

趣头条基于Flink+ClickHouse 构建实时数据平台

王金海
趣头条数据平台负责人

FLINK FORWARD # ASIA

实时即未来 # Real-time Is The Future

**FLINK
FORWARD**

自我介绍

About me

- 01/ 9年互联网历练，先后在唯品会、饿了么负责大数据开发架构工作
- 02/ 现为趣头条数据中心平台负责人。负责大数据基础设施建设（spark、presto、flink、clickhouse）、数据产品化输出（qe、horizon、kepler）、团队建设

目录

Content



业务场景与现状分析

Business scenario and
current situation analysis



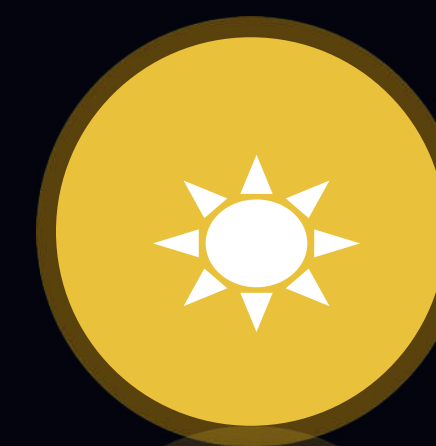
Flink-to-Hive小时级场景

Flink-to-Hive



Flink-to-ClickHouse秒级场景

Flink-to-ClickHouse



未来发展与思考

Future

业务场景与现状分析

Business scenario and current situation Analysis

01

业务场景与现状分析

Business scenario and current situation Analysis

The image displays two screenshots of the Flink Forward web interface, illustrating the process of selecting data sources for analysis.

Top Screenshot (Offline Data):

- 数据来源 (Data Source):** 离线数据 (Offline Data)
- Databases (204):** A list of databases is shown, with a red arrow pointing to the '离线数据' dropdown.
- Limit:** 100
- Database Selection:** A dropdown menu is open, showing options: presto, spark, and hive. A red arrow points to 'presto'.
- adhoc:** A dropdown menu is open, showing options: presto, spark, and hive. A red arrow points to 'presto'.

Bottom Screenshot (Real-time Data):

- 数据来源 (Data Source):** 实时数据 (Real-time Data)
- Databases (7):** A list of databases is shown, with a red arrow pointing to the '实时数据' dropdown.
- Limit:** 100
- Database Selection:** A dropdown menu is open, showing options: clickhouse, spark, and hive. A red arrow points to 'clickhouse'.
- adhoc:** A dropdown menu is open, showing options: clickhouse, spark, and hive. A red arrow points to 'clickhouse'.

