不超过 max_rows 行, max_rows=0 表示“无限制”)  
2  
3 Syntax: set(max_rows)
```

```
1 ✓ Bloom Filter — false_positive 默认 0.025.  
2 Supported data types: Int*, UInt*, Float*, Enum, Date, DateTime, String, FixedString, Array,  
LowCardinality, Nullable, UUID and Map.  
3  
4 Syntax: bloom_filter([false_positive])
```

```
1 ✓ N-gram Bloom Filter  
2 Only works with datatypes: String, FixedString and Map.  
3  
4 Syntax: ngrambf_v1(n, size_of_bloom_filter_in_bytes, number_of_hash_functions, random_seed)
```

```
1 ✓ Token Bloom Filter  
2 The same as ngrambf_v1, but stores tokens instead of ngrams. Tokens are sequences separated by non-  
alphanumeric characters.  
3  
4 Syntax: tokenbf_v1(size_of_bloom_filter_in_bytes, number_of_hash_functions, random_seed)
```

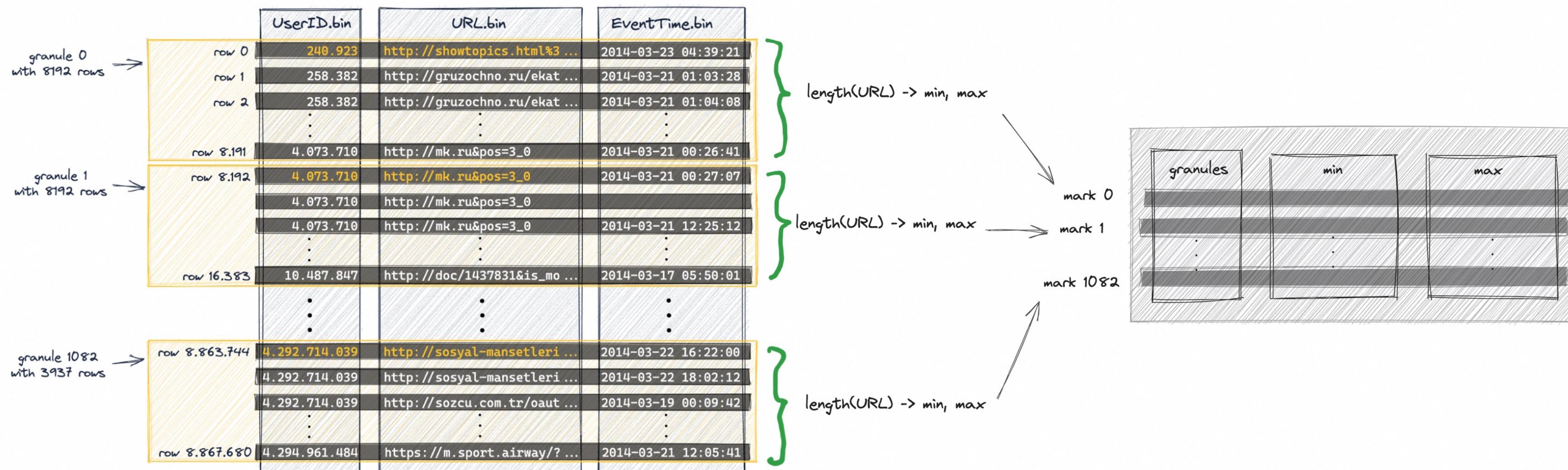
2. Skipping Indexes - minmax

```

1 ✓ ALTER TABLE hits_UserID_URL ADD INDEX url_len_minmax length(URL) TYPE minmax GRANULARITY 1
2
3 ALTER TABLE hits_UserID_URL MATERIALIZE INDEX url_len_minmax;
  
```

```

1 ✓ SELECT
2     UserID,
3     count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE length(URL) > 1000
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
  
```



Note:

1. minmax 适用于局部区间内的极值，比如日期、数量、长度等

```

1 ✓ aae564a16) (SelectExecutor): Filtering marks by primary and secondary keys
2 aae564a16) (SelectExecutor): Index `url_skipping_index_len_minmax` has dropped 166/1083 granules.
3 aae564a16) (SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1083/1083 marks by primary key, 917 marks to read from 139 ranges
4 aae564a16) (SelectExecutor): Spreading mark ranges among streams (default reading)
5 aae564a16) (SelectExecutor): Reading approx. 7507809 rows with 8 streams
  
```

2. Skipping Indexes - set

```

1 ✓ ALTER TABLE hits_UserID_URL ADD INDEX url_set_len length(URL) TYPE set(2000) GRANULARITY 4
2 ALTER TABLE hits_UserID_URL MATERIALIZE INDEX url_set_len;
  
```

```

1 ✓ SELECT
2     UserID,
3     count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE length(URL) = 500
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
  
```

Diagram illustrating the ClickHouse storage structure for the `hits_UserID_URL` table:

The table consists of three columns: `UserID.bin`, `URL.bin`, and `EventTime.bin`. The data is organized into granules:

- granule 0** (with 8192 rows): Contains rows 0 to 8191.
- granule 1** (with 8192 rows): Contains rows 8192 to 16383.
- granule 1082** (with 3937 rows): Contains rows 8863 to 8867.

A green bracket on the right indicates that the condition `length(URL) > set(2000)` filters out rows from granule 1. An orange arrow points from the text above to this bracket.



Note:

1. set 适用于局部区间内重复率较多的固定数据，比如城市、报错码等

```

1 ✓ 9aae564a16) (SelectExecutor): Key condition: unknown
2 9aae564a16) (SelectExecutor): Filtering marks by primary and secondary keys
3 9aae564a16) (SelectExecutor): Index `url_set_len` has dropped 459/1083 granules.
4 9aae564a16) (SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1083/1083 marks by primary key, 624 marks to read from 64 ranges
5 9aae564a16) (SelectExecutor): Spreading mark ranges among streams (default reading)
6 9aae564a16) (SelectExecutor): Reading approx. 5111808 rows with 8 streams
  
```

2. Skipping Indexes – bloom filter



```
1 ✓ ALTER TABLE hits_UserID_URL ADD INDEX url_bf URL TYPE bloom_filter GRANULARITY 1  
2 ALTER TABLE hits_UserID_URL MATERIALIZE INDEX url_bf;
```

```
1 v SELECT
2     UserID,
3     count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE URL = 'http://public/?hash=8xqk3hMN9aDmK6l+K2HMuVhTDUnF0T5ugr7ITsXOPDI%3DfdSMTQ0ZDEwMDkuaHRtbD9
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
```

The diagram illustrates the storage and retrieval process of a dataset using granules and bloom filters.

Dataset Structure:

- Granule 0:** 8192 rows (highlighted in orange)
- Granule 1:** 8192 rows (highlighted in orange)
- Granule 1082:** 3937 rows (highlighted in orange)

Hashing and Bloom Filter:

- URL values are hashed (indicated by green curly braces).
- The resulting hash values are used to mark bits in a **bloom filter**.
- The bloom filter consists of two parts: **granules** and **bloom filter**.
- mark 0:** Corresponds to the first bit in the bloom filter.
- mark 1:** Corresponds to the second bit in the bloom filter.
- mark 1082:** Corresponds to the 1083rd bit in the bloom filter.

Retrieval Process:

When a URL is hashed, the resulting bit patterns are checked against the bloom filter to determine which granules contain the URL. The diagram shows arrows pointing from the hash values in the tables to the corresponding bits in the bloom filter.

Note:

1. Bloom filter 适用于局部区间内重复率较多的数据，比如手机号、query、日志等

```
1 ✓99aae564a16) (SelectExecutor): Filtering marks by primary and secondary keys
2 99aae564a16) (SelectExecutor): Used generic exclusion search over index for part all_1_9_2_48 with 1537 steps
3 99aae564a16) (SelectExecutor): Index `url_bf` has dropped 1050/1076 granules.
4 99aae564a16) (SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1076/1083 marks by primary key, 26 marks to read from 26 ranges
5 99aae564a16) (SelectExecutor): Spreading mark ranges among streams (default reading)
6 99aae564a16) (SelectExecutor): Reading approx. 212992 rows with 4 streams
```

2. Skipping Indexes – tokenbf_v1, ngrambf_v1



- 1 ✓ Token Bloom Filter
- 2 The same as ngrambf_v1, but stores tokens instead of ngrams. Tokens are sequences separated by non-alphanumeric characters. - 非字母数字的字符
- 3
- 4 Syntax: tokenbf_v1(size_of_bloom_filter_in_bytes, number_of_hash_functions, random_seed)

1 ✓ N-gram Bloom Filter
2 Only works with datatypes: String, FixedString and Map.
3
4 Syntax: ngrambf_v1(n, size_of_bloom_filter_in_bytes, number_of_hash_functions, random_seed)

