# From ElasticSearch To ClickHouse, Log Analysis Practice in bilibili

Jiaqi Shu(東家麒)
bilibili INFRA

Tao Ling (凌涛)
bilibili OLAP

# Agenda

❖ Log Analysis System based on ClickHouse

  ➢ Challenge & problems

  ➢ Why ClickHouse?

  ➢ How we build

❖ Optimization & Customization of ClickHouse for Log Analysis

  ➢ Challenges for ClickHouse

  ➢ Overall Design

  ➢ Implicit Map

# Log Analysis System based on ClickHouse

# Challenge & problems

- PB level data, Based on Elasticsearch, 1k+. nodes

- Write performance bottleneck

- Heavy cost on index

- Operation and maintenance costs

# Why ClickHouse?

## Requirements

- Low cost

- High performance

- Support semi-structure

## ClickHouse

- Cost saving over 60%
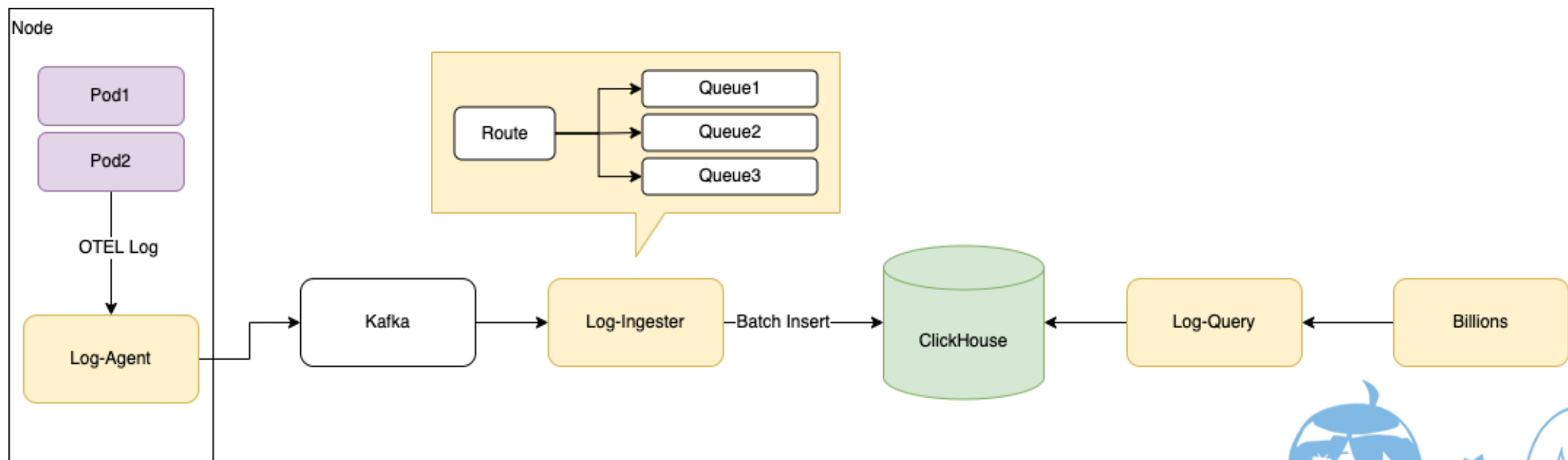
- Avg Query latency < 1s