业务场景与现状分析

Business scenario and current situation Analysis



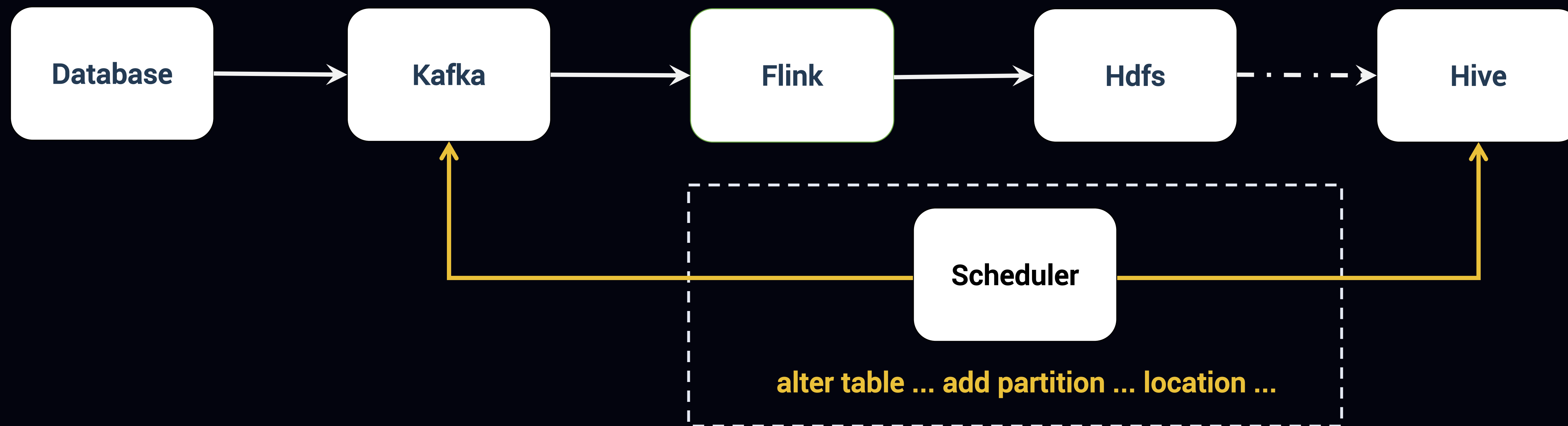
Flink-to-Hive小时级场景

Flink-to-Hive

02

小时级实现架构图

Hour level architecture



实现原理

StreamingFileSink API



```
StreamingFileSink.forBulkFormat(new Path("hdfs://data/xxx.db/yyyy/"),  
    ParquetAvroWriters.forGenericRecord(schema))  
    .withBucketAssigner(new EventTimeBucketAssigner())  
    .withBucketCheckInterval(10 * 60 * 1000L)  
    .build();
```