```
1 ▼ SELECT tokens('https://www.zanbil.ir/m/filter/b113')
2
3 tokens
4 ┌───┐
5 │ ['https', 'www', 'zanbil', 'ir', 'm', 'filter', 'b113'] ──┐
6
7
8
9 SELECT ngrams('https://www.zanbil.ir/m/filter/b113', 3)
10
11 ngrams('https://www.zanbil.ir/m/filter/b113', 3)
12 ┌───┐
13 │ ['htt', 'ttp', 'tps', 'ps:', 's:/', '://', '//w', '/ww', 'www', 'ww.', 'w.z', '.za', 'zan', 'anb', 'nbi', 'bil', 'i'] ──┐
```

Note:

1. 当需要执行**文本搜索**时， tokenbf_v1, ngrambf_v1 类型的 skipping index 比较有用，操作符包括 LIKE, IN, hasToken
 2. 基于 token 的索引使用非字母数字字符作为分隔符生成 token，这在查询时只能匹配 token (或者说整个单词)
 3. 对于更多的粒度匹配，可以使用 N-gram Bloom Filter，把字符串分割成指定大小的子串，从而完成子串匹配
 4. 适用于文本搜索

2. Skipping Indexes – tokenbf_v1, ngrambf_v1



```
1 ✓ ALTER TABLE hits_UserID_URL ADD INDEX url_token URL TYPE tokenbf_v1(10480, 20, 20) GRANULARITY 1
2   OPTIMIZE TABLE hits_UserID_URL FINAL
```

```
1 ✓ SELECT
2   UserID,
3   count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE hasToken(URL, 'c291cmNlPWxjJnV0bV9jYW1wYW5lbWEudWE1MkY')
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
```

```
1 ✓ AS Count FROM hits_UserID_URL WHERE hasToken(URL, 'c291cmNlPWxjJnV0bV9jYW1wYW5lbWEudWE1MkY') GROUP BY UserID ORDER BY Count DESC LIMIT 10 (stage: Complete)
2
3
4 i for moving to the tail of PREWHERE is -1
5 'WxjJnV0bV9jYW1wYW5lbWEudWE1MkY'_String) moved to PREWHERE
6 i) (SelectExecutor): Key condition: unknown
7 i) (SelectExecutor): Filtering marks by primary and secondary keys
8 i) (SelectExecutor): Index `url_token` has dropped 807/1083 granules. ↗
9 i) (SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1083/1083 marks by primary key, 276 marks to read from 211 ranges
10 i) (SelectExecutor): Spreading mark ranges among streams (default reading)
11 i) (SelectExecutor): Reading approx. 2260992 rows with 8 streams ↗
```

2. Skipping Indexes – 小结



```
1 SELECT
2     UserID,
3     count(UserID) AS Count
4 FROM hits_UserID_URL
5 WHERE URL = 'http://public_search'
6 GROUP BY UserID
7 ORDER BY Count DESC
8 LIMIT 10
```

Question:

Q: 怎么加速这种查询?

A: Secondary Table、MV、Projection、Skipping Index。

```
1 \v INDEX index_name expr TYPE type(...) [GRANULARITY granularity_value]
2   - MinMax
3   - set(max_rows)
4   - bloom_filter([false_positive])
5   - tokenbf_v1(size_of_bloom_filter_in_bytes, number_of_hash_functions, random_seed)
6   - ngrambf_v1(n, size_of_bloom_filter_in_bytes, number_of_hash_functions, random_seed)
```



03 ClickHouse Explain

3. EXPLAIN Statement



