

# Clickhouse in Telecom ( From 0 to 1 )

— Dataliance

# Agenda

Application Scenarios in G-Net data analysis

Evolution of G-Net bigdata platform architecture

Comparison with Clickhouse and Hadoop

Make a Migration to Clickhouse

Technology Architecture After Clickhouse

Clickhouse on public cloud

# 中国电信G网数据分析典型应用场景

## ■多维度的用户行为特征分析

对数据业务流量按区域、区域类型、区域场景、业务、终端、SP等多维度进行组合分析，以便掌握用户行为特征



## ■基于用户个体的业务消费模型分析

通过业务模型分析、终端业务分析、以及用户区域分析，建立起从业务服务提供端至用户终端的分析手段，再结合经分、BOSS等系统中的业务信息以及用户信息后，就能够实现基于用户个体的业务消费模型分析，进而达到为市场实现精细化营销的目的



# 中国电信G网数据分析总体技术架构



**ETL层**

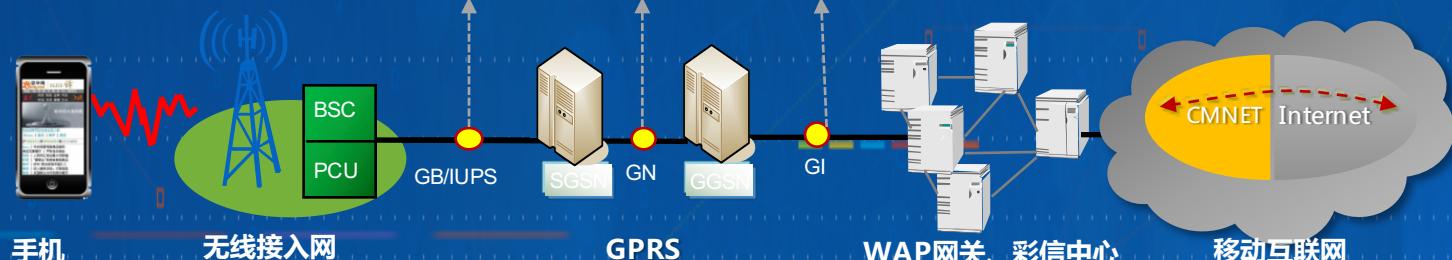


按照业务规则清洗数据，完成各类业务数据的抽取、转化、加载

**采集层**

数据捕获 → 会话管理 → 信令解析 → 业务识别 → CDR合成

**网络层**



# Data Scale in telecom industry

## Data processing volume:

- Ingesting from different data sources e.g. base station、 monitor、 backbone network
- 5billions Entries, ~700G/Day now, to 20TB/Day in the furture
- Generate realtime analysis report

# 基于位置的服务, 网络优化



- 网优
  - 例如. 重新路由来电到另外一个基站, 如果检测到有网络拥塞存在
- 基于位置的营销
  - 匹配点击事件到订阅者资料; 如果匹配 说明是位置敏感性广告
- 挑战: 交互式实时控制台
  - 简单的规则 -  $(CallDroppedCount > threshold)$  然后告警
  - 或者, 复杂 (OLAP 查询)
  - $TopK$ , 趋势分析, *Join*查询, 与历史数据关联

目前的查询场景

需要强大的查询分析引擎

# Comparison with Clickhouse and Hadoop

## Why Choose Clickhouse? Drop Hadoop

- Hadoop Cluster has a poor performance , that is too slow to be valid.
- Hadoop Cluster is fat.
- Cant execute to query data (PB) in real time.

## Clickhouse

- Rich functions
- Perfect performance
- Structure flat not fat
- Flexible way in execution