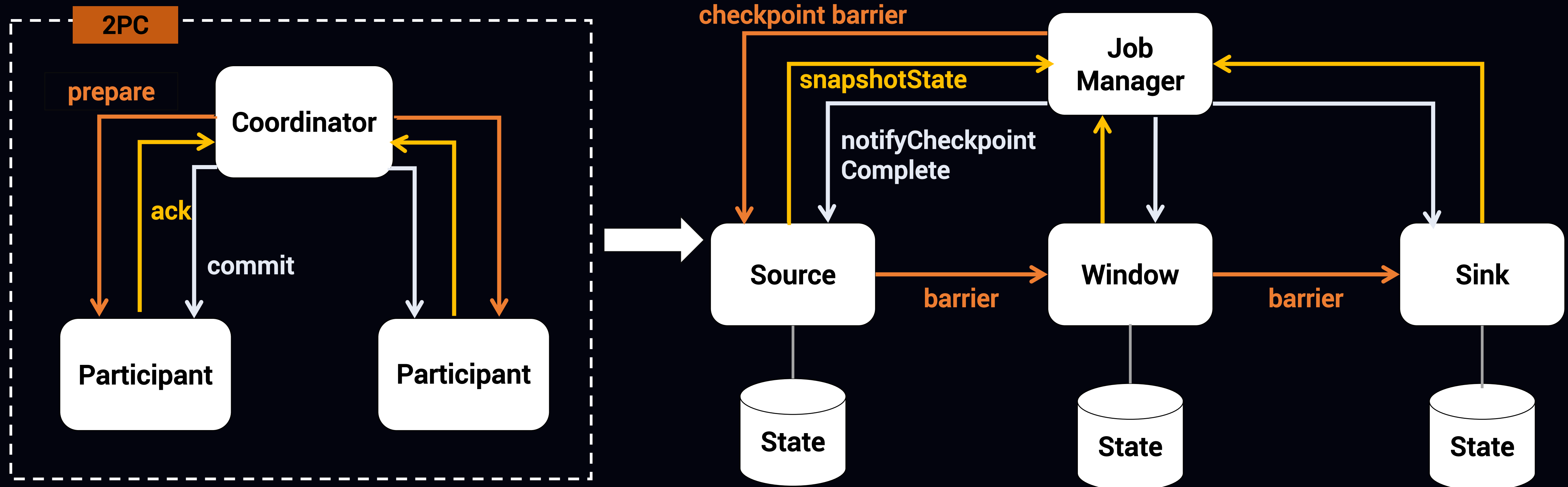
1. **forBulkFormat**支持avro、parquet格式

2. **withBucketAssigner**自定义按数据时间分桶

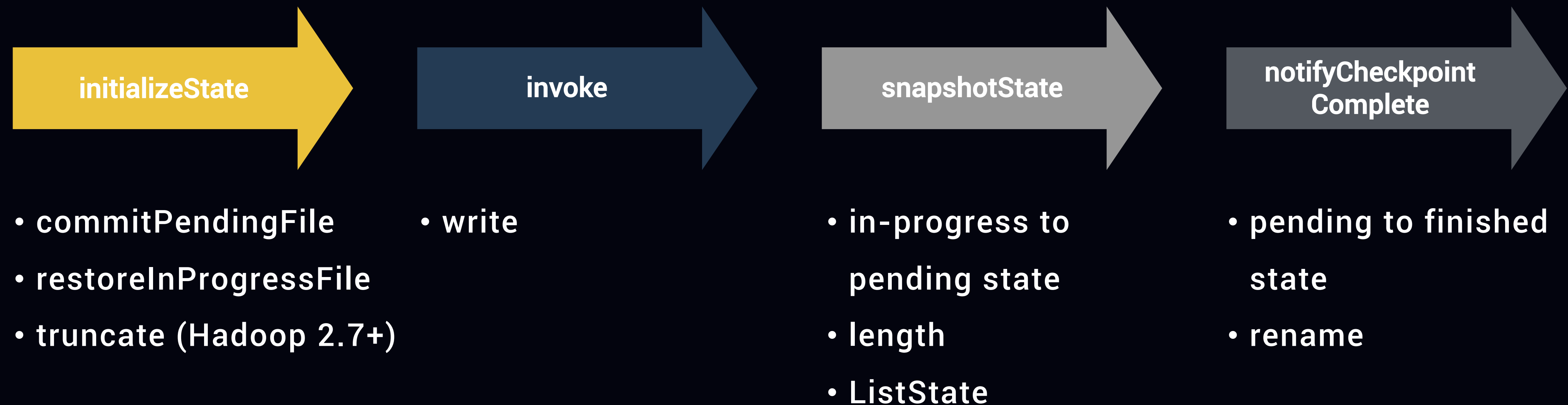
3. **OnCheckpointRollingPolicy**

4. **Exactly-Once**语义实现

Exactly-Once



StreamingFileSink implements CheckpointedFunction, CheckpointListener



跨集群多nameservices

Cross-cluster multi-nameservices

```
<?xml version="1.0" encoding="UTF-8"?>
<configuration>
  <property>
    <name>dfs.nameservices</name>
    <value>stream,data</value>
    <final>true</final>
  </property>

  <!-- stream实时集群的namenode ha配置 -->

  <!-- data离线集群的namenode ha配置 -->
</configuration>
```


多用户写入权限

Multi-user write permission

HadoopFsFactory.java

```
@Override
public FileSystem create(Uri fsUri, String user) throws IOException {
    try {
        return UserGroupInformation.createRemoteUser(user).doAs(new
PrivilegedExceptionAction<FileSystem>() {
            @Override
            public FileSystem run() throws Exception {
                return create(fsUri);
            }
        });
    } catch (InterruptedException e) {
        throw new IOException(e);
    }
}
```

```
StreamingFileSink.forBulkFormat(new Path("hdfs://data/xxx.db/yyyy/"),
    ParquetAvroWriters.forGenericRecord(schema))
    .withBucketAssigner(new EventTimeBucketAssigner())
    .withBucketCheckInterval(10 * 60 * 1000L)
    .withBucketUser("stream")
    .build();
```