```
1 SELECT
2     URL,
3     count(URL) AS Count
4 FROM hits_UserID_URL
5 WHERE UserID = 749927693
6 GROUP BY URL
7 ORDER BY Count DESC
8 LIMIT 10
9
10 Query id: 85706ec8-14cb-40b8-b851-85ef0de3f5eb
11
12   URL          Count
13 1. http://auto.ru/chatay-barana... 170
14 2. http://auto.ru/chatay-id=371... 52
15 3. http://public_search           45
16 4. http://kovrik-medvedevushku... 36
17 5. http://forumal               33
18 6. http://korablitz.ru/L_10FFER... 14
19 7. http://auto.ru/chatay-id=371... 14
20 8. http://auto.ru/chatay-john-D... 13
21 9. http://auto.ru/chatay-john-D... 10
22 10. http://wot/html?page/23600_m... 9
23
24
25 10 rows in set. Elapsed: 0.007 sec. Processed 8.19 thousand rows, 740.18 KB (1.10 million rows/s., 99.05 MB/s.)
26 Peak memory usage: 243.55 KiB.
```

3. EXPLAIN Statement

Question:

1. 怎么了解 SQL 命中了哪些索引呢？索引过滤了多少数据呢？

2. 怎么看 SQL 的并行查询情况呢？索引慢是不是因为计算资源不够？

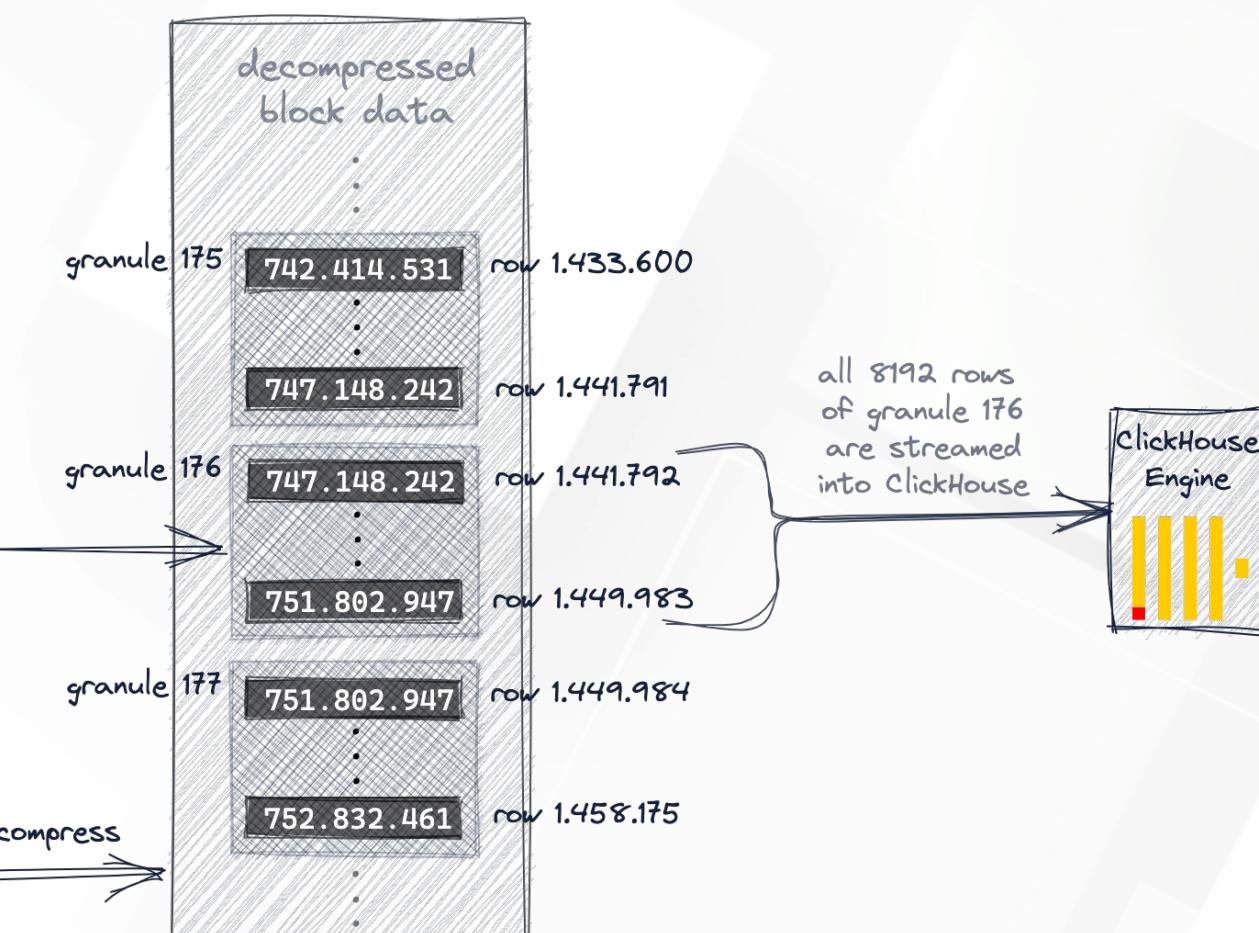
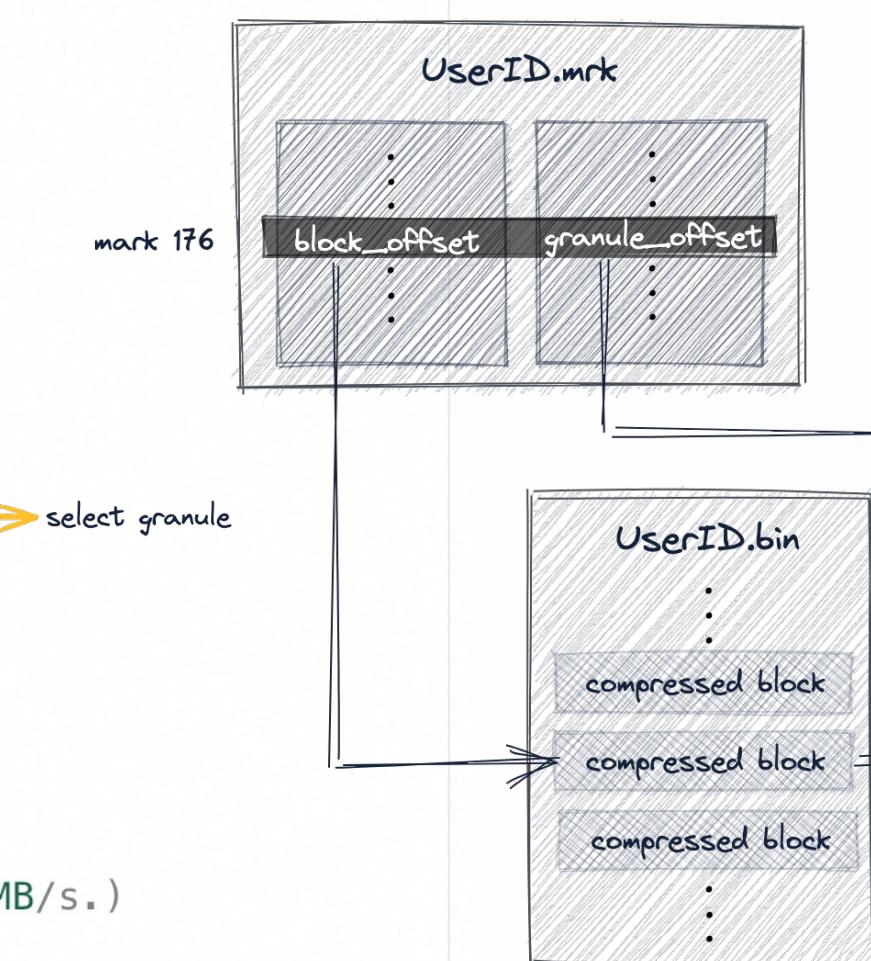
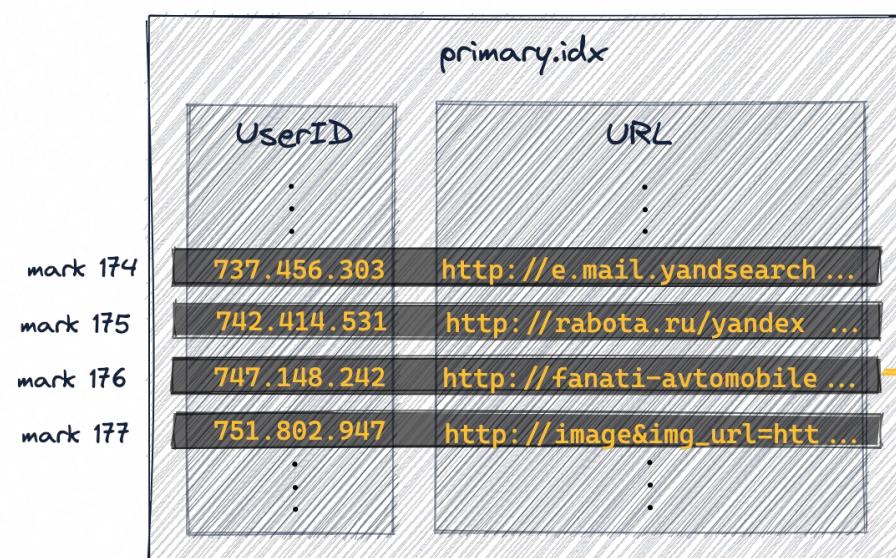
3. 怎么优化这些 SQL 呢？是增加跳数索引、MV、或者重新给主键重新排序？

→ EXPLAIN Statement

```

1 SELECT
2   URL,
3   count(URL) AS Count
4   FROM hits_UserID_URL
5 WHERE UserID = 749927693
6 GROUP BY URL
7 ORDER BY Count DESC
8 LIMIT 10
9
10 Query id: 85706ec8-14cb-40b8-b851-85ef0de3f5eb
11
12   URL      Count
13  1. http://auto.ru/chatay-barana... 170
14  2. http://auto.ru/chatay-id=371... 52
15  3. http://public_search          45
16  4. http://kovrik-medvedevushku... 36
17  5. http://forumal...            33
18  6. http://korablitz.ru/L_10FFER... 14
19  7. http://auto.ru/chatay-id=371... 14
20  8. http://auto.ru/chatay-john-D... 13
21  9. http://auto.ru/chatay-john-D... 10
22 10. http://wot/html?page/23600_m... 9
23
24
25 10 rows in set. Elapsed: 0.007 sec. Processed 8.19 thousand rows, 740.18 KB (1.10 million rows/s., 99.05 MB/s.)
26 Peak memory usage: 243.55 KiB.

```



3. EXPLAIN Statement



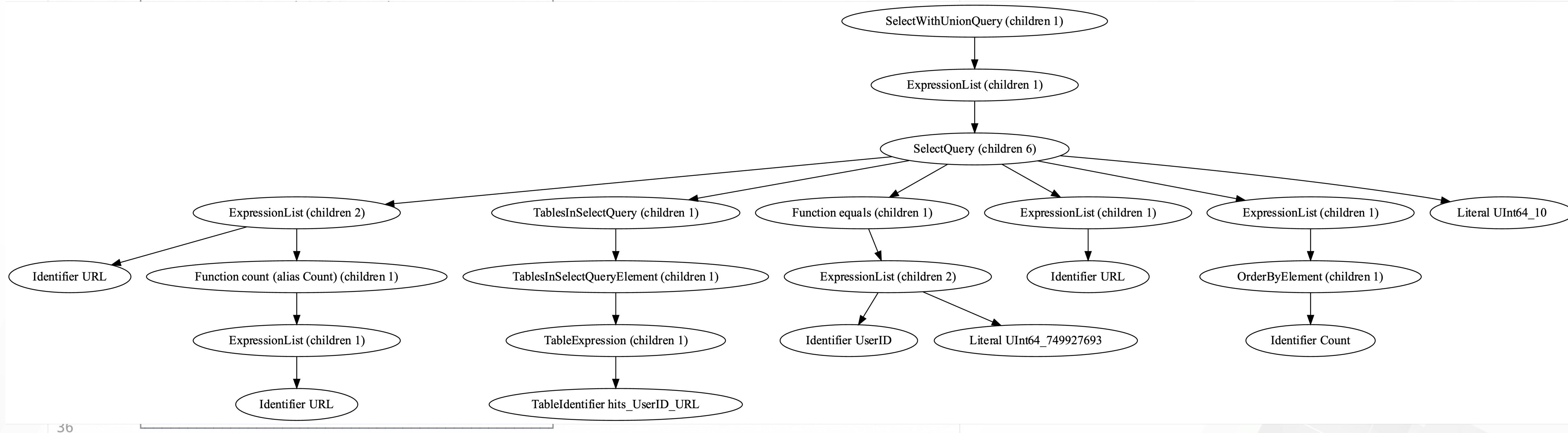
```
1 ✓ EXPLAIN [AST | SYNTAX | QUERY TREE | PLAN | PIPELINE | ESTIMATE | TABLE OVERRIDE] [setting = value, ...]
2 [
3     SELECT ... |
4     tableFunction(...) [COLUMNS (...)] [ORDER BY ...] [PARTITION BY ...] [PRIMARY KEY] [SAMPLE BY ...] [TTL ...]
5 ]
6 [FORMAT ...]
```