# How we build

# How we build

```
Create Table <log_app_name>
{
    timestamp Datetime64,
    hostname  String CODC(ZSTD(1)),
    zone        String CODC(ZSTD(1)),
    log          String CODC(ZSTD(1)),
    ....
    string_map  MapV2(String, Nullable (String))
                            CODEC(ZSTD(1))
    number_map  MapV2(String, Nullable (Float64))
                            CODEC(ZSTD(1))
    bool_map  MapV2(String, Nullable (UInt8))
}
....
```
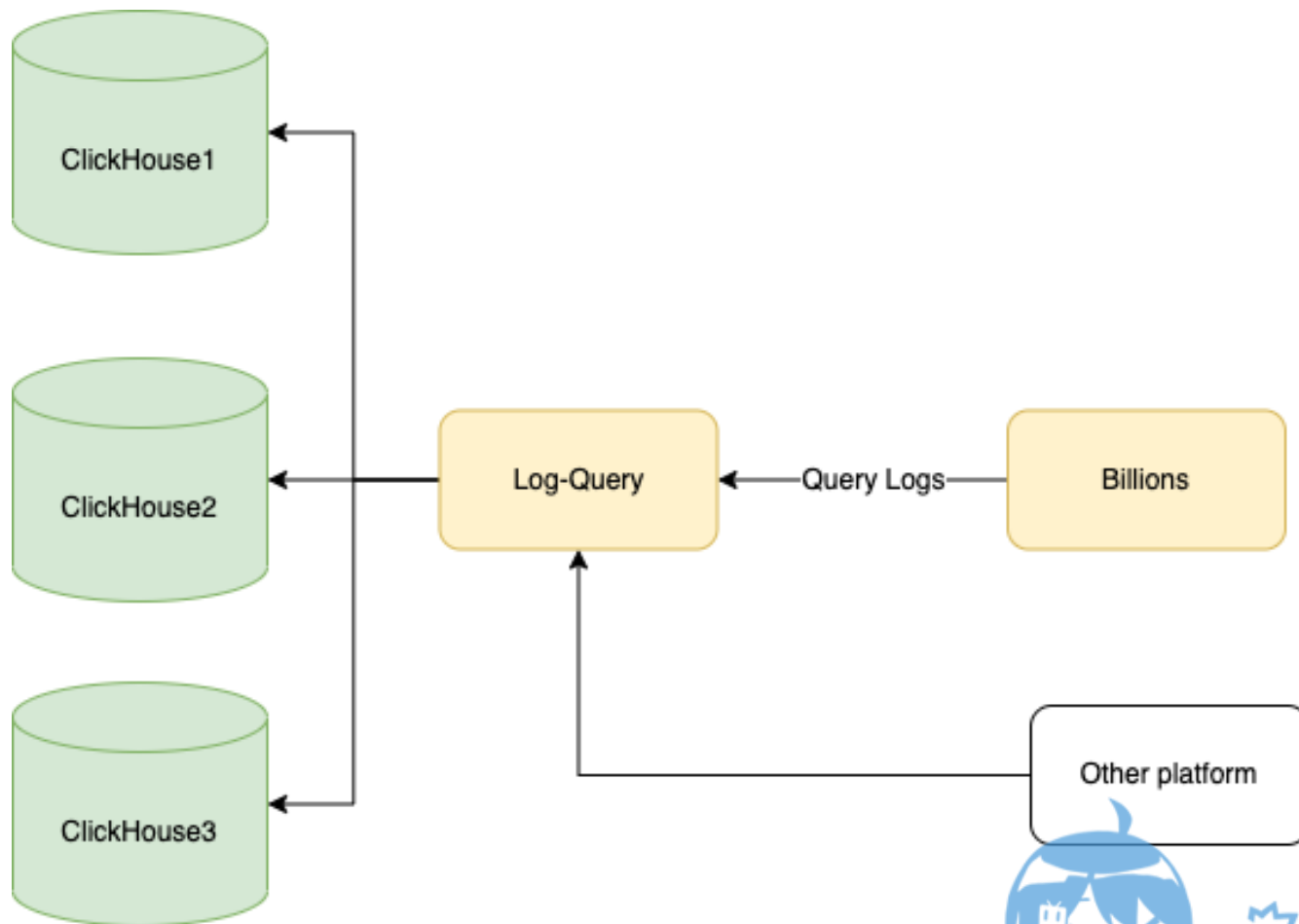
# How we build

## Log Query

- High level query

- Query routing

- Rate Limit, Access log...