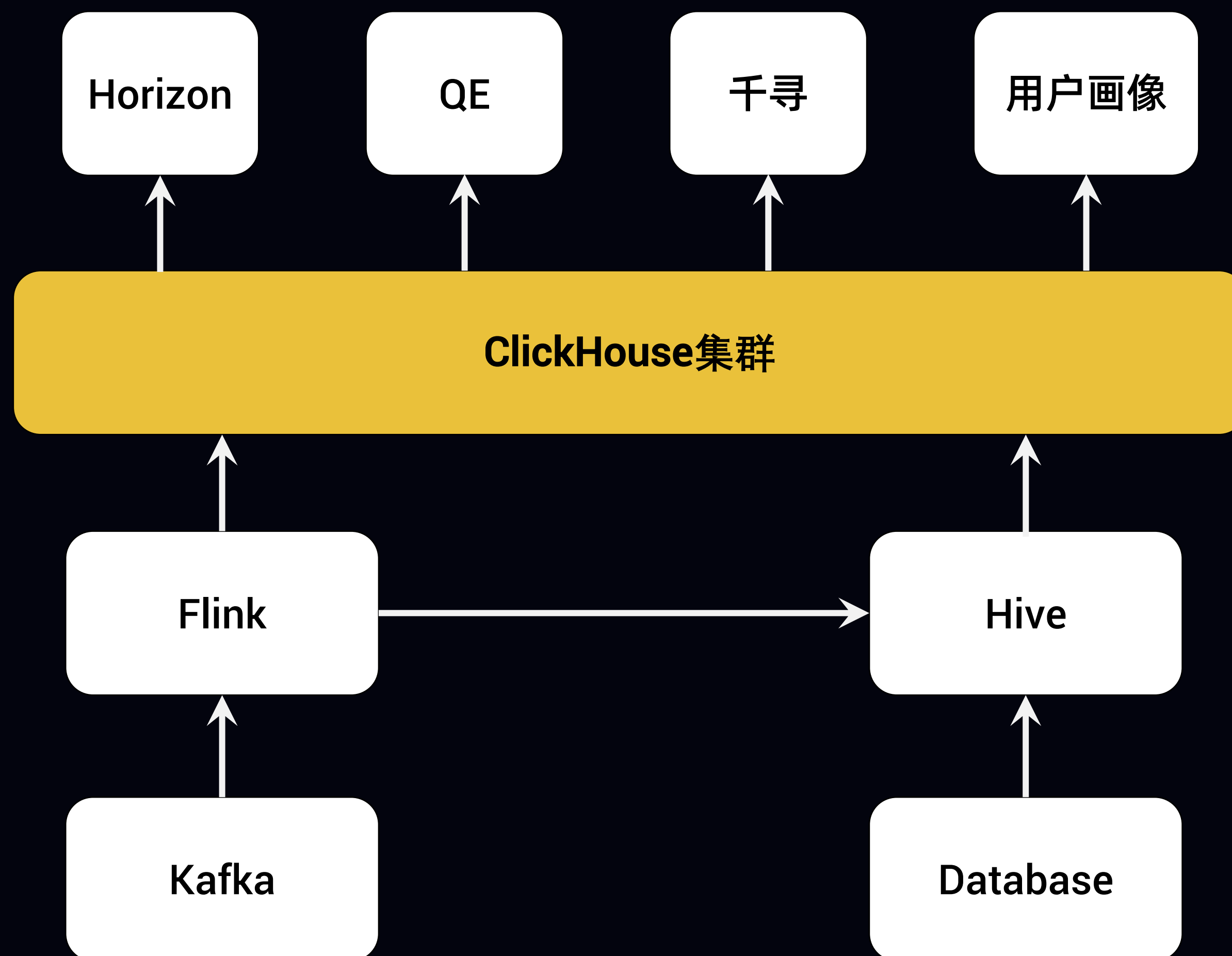
Flink-to-ClickHouse秒级场景

Flink-to-ClickHouse

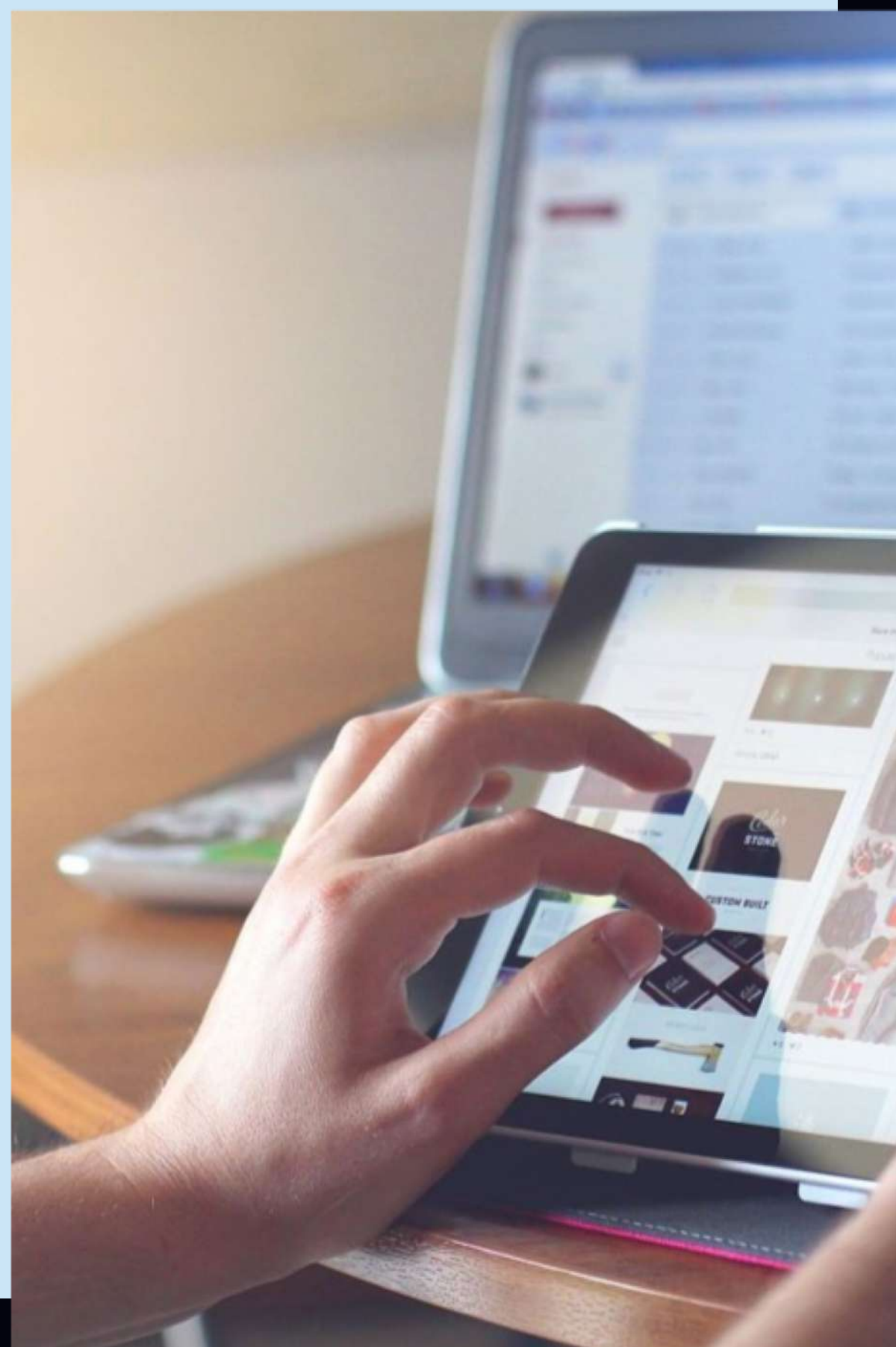
03

秒级实现架构图

Second level architecture



Why Flink+ClickHouse



1

指标实现支持sql化描述

2

指标的上下线互不影响

3

数据可回溯，方便异常排查

4

计算快，一个周期内完成所有指标计算

5

支持实时流，分布式部署、运维简单

100台

100+台32核128G
3.5T SSD

2000亿

2000+亿/天