```
1 ✓ EXPLAIN Types
2 AST – Abstract syntax tree.
3 SYNTAX – Query text after AST-level optimizations.
4 PLAN – Query execution plan.
5 PIPELINE – Query execution pipeline.
```

3.1 EXPLAIN AST



```
1 v EXPLAIN AST
2 SELECT
3     URL,
4     count(URL) AS Count
5 FROM hits_UserID_URL
6 WHERE UserID = 749927693
7 GROUP BY URL
8 ORDER BY Count DESC
9 LIMIT 10
10
11 Query id: a238bed9-b0fa-4db2-8d07-c81a91e0f84a
12
13 explain
14 1. [SelectWithUnionQuery (children 1)]
```



3.2 EXPLAIN SYNTAX



```
1 v EXPLAIN SYNTAX
2 SELECT
3     URL,
4     count(URL) AS Count
5 FROM hits_UserID_URL
6 WHERE UserID = 749927693
7 GROUP BY URL
8 ORDER BY Count DESC
9 LIMIT 10
10
11 Query id: f732feb6-0bc1-4d91-b90d-10079d1308af
12
13 explain
14 1. SELECT
15     URL,
16     count(URL) AS Count
17     FROM hits_UserID_URL
18     WHERE UserID = 749927693
19     GROUP BY URL
20     ORDER BY Count DESC
21     LIMIT 10
22
23
24 8 rows in set. Elapsed: 0.001 sec.
```

```
1 v EXPLAIN SYNTAX
2 SELECT *
3 FROM hits_UserID_URL AS a, hits_UserID_URL AS b
4 LIMIT 10
5
6 Query id: b829740a-f89f-4d22-b9c9-4a8c0877b6fd
7
8 explain
9 1. SELECT
10    UserID,
11    URL,
12    EventTime,
13    b.UserID,
14    b.URL,
15    b.EventTime
16    FROM hits_UserID_URL AS a
17    CROSS JOIN hits_UserID_URL AS b
18    LIMIT 10
19
20
21 10 rows in set. Elapsed: 0.002 sec.
```

3.3 EXPLAIN PLAN

Plan 就是执行计划，默认的 explain 语句就是执行的 explain plan，查看 SQL 执行计划。

```

1 v EXPLAIN
2   SELECT
3     URL,
4       count(URL) AS Count
5   FROM hits_UserID_URL
6 WHERE UserID = 749927693
7 GROUP BY URL
8 ORDER BY Count DESC
9 LIMIT 10
10
11 Query id: 34877b39-7118-41f4-a246-16eda100404b
12
13   explain
14   1. Expression (Project names)
15   2. Limit (preliminary LIMIT (without OFFSET))
16   3.   Sorting (Sorting for ORDER BY)
17   4.     Expression ((Before ORDER BY + Projection))
18   5.     Aggregating
19   6.       Expression (Before GROUP BY)
20   7.         Expression
21   8.           ReadFromMergeTree (default.hits_UserID_URL)
22
23
24 8 rows in set. Elapsed: 0.002 sec.

```

```

1 v EXPLAIN actions = 1
2   SELECT
3     URL,
4       count(URL) AS Count
5   FROM hits_UserID_URL
6 WHERE UserID = 749927693
7 GROUP BY URL
8 ORDER BY Count DESC
9 LIMIT 10
10 Query id: f029496d-5ad8-4c59-8130-01f02e0c2839
11   explain
12   1. Expression (Project names)
13   2. Actions: INPUT : 0 -> count(__table1.URL) UInt64 : 0
14   3.   INPUT : 1 -> __table1.URL String : 1
15   4.     ALIAS count(__table1.URL) :: 0 -> Count UInt64 : 2
16   5.     ALIAS __table1.URL :: 1 -> URL String : 0
17   6. Positions: 0 2
18   7.     Limit (preliminary LIMIT (without OFFSET))
19   8.       Limit 10
20   9.       Offset 0
21   10.      Sorting (Sorting for ORDER BY)
22   11.        Sort description: count(__table1.URL) DESC
23   12.        Limit 10
24   13.          Expression ((Before ORDER BY + Projection))
25   14.            Actions: INPUT :: 0 -> count(__table1.URL) UInt64 : 0
26   15.              INPUT :: 1 -> __table1.URL String : 1
27   16.              Positions: 0 1
28   17.                Aggregating
29   18.                  Keys: __table1.URL
30   19.                  Aggregates:
31   20.                    count(__table1.URL)
32   21.                      Function: count(String) -> UInt64
33   22.                      Arguments: __table1.URL
34   23. Skip merging: 0
35   24. Expression (Before GROUP BY)
36   25. Actions: INPUT :: 0 -> __table1.URL String : 0
37   26. Positions: 0
38   27.   Expression
39   28.     Actions: INPUT : 0 -> URL String : 0
40   29.       INPUT :: 1 -> UserID UInt32 : 1
41   30.         ALIAS URL :: 0 -> __table1.URL String : 2
42   31. Positions: 2
43   32.   ReadFromMergeTree (default.hits_UserID_URL)
44   33.   ReadType: Default
45   34.   Parts: 1
46   35.   Granules: 1
47   36.   Prewhere info
48   37.   Need filter: 1
49   38.     Prewhere filter
50   39.       Prewhere filter column: equals(__table1.UserID, 749927693 UInt32) (removed)
51   40.     Actions: INPUT : 0 -> UserID UInt32 : 0
52   41.       COLUMN Const(UInt32) -> 749927693 UInt32 UInt32 : 1
53   42.         FUNCTION equals(UserID : 0, 749927693 UInt32 :: 1) -> equals(__table1.UserID, 749927693 UInt32) !
54   43. Positions: 0 2
55

```

3.3 EXPLAIN PLAN

Plan 就是执行计划，默认的 explain 语句就是执行的 explain plan，查看 SQL 执行计划。

```

1 v EXPLAIN
2   SELECT
3     URL,
4       count(URL) AS Count
5   FROM hits_UserID_URL
6   WHERE UserID = 749927693
7   GROUP BY URL
8   ORDER BY Count DESC
9   LIMIT 10
10
11 Query id: 34877b39-7118-41f4-a246-16eda100404b
12
13   explain
14   1. Expression (Project names)
15   2. Limit (preliminary LIMIT (without OFFSET))
16   3.   Sorting (Sorting for ORDER BY)
17   4.     Expression ((Before ORDER BY + Projection))
18   5.     Aggregating
19   6.       Expression (Before GROUP BY)
20   7.     Expression
21   8.       ReadFromMergeTree (default.hits_UserID_URL)
22
23
24 8 rows in set. Elapsed: 0.002 sec.

```

```

1 v EXPLAIN indexes = 1
2   SELECT
3     URL,
4       count(URL) AS Count
5   FROM hits_UserID_URL
6   WHERE UserID = 749927693
7   GROUP BY URL
8   ORDER BY Count DESC
9   LIMIT 10
10
11 Query id: 43bbf6d0-48a5-4f64-9f4e-39c3aa3809d6
12
13   explain
14   1. Expression (Project names)
15   2. Limit (preliminary LIMIT (without OFFSET))
16   3.   Sorting (Sorting for ORDER BY)
17   4.     Expression ((Before ORDER BY + Projection))
18   5.     Aggregating
19   6.       Expression (Before GROUP BY)
20   7.     Expression
21   8.       ReadFromMergeTree (default.hits_UserID_URL)
22   9.       Indexes:
23   10.         PrimaryKey
24   11.           Keys:
25   12.             UserID
26   13.               Condition: (UserID in [749927693, 749927693])
27   14.               Parts: 1/1
28   15.               Granules: 1/1083
29

```