# How we build

```json
{
    "app_id":"ops.billions.elastalert-worker",
    "query":"level = 'ERROR' AND cluster = 'shylf'",
    "from":"2022-03-02 06:00:00",
    "to":"2022-03-02 06:15:00"
}
```

```sql
SELECT *
FROM ops_billions_elastalert_worker
WHERE
  level = 'ERROR'
AND
  string_map['cluster'] = 'shylf'
AND
  timestamp >= '2022-03-02 06:00:00'
AND
  timestamp >= '2022-03-02 06:15:00'
ORDER BY
  timestamp desc
LIMIT 1000
```

```json
[
    {
        "timestamp":"2022-03-02 06:01:03",
        "level":"ERROR",
        "log":"fail to get connection",
        "cluster":"shylf"
    },
    {
        "timestamp":"2022-03-02 06:01:02",
        "level":"ERROR",
        "log":"io error",
        "cluster":"shylf"
    }
]
```
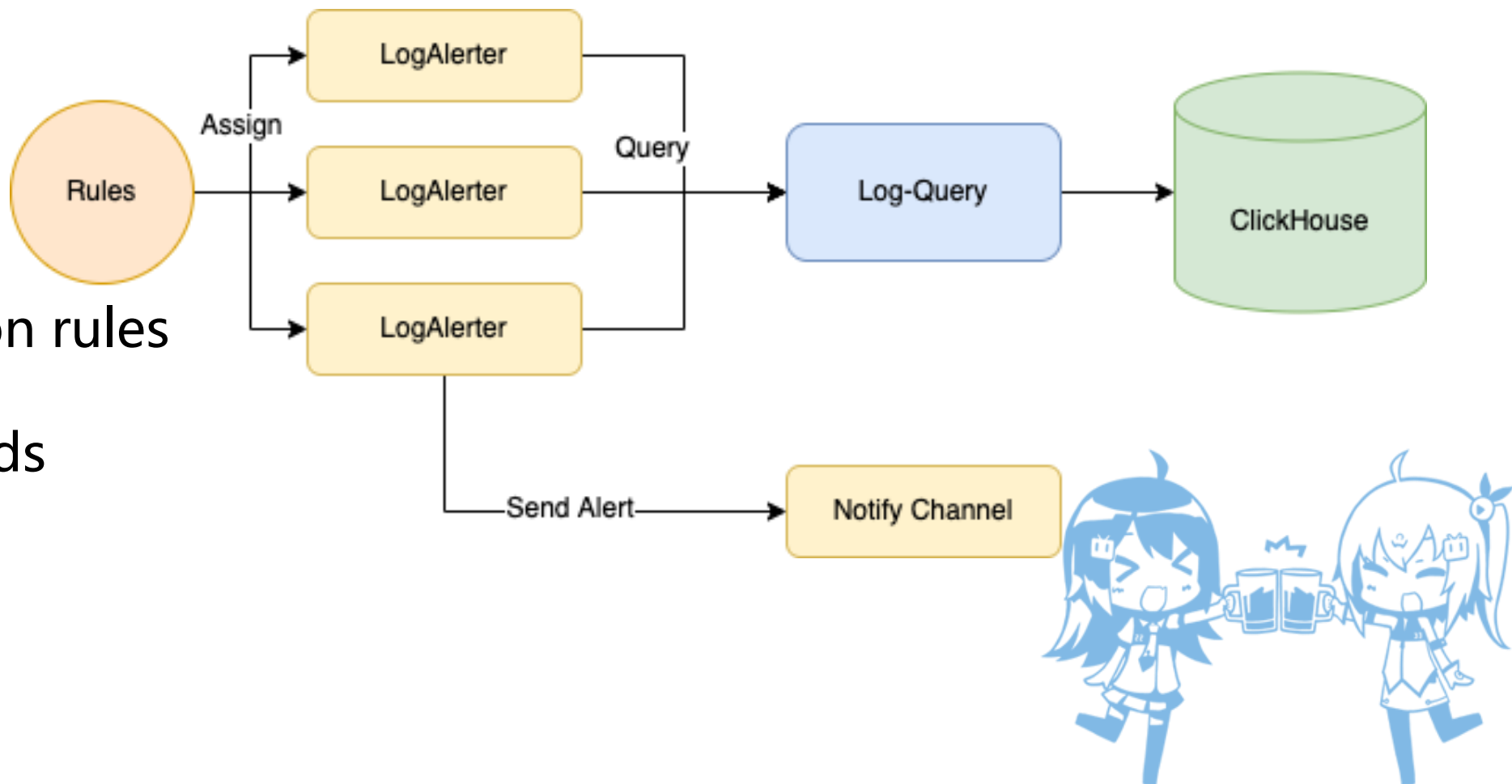
| timestamp | level | instance_name | log | string_map |
|---|---|---|---|---|
| 2022-03-02 06:01:03 | ERROR | shylf-opslog-02 | fail to get connection | {"cluster":"shylf"} |
| 2022-03-02 06:01:02 | ERROR | shylf-opslog-01 | io error | {"cluster":"shylf"} |