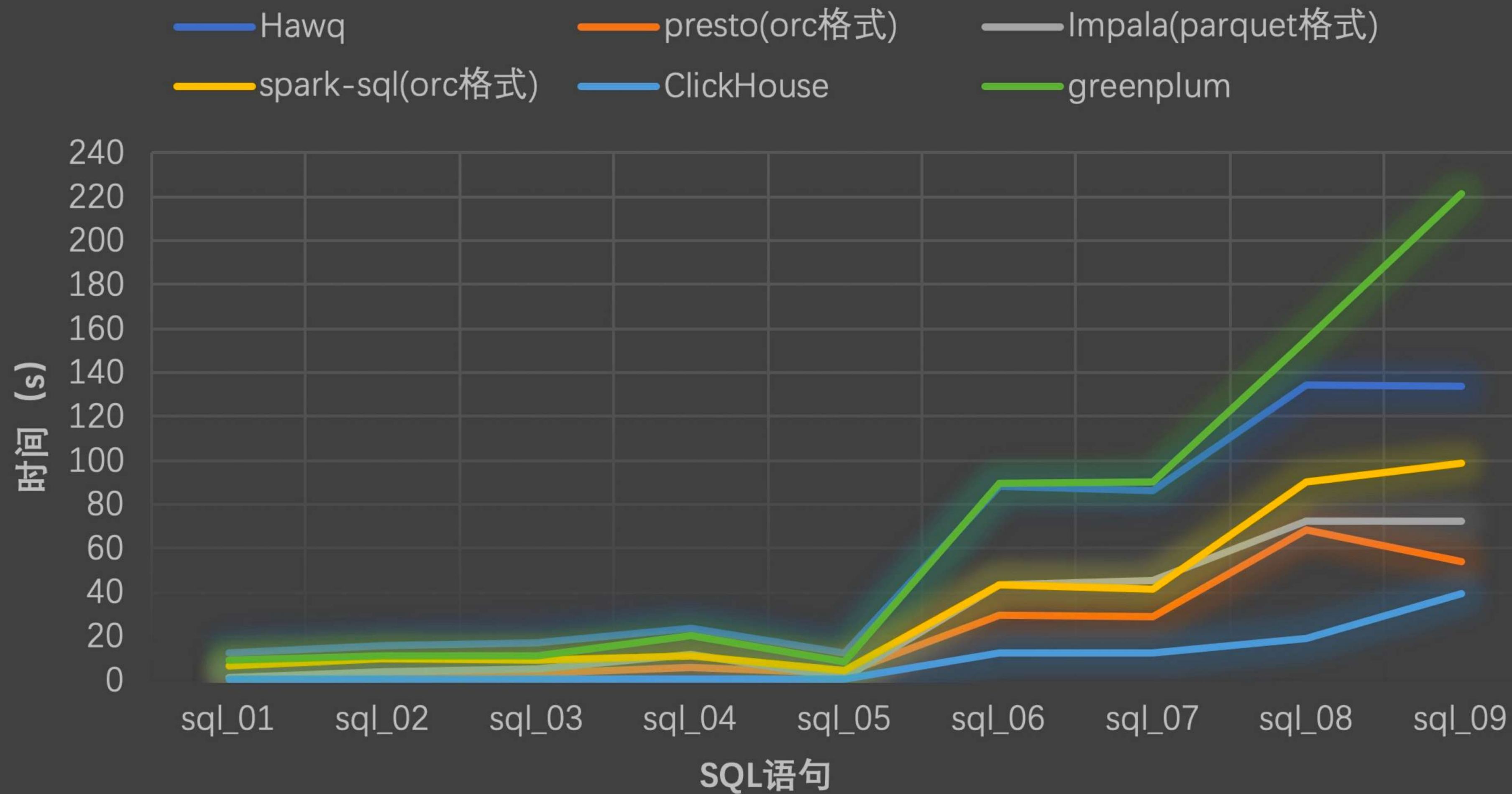
21万

21+万次查询/天

1S

80%查询1S内完成

单表测试




```
bjg-data-platform-clickhouse-01 :) select cmd, count(1) as cnt from quiver_client_collect_v1_local_20191106 group by cmd order by cnt desc limit 10;
```

```
SELECT  
  cmd,  
  count(1) AS cnt  
FROM quiver_client_collect_v1_local_20191106  
GROUP BY cmd  
ORDER BY cnt DESC  
LIMIT 10
```

cmd	cnt
1001	54003480
3001	4277208
5005	301731750
15000	17546291
1000	166471366
4001	13414222
5008	129699185
8002	113524703
2002	41265766
8007	40227632

26亿
count + group by + order by
3.6S

10 rows in set. Elapsed: 3.624 sec. Processed 2.62 billion rows, 34.65 GB (722.05 million rows/s., 9.56 GB/s.)