# Clickhouse功能特点与优势

## 数据库内压缩

采用了业内领先的压缩技术，提高性能的同时，显著地减少存储数据所需的空间。客户可以将所用空间减少3-10倍，并提高有效的I/O性能。

## 千万亿字节规模的数据加载操作

高性能的并行数据装载器可以在所有节点上同步执行操作，装载速度超过50W条/秒。

## 随地访问数据

不管数据的位置、格式或存储介质如何，都可以从数据库向外部数据源执行查询操作，并行向数据库返回数据。

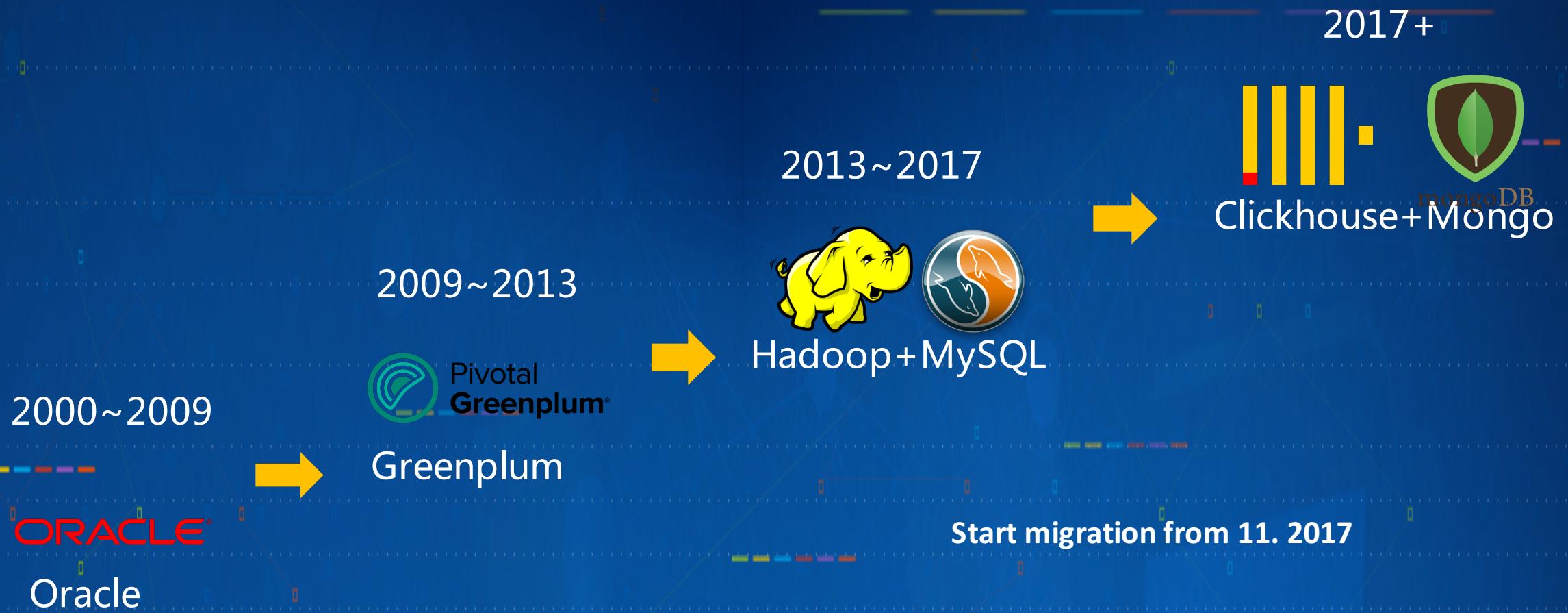
## 动态扩展

对数据仓库进行便捷的小规模或大规模扩展，同时避免高成本的设备或SMP服务器升级。

## 集中管理

提供集群级管理工具和资源，帮助管理人员像管理一台服务器一样管理整个多维实时分析平台。

# Evolution of G-Net bigdata platform architecture



# Technology Architecture Before Clickhouse

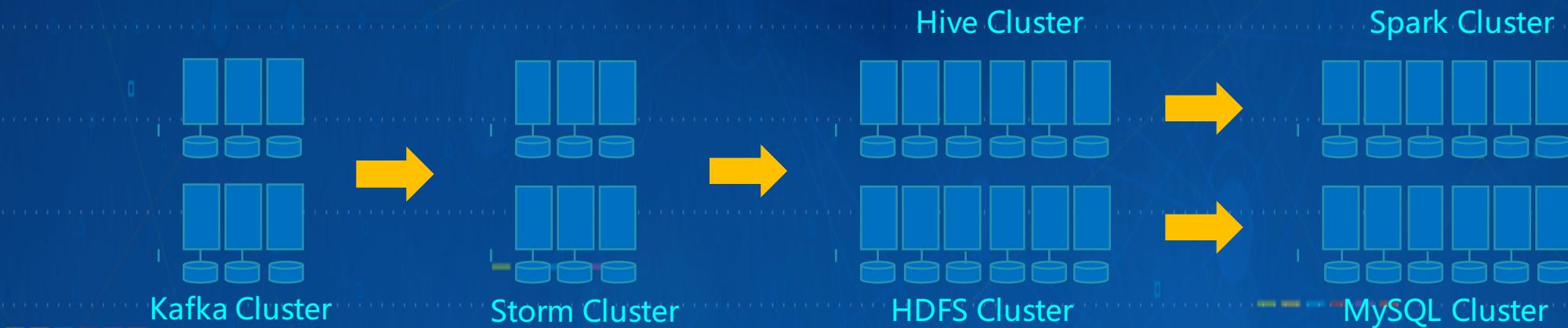
Legacy Architecture: Kafka+Storm+Hive+Spark+MySQL

Kafka: collect and aggregate data

Storm: wrangle data

Hive: ad-hot query data

Spark: analysis offline



# Make a Migration to Clickhouse

Remove HDFS、Hive、Spark Solutions

- Slow point—full scan、data sorting
- Not suitable for Aggregated online analysis
- Not suitable for Ad-hot query
- Bad interact experience

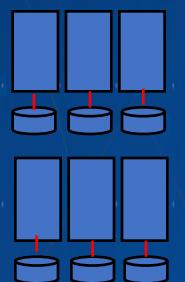


Speed up ~560X!

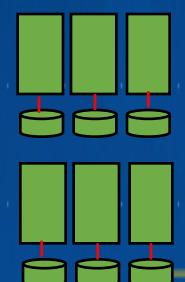
Elapsed Time 80s -0.3s

# Technology Architecture After Clickhouse

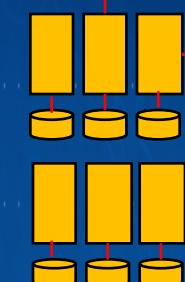
- Memory Table
- ReplicatedMergeTree Table
- Distributed Table



Kafka Cluster

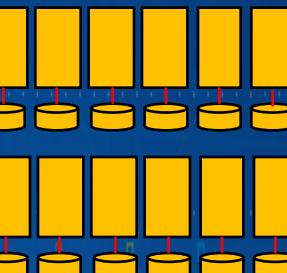
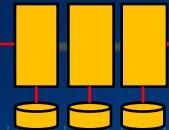


Mongo Cluster

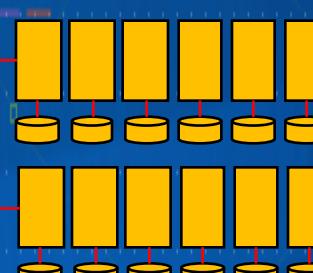


CK Client

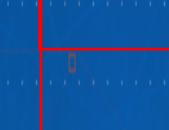