# How we build

# How we build

## Log Alerter

- Distributed

- Various calculation rules

- Alert in 30 seconds

# Optimization & Customization of ClickHouse for Log Analysis

# Challenges for ClickHouse

## Data writing:

- Large number of applications
- Low write latency
- Dynamic schema for private attribute
- Specify codec to save more disk

## Data querying:

- Query based on application use different conditions
- Clickhouse native map type is underperforming

# Overall Design

- ## Each application has a separate table

  We can specify different settings based on application at table level.
  1. Skipping indexes
  2. Max implicit columns size

- ## Storage policy for hot and cold data

  Specify storage policy in table DDL
  1. NVME for hot data
  2. SATA for cold data
  3. ZSTD codec better than default codec (ZSTD can save up to 60% more space in some scenarios)

- ## Use Map type for dynamic attribute

# Implicit Map

## Native Map Type

1. Tuple of Arrays

2. Queries read redundant data

3. Skipping indexes don't support Map type

Array of Keys    Array of values

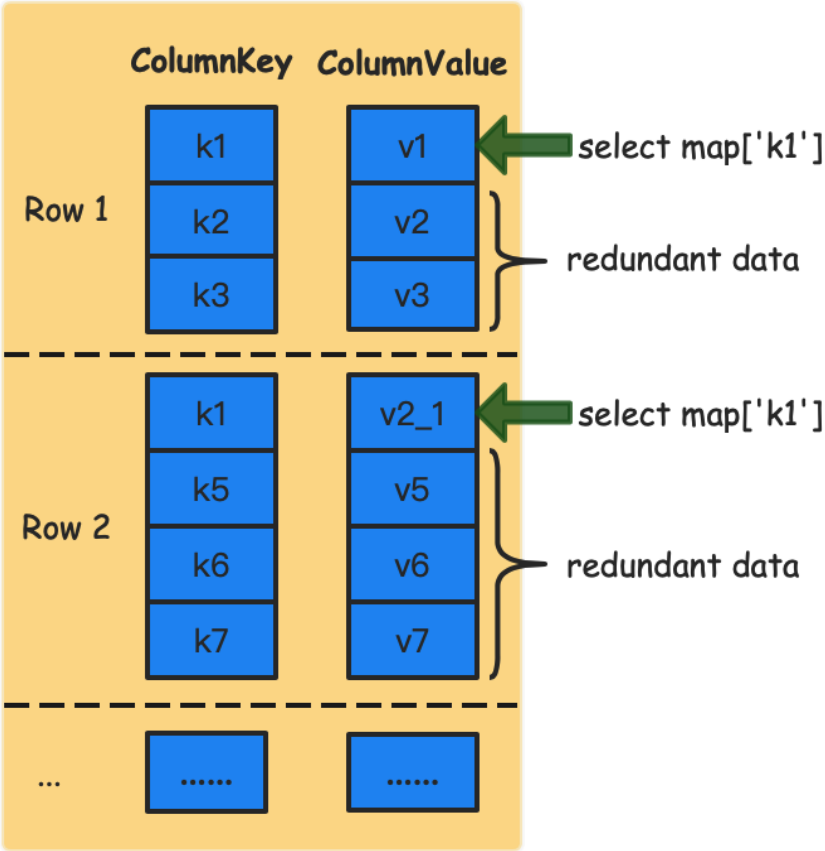| k1 | | v1 |
| k2 | | v2 |
| k3 | | v3 |
| k4 | | v4 |
| k5 | | v5 |
| k6 | | v6 |
| k7 | | v7 |
| ... | | ... |