```

6 8-a48d-c99aae564a16) (SelectExecutor): Key condition: (column 0 in [749927693, 749927693])
7 8-a48d-c99aae564a16) (SelectExecutor): Filtering marks by primary and secondary keys
8 8-a48d-c99aae564a16) (SelectExecutor): Running binary search on index range for part all_1_9_2 (1083 marks)
9 8-a48d-c99aae564a16) (SelectExecutor): Found (LEFT) boundary mark: 176
10 8-a48d-c99aae564a16) (SelectExecutor): Found (RIGHT) boundary mark: 177
11 8-a48d-c99aae564a16) (SelectExecutor): Found continuous range in 19 steps
12 8-a48d-c99aae564a16) (SelectExecutor): Selected 1/1 parts by partition key, 1 parts by primary key, 1/1083 marks by primary key, 1 marks to read from 1 ranges
13 8-a48d-c99aae564a16) (SelectExecutor): Spreading mark ranges among streams (default reading)
14 8-a48d-c99aae564a16) (SelectExecutor): Reading 1 ranges in order from part all_1_9_2, approx. 8192 rows starting from 1441792

```

3.4 EXPLAIN PIPELINE



```
1 ✓ EXPLAIN PIPELINE ←
2 SELECT
3     URL,
4     count(URL) AS Count
5 FROM hits_UserID_URL
6 WHERE UserID = 749927693
7 GROUP BY URL
8 ORDER BY Count DESC
9 LIMIT 10
10
11 Query id: 1aa3dc5a-aae0-434e-a9f0-075f3a4705a2
12
13 ┌─explain─┐
14 1. (Expression)
15 2. ExpressionTransform
16 3. (Limit)
17 4. Limit
18 5. (Sorting)
19 6. MergingSortedTransform 8 → 1
20 7. MergeSortingTransform × 8
21 8. LimitsCheckingTransform × 8
22 9. PartialSortingTransform × 8
23 10. (Expression)
24 11. ExpressionTransform × 8
25 12. (Aggregating)
26 13. Resize 1 → 8
27 14. AggregatingTransform
28 15. (Expression)
29 16. ExpressionTransform
30 17. (Expression)
31 18. ExpressionTransform
32 19. (ReadFromMergeTree)
33 20. MergeTreeSelect(pool: ReadPoolInOrder, algorithm: InOrder) 0 → 1
34
35
36 20 rows in set. Elapsed: 0.003 sec.
```

3.4 EXPLAIN PIPELINE

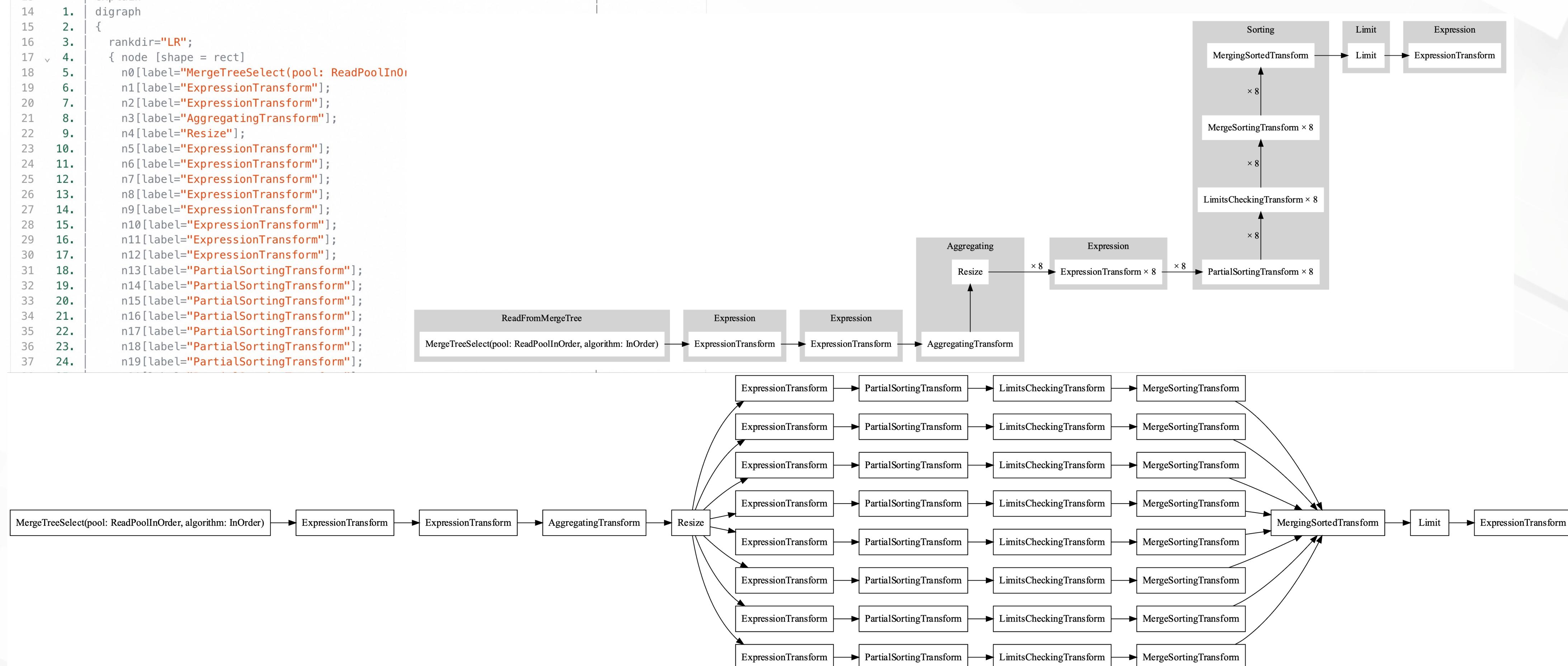
```