Why ClickHouse so fast

01/ 列式存储 + LZ4、ZSTD数据压缩

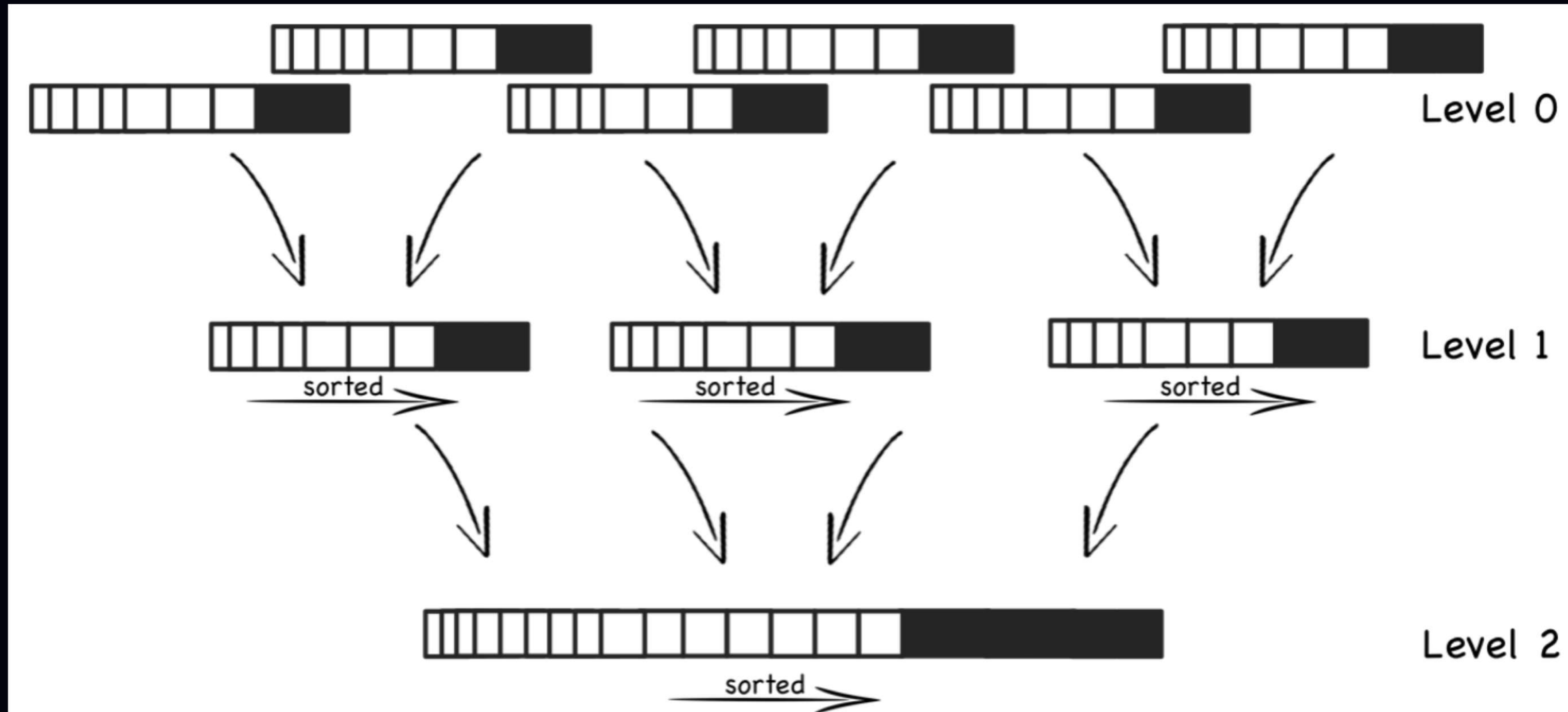
02/ 计算存储本地化 + 向量化执行

03/ LSM merge tree + Index

04/ SIMD + LLVM优化

05/ SQL语法、UDF完善

MergeTree



Compaction continues creating fewer, larger and larger files

ClickHouse Connector

n1

写Local table;
读Distributed table

n2

5~10W batch;
5s interval

n2