# Implicit Map

## Redundant data read

*SELECT map['k1'] FROM table*

- 2 arrays will be read

- All elements of map will be read

- Every rows of 2 arrays will be read

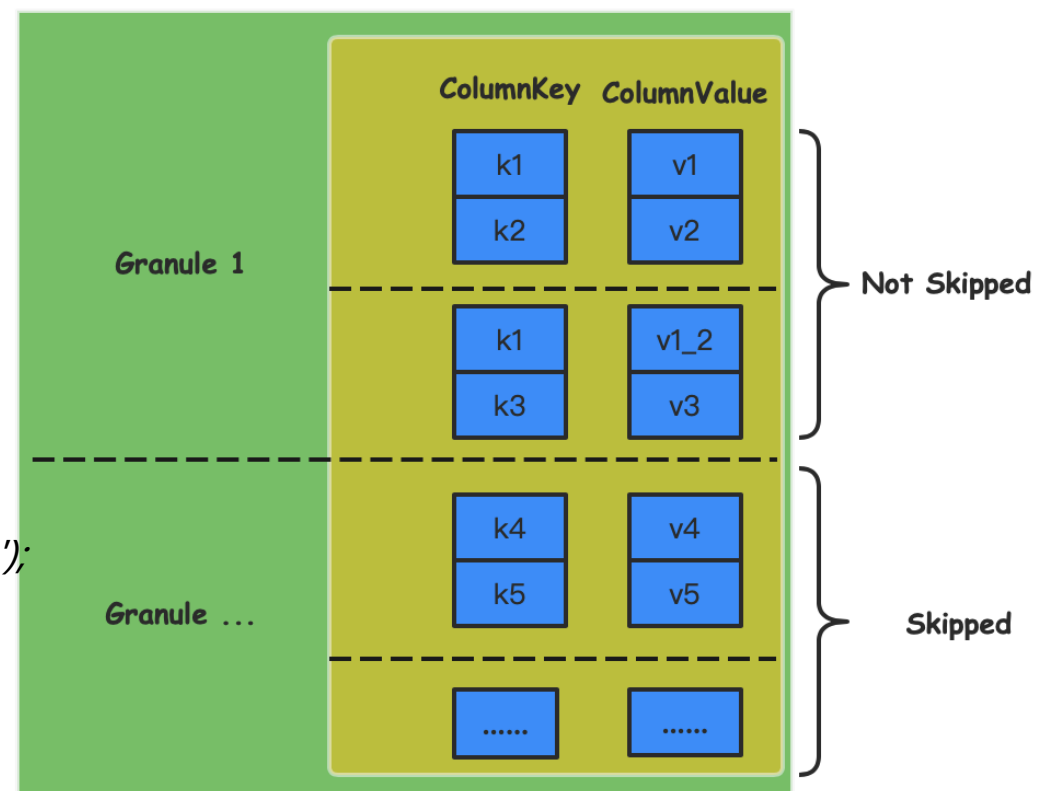# Skip Indices for Map

## Skip Indices for Map

```
CREATE TABLE bloom_filter_map
(
    `id` UInt32,
    `map` Map(String, String),
    INDEX map_index map TYPE tokenbf_v1(128, 3, 0) GRANULARITY 1
)
ENGINE = MergeTree
ORDER BY id
SETTINGS index_granularity = 2
```

```
SELECT map['k1'] FROM bloom_filter_map WHERE mapContains(map, 'k1');
```