1 ✓ EXPLAIN PIPELINE graph = 1, compact = 0
2 SELECT
3   URL,
4     count(URL) AS Count
5 FROM hits_UserID_URL
6 WHERE UserID = 749927693
7 GROUP BY URL
8 ORDER BY Count DESC
9 LIMIT 10
10
11 Query id: 43ed1457-9cb9-450a-845e-25b70163ecfe
12
13   explain
14   1. digraph
15   2.
16   3. rankdir="LR";
17   4. { node [shape = rect]
18   5.   n0[label="MergeTreeSelect(pool: ReadPoolIn0];
19   6.   n1[label="ExpressionTransform"];
20   7.   n2[label="ExpressionTransform"];
21   8.   n3[label="AggregatingTransform"];
22   9.   n4[label="Resize"];
23   10.  n5[label="ExpressionTransform"];
24   11.  n6[label="ExpressionTransform"];
25   12.  n7[label="ExpressionTransform"];
26   13.  n8[label="ExpressionTransform"];
27   14.  n9[label="ExpressionTransform"];
28   15.  n10[label="ExpressionTransform"];
29   16.  n11[label="ExpressionTransform"];
30   17.  n12[label="ExpressionTransform"];
31   18.  n13[label="PartialSortingTransform"];
32   19.  n14[label="PartialSortingTransform"];
33   20.  n15[label="PartialSortingTransform"];
34   21.  n16[label="PartialSortingTransform"];
35   22.  n17[label="PartialSortingTransform"];
36   23.  n18[label="PartialSortingTransform"];
37   24.  n19[label="PartialSortingTransform"];

```

1.Explain pipeline graph=1

2.Explain pipeline graph=1, compact=0

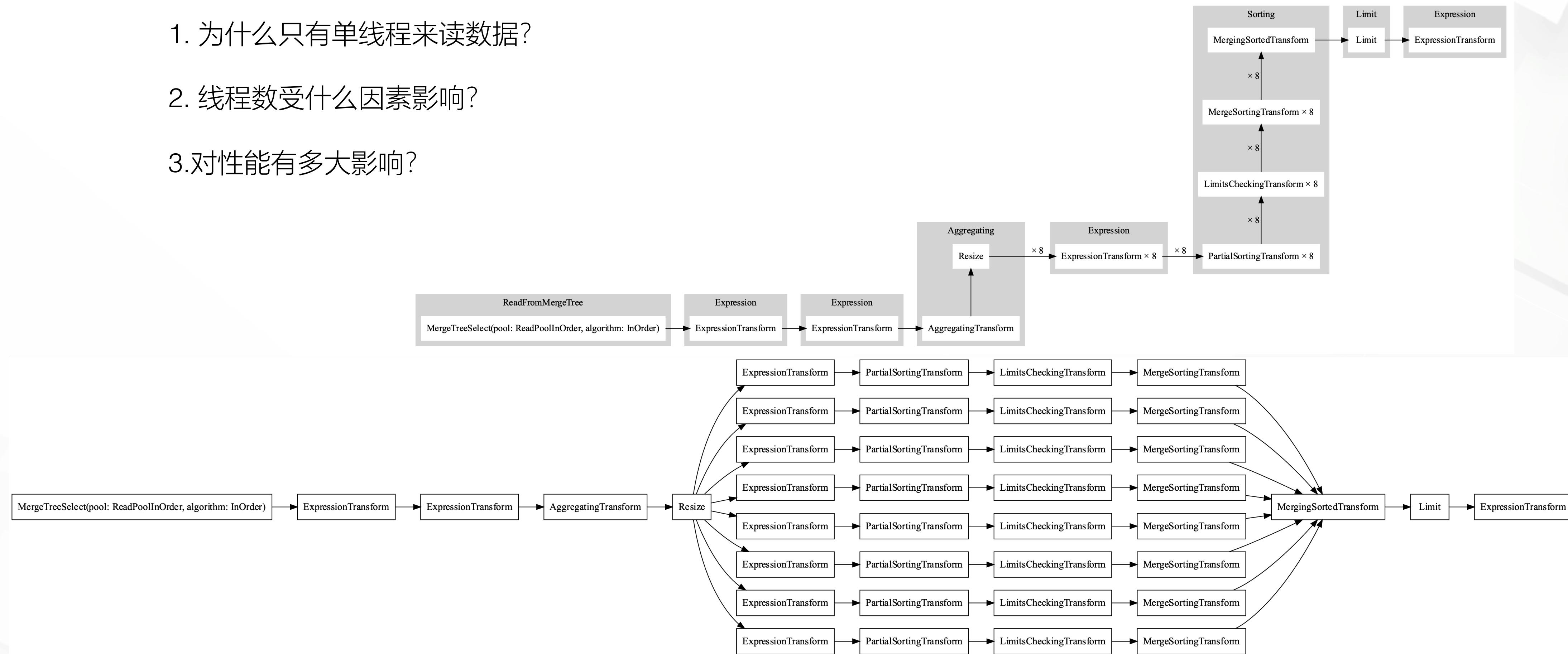


3.4 EXPLAIN PIPELINE



Question:

1. 为什么只有单线程来读数据?
2. 线程数受什么因素影响?
3. 对性能有多大影响?

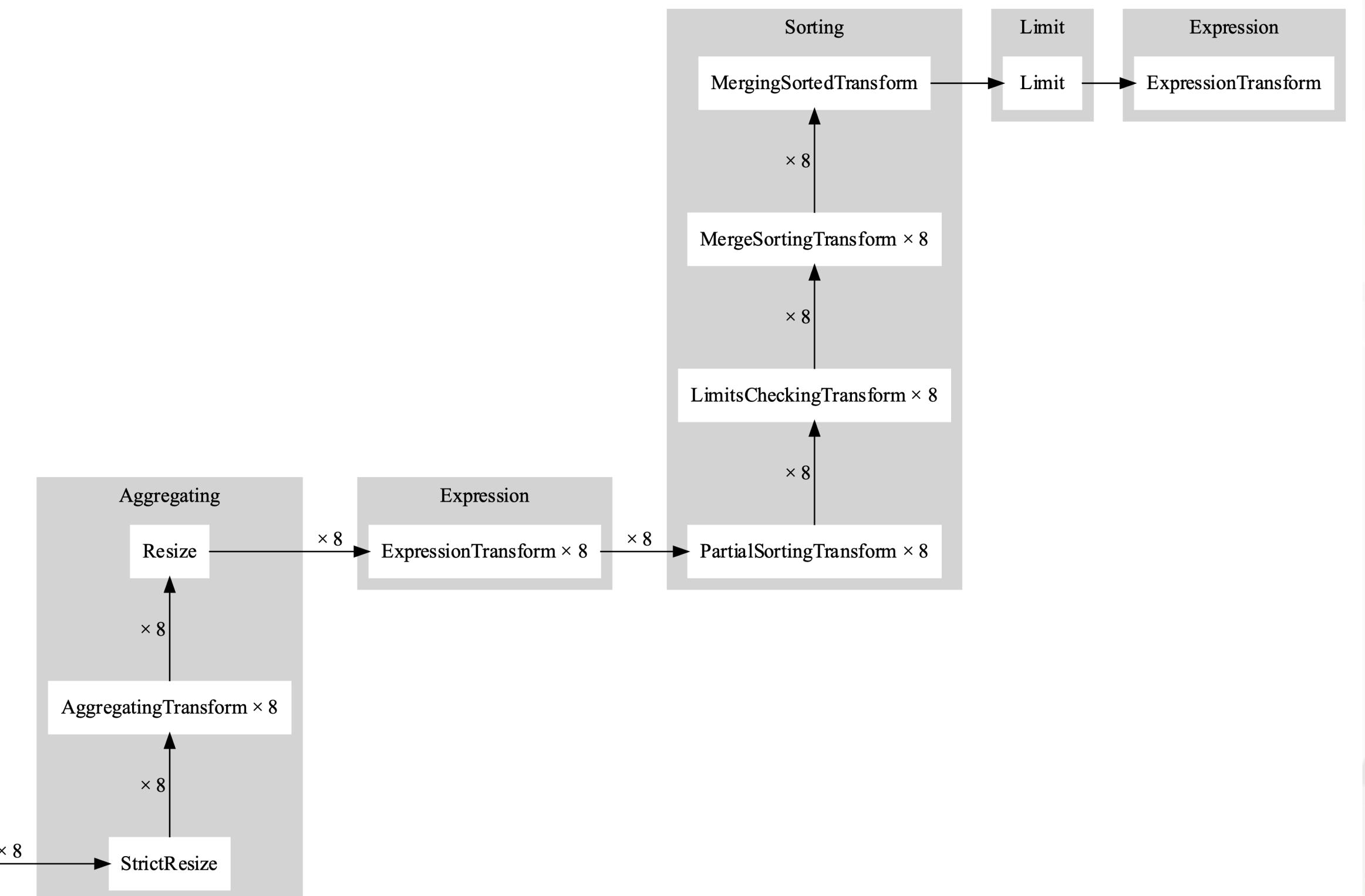


3.4 EXPLAIN PIPELINE



```
1 v EXPLAIN indexes = 1
2 SELECT
3     URL,
4     count(URL) AS Count
5 FROM hits_UserID_URL
6 WHERE UserID > 749927693    ←
7 GROUP BY URL
8 ORDER BY Count DESC
9 LIMIT 10
10
11 Query id: 2f9bd5cd-ada3-4132-92e3-8199324a5b2b
```

```
12
13     explain
14     1. Expression (Project names)
15     2. Limit (preliminary LIMIT (without OFFSET))
16     3. Sorting (Sorting for ORDER BY)
17     4. Expression ((Before ORDER BY + Projection))
18     5. Aggregating
19     6. Expression (Before GROUP BY)
20     7. Expression
21     8. ReadFromMergeTree (default.hits_UserID_URL)
22     9. Indexes:
23     10. PrimaryKey
24     11. Keys:
25     12. UserID
26     13. Condition: (UserID in [749927694, +Inf))
27     14. Parts: 1/1
28     15. Granules: 907/1083
29
```



3.4 EXPLAIN PIPELINE