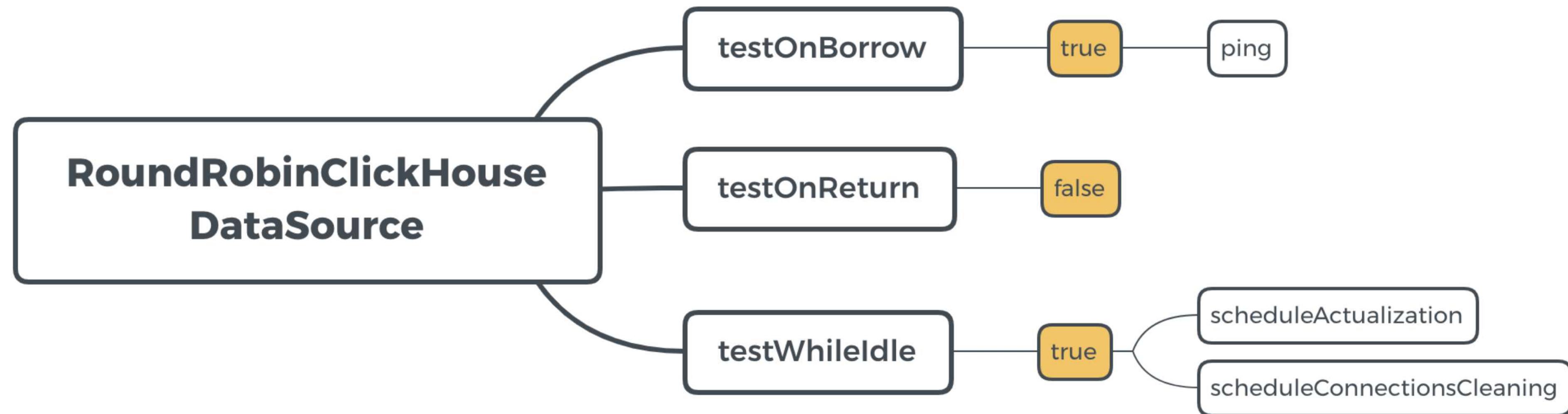
RoundRobinClickHouse
DataSource

BalancedClickHouseDataSource

01/ `scheduleActualization`

02/ `scheduleConnectionsCleaning`

RoundRobinClickHouseDataSource



Backfill



01/ Flink任务小时级容错

单击此处添加文本具体内容

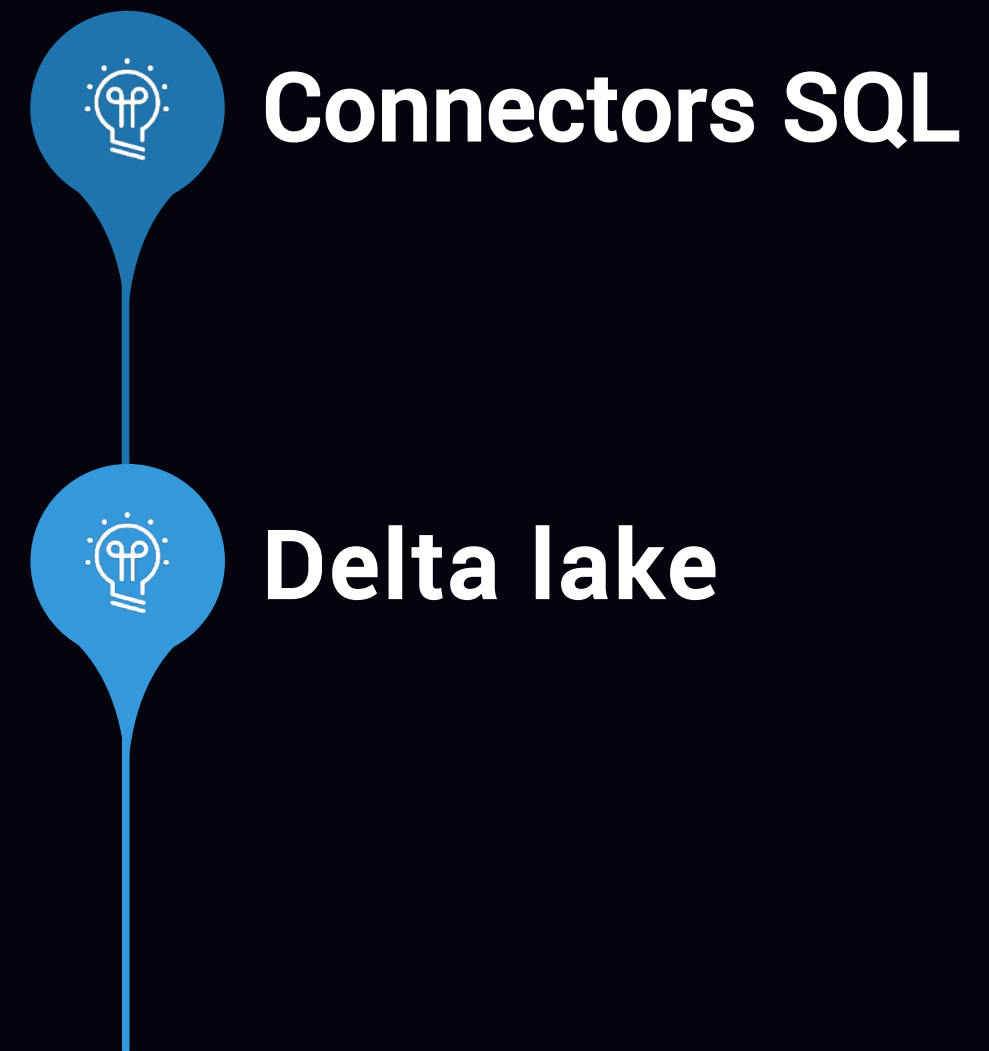
02/ ClickHouse集群小时级容错

单击此处添加文本具体内容

未来发展与思考

New Future on Cloud

04





THANKS