🔍    Avoid reading unnecessary granules
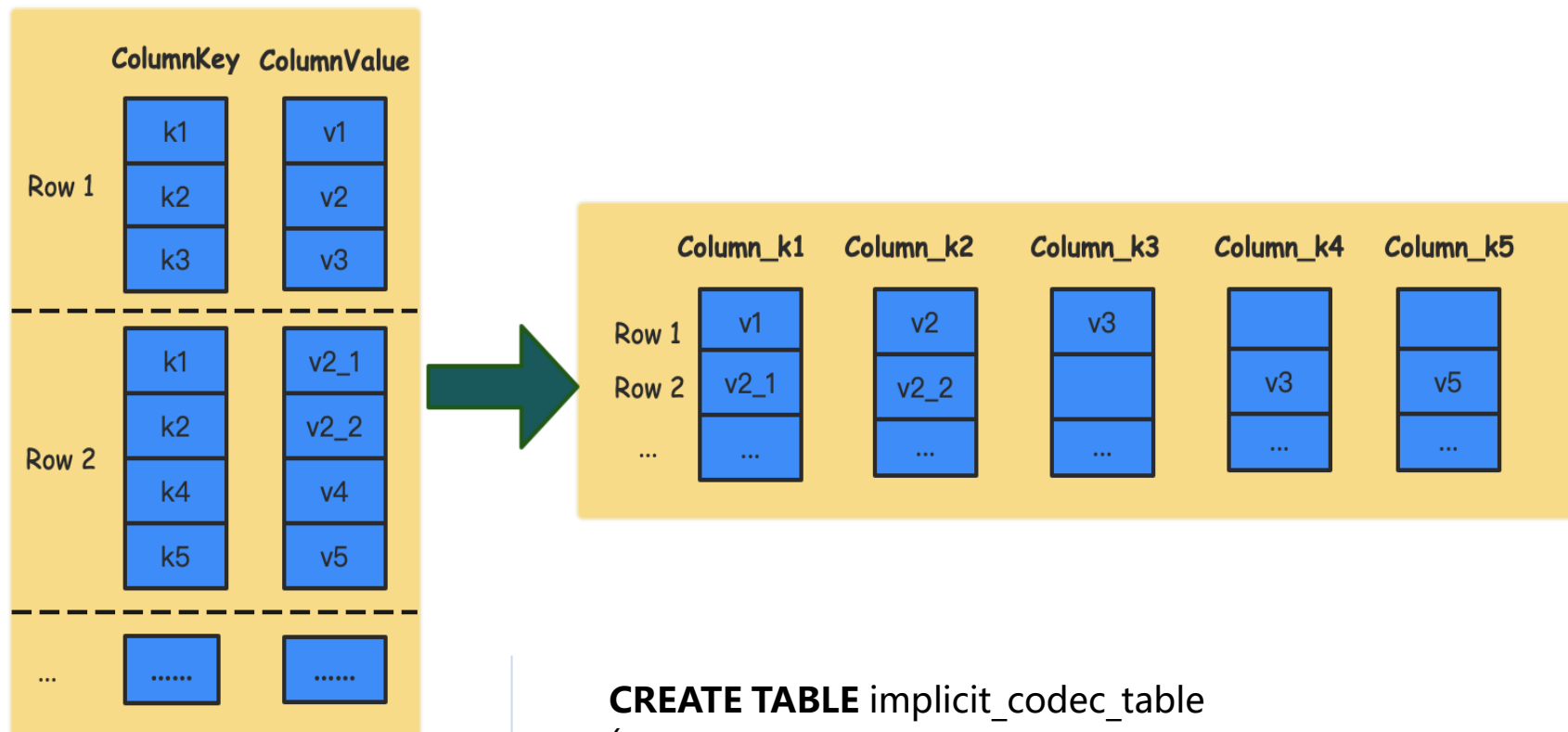
Still read all elements in map

# Implicit Map

## What is Implicit Map?

Extract every uniques key as a base column. This process is transparent to the users, so we called it as implicit columns.



```
CREATE TABLE implicit_codec_table
(
    `id` UInt32,
    `map` MapV2(String, Nullable(String)) CODEC(ZSTD(1)),
    `map2` MapV2(String, Nullable(String))
)
ENGINE = MergeTree
ORDER BY id
SETTINGS index_granularity = 8192
```

# Implicit Map

## Performance comparison

The data of two table is totally same

- ## Bytes on disk

| Table Name | Rows | Bytes on disk |
|---|---|---|
| normal_map_table | 10,000,000 | 3.19 GiB |
| implicit_map_table | 10,000,000 | 3.11 GiB |

- ## Query specified key

| Table Name | Elasped | Processed | RPS |
|---|---|---|---|
| normal_map_table | 10.460 sec | 5.07 GB | 956.04 thousand rows/s 484.87 MiB/s |
| implicit_map_table | 1.450 sec | 551.34 MB | 16.89 million rows/s 380.14 MiB/s. |

## Table DDL:

```
CREATE TABLE test.implicit_map_table
(
    `id` UInt64,
    `map` MapV2(String, String)
)
ENGINE = MergeTree
ORDER BY id
```

## Query:

```
SELECT uniqExact(op)
FROM (
    SELECT map['OnePiece'] AS op
    FROM
    normal_map_table
)
```

# Implicit Map

## select all for Implicit Map

**We have tried:**

- Construct Map data from all implicit columns during query process

- We create implicit columns and also store native map
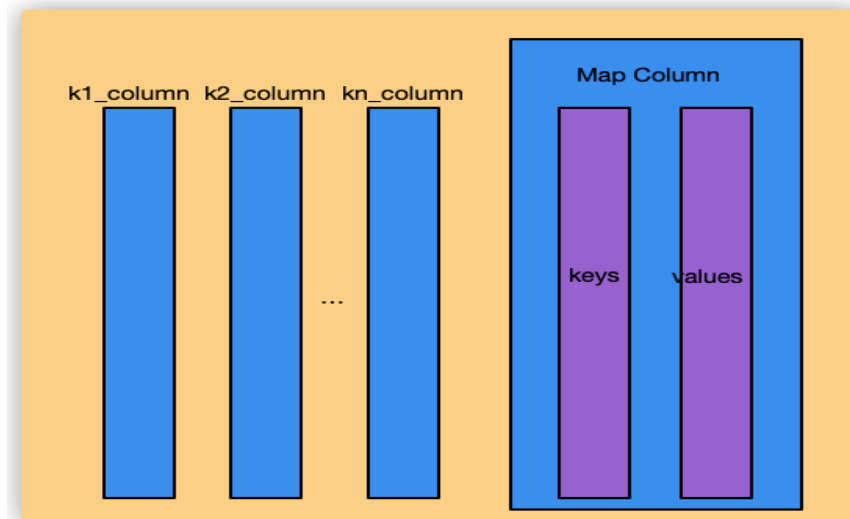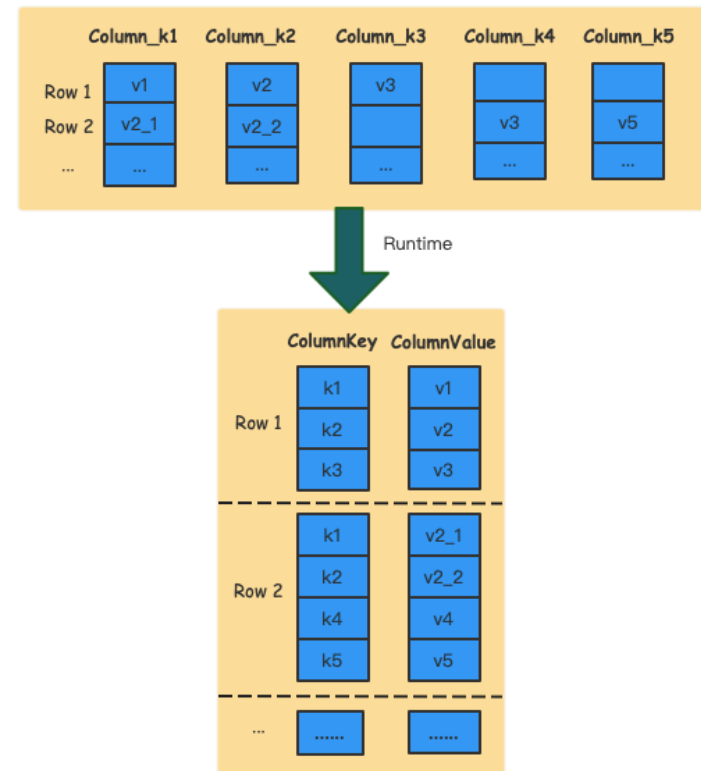
  Query with specified key will read implicit column. If users want get whole map, it will read native map column as usual.

- Specify every key

  SELECT map FROM map_table

  SELECT map['k1'], map['k2'] … map['kn'] FROM map_table

  We can get all implicit column name from system.parts.

# Thanks