```

13   explain
14 1. Expression (Project names)
15 2. Limit (preliminary LIMIT (without OFFSET))
16 3. Sorting (Sorting for ORDER BY)
17 4. Expression ((Before ORDER BY + Projection))
18 5. Aggregating
19 6. Expression (Before GROUP BY)
20 7. Expression
21 8. ReadFromMergeTree (default.hits_UserID_URL)
22 9. Indexes:
23 10. PrimaryKey
24 11. Keys:
25 12. UserID
26 13. Condition: (UserID in [749927694, +Inf))
27 14. Parts: 1/1
28 15. Granules: 907/1083
  
```

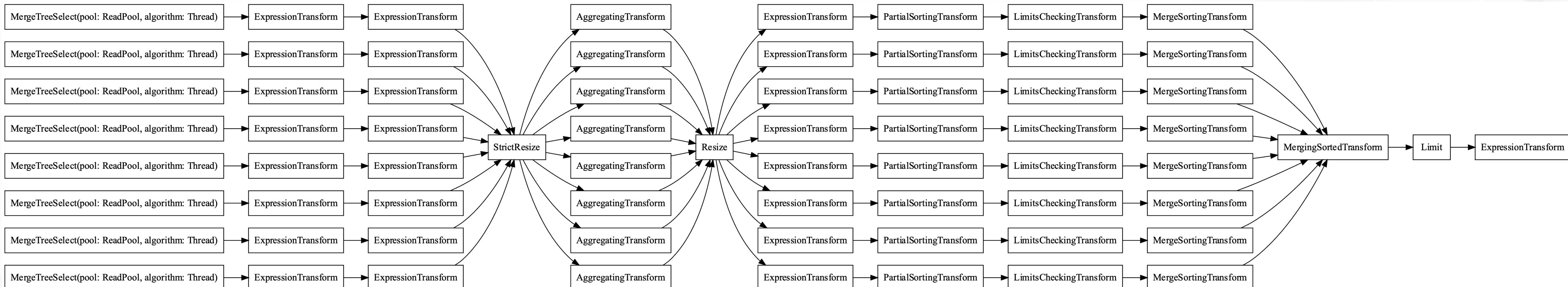
Question:

1. 为什么只有单线程来读数据?

2. 线程数受什么因素影响?

3. 对性能有多大影响?

1. merge_tree_min_rows_for_concurrent_read [default 163840 - 20 granules]
2. merge_tree_min_bytes_for_concurrent_read [default 251658240 - 240MB]



3.4 EXPLAIN PIPELINE



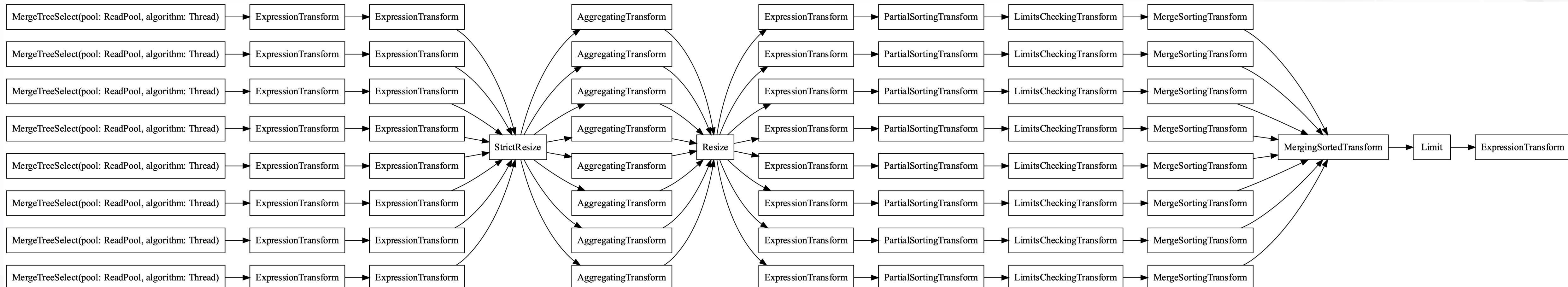
```
13 explain
14 1. Expression (Project names)
15 2. Limit (preliminary LIMIT (without OFFSET))
16 3. Sorting (Sorting for ORDER BY)
17 4. Expression ((Before ORDER BY + Projection))
18 5. Aggregating
19 6. Expression (Before GROUP BY)
20 7. Expression
21 8. ReadFromMergeTree (default.hits_UserID_URL)
22 9. Indexes:
23 10. PrimaryKey
24 11. Keys:
25 12. UserID
26 13. Condition: (UserID in [749927694, +Inf))
27 14. Parts: 1/1
28 15. Granules: 907/1083
29
```

Question:

1. 为什么只有单线程来读数据?
2. 线程数受什么因素影响?
3. 对性能有多大影响?

1. max_threads [default CPU cores]

```
1 SHOW SETTINGS LIKE 'max_threads'
2
3 Query id: b195e9ac-5a0f-49e3-a7da-2454f268305f
4
5 1. name type value
6 1. max_threads MaxThreads auto(8)
7
```



3.4 EXPLAIN PIPELINE

```

1 ✓ SELECT
2   URL,
3   count(URL) AS Count
4 FROM hits_UserID_URL
5 WHERE UserID > 749927693
6 GROUP BY URL
7 ORDER BY Count DESC
8 LIMIT 10 settings max_threads = N;
  
```

```

1 ✓ SETTINGS max_threads = 8
2 10 rows in set. Elapsed: 0.559 sec. Processed 7.43 million rows, 672.92 MB (13.28 million rows/s, 1.20 GB/s)
3 Peak memory usage: 623.05 MiB.

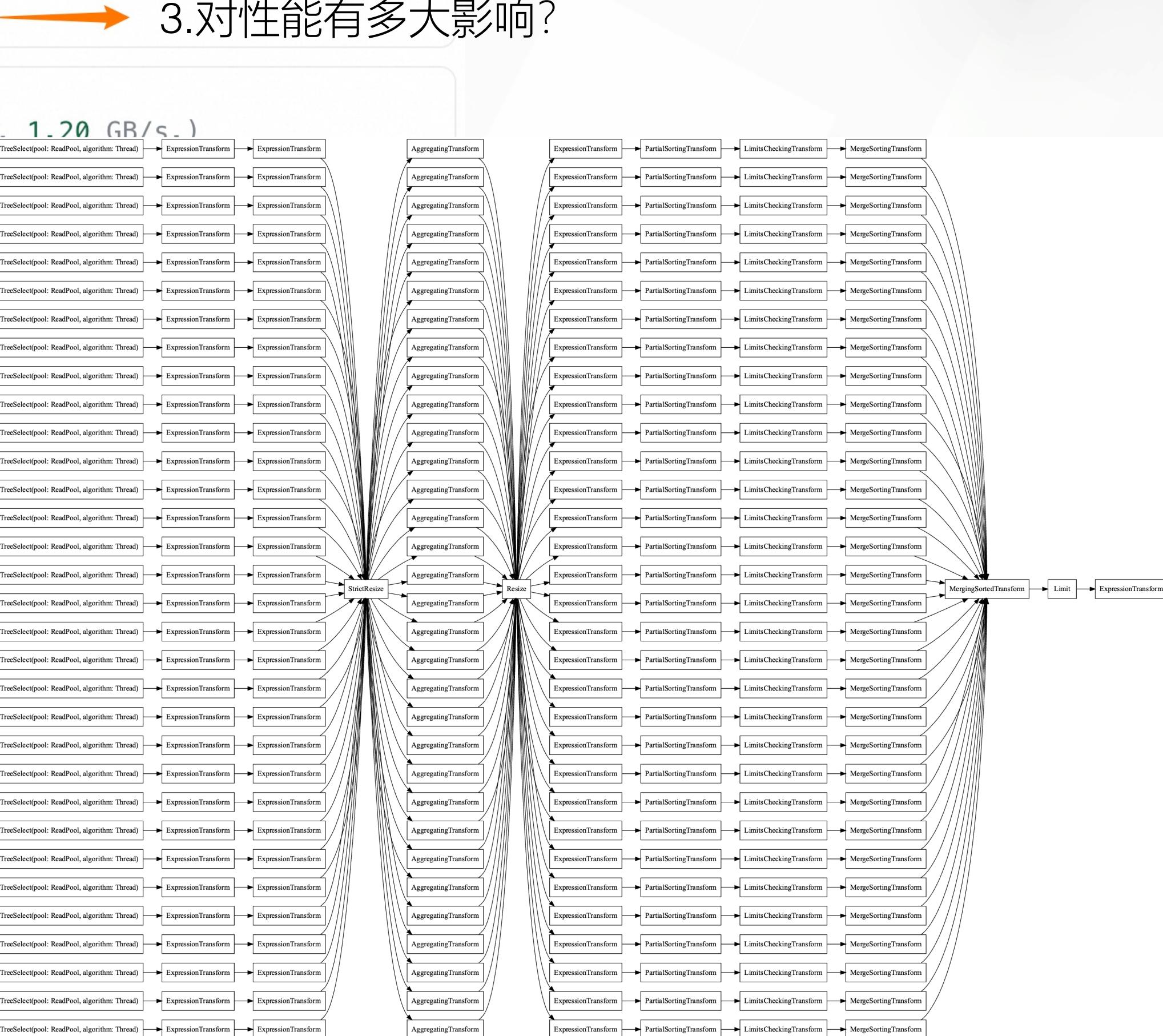
4
5 SETTINGS max_threads = 16
6 10 rows in set. Elapsed: 0.535 sec. Processed 7.43 million rows, 672.92 MB (13.89 million rows/s, 1.20 GB/s)
7 Peak memory usage: 666.82 MiB.

8
9 SETTINGS max_threads = 32
10 10 rows in set. Elapsed: 0.586 sec. Processed 7.43 million rows, 672.92 MB (12.67 million rows/s, 1.20 GB/s)
11 Peak memory usage: 774.41 MiB.

12
13 SETTINGS max_threads = 64
14 10 rows in set. Elapsed: 0.622 sec. Processed 7.43 million rows, 672.92 MB (11.94 million rows/s, 1.20 GB/s)
15 Peak memory usage: 854.14 MiB.
  
```

Question:

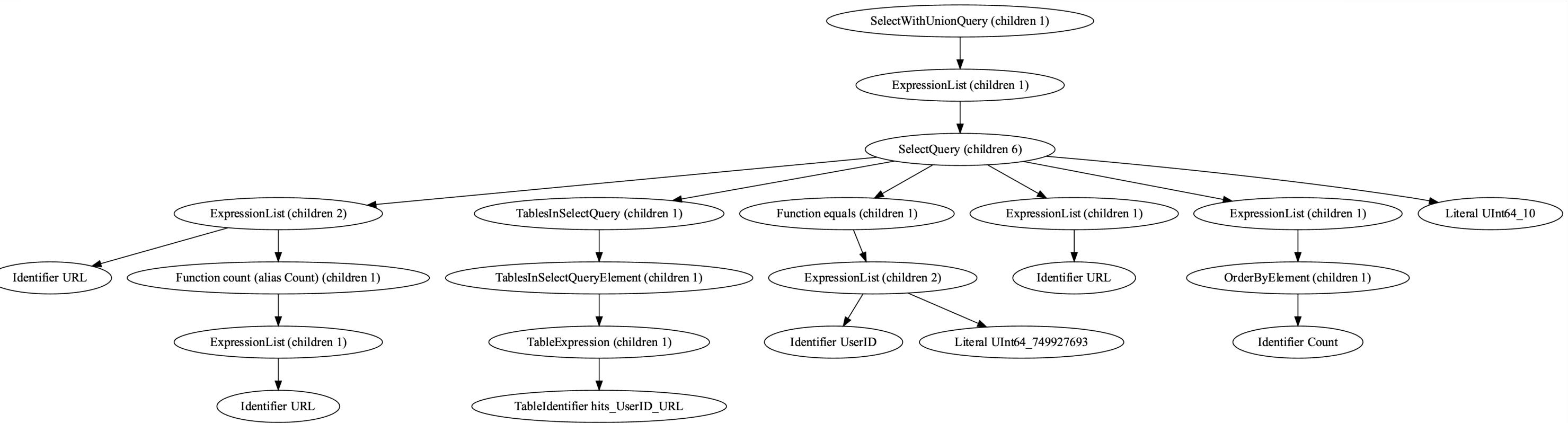
1. 为什么只有单线程来读数据?
2. 线程数受什么因素影响?
3. 对性能有多大影响?



1. max_threads 越大，并行执行线程越多；
2. 并行线程越多，同一时刻读到内存的数据越多，使用内存量越大；
3. Max_threads 超过 CPU cores 越多，线程上下文切换带来的影响越大，执行效率越低。

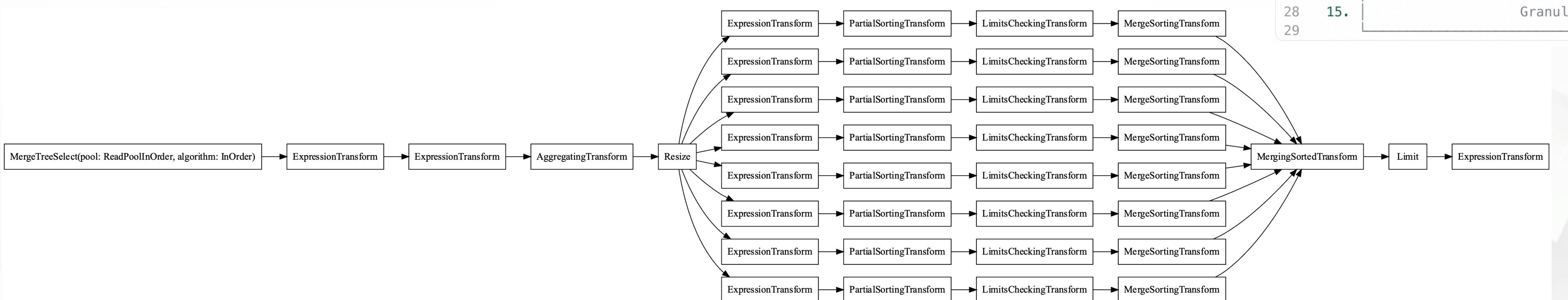
3. EXPLAIN Statement 小结

- 1 ✓ EXPLAIN Types
- 2 AST – Abstract syntax tree.
- 3 SYNTAX – Query text **after** AST-level optimizations.
- 4 PLAN – Query execution plan.
- 5 PIPELINE – Query execution pipeline.



```

1 ✓ EXPLAIN indexes = 1
2 SELECT
3   URL,
4     count(URL) AS Count
5   FROM hits_UserID_URL
6   WHERE UserID = 749927693
7   GROUP BY URL
8   ORDER BY Count DESC
9   LIMIT 10
10
11 Query id: 43bbf6d0-48a5-4f64-9f4e-39c3aa3809d6
12
13 explain
14   Expression (Project names)
15   Limit (preliminary LIMIT (without OFFSET))
16   Sorting (Sorting for ORDER BY)
17     Expression ((Before ORDER BY + Projection))
18   Aggregating
19     Expression (Before GROUP BY)
20   Expression
21     ReadFromMergeTree (default.hits_UserID_URL)
22   Indexes:
23     PrimaryKey
24   Keys:
25     UserID
26   Condition: (UserID in [749927693, 749927693])
27   Parts: 1/1
28   Granules: 1/1083
29
  
```





04 ClickHouse Optimization

4. Optimization



步骤：

1. 根据最频繁的业务场景，调整 order by key 的顺序；
2. 根据数据分布，增加 skipping indexes；
3. 根据特殊业务场景，增加第二张表、物化视图、Projection；
4. 根据 query plan，查看实际执行步骤、索引命中情况；
5. 根据 query pipeline，查看 plan 中每一步的并行执行情况。

```
1 ✓ EXPLAIN indexes = 1
2 SELECT
3     UserID,
4     count(UserID) AS Count
5 FROM hits_UserID_URL
6 WHERE URL = 'http://public_search'
7 GROUP BY UserID
8 ORDER BY Count DESC
9 LIMIT 10
10
11 Query id: 600f5e3d-a8e6-4a84-b36f-7015276248ec
12
13 explain
14 1. Expression (Project names)
15 2. Limit (preliminary LIMIT (without OFFSET))
16 3. Sorting (Sorting for ORDER BY)
17 4. Expression ((Before ORDER BY + Projection))
18 5. Aggregating
19 6. Expression (Before GROUP BY)
20 7. Expression
21 8. ReadFromMergeTree (default.hits_UserID_URL)
22 9. Indexes:
23 10. PrimaryKey
24 11. Keys:
25 12. URL
26 13. Condition: (URL in ['http://public_search', 'http://public_search'])
27 14. Parts: 1/1
28 15. Granules: 1076/1083
29 16. Skip
30 17. Name: url_skipping_index
31 18. Description: minmax GRANULARITY 4
32 19. Parts: 1/1
33 20. Granules: 1076/1076
34
35
36 20 rows in set. Elapsed: 0.004 sec.
```

0 Q&A 5



奥运会全球指定云服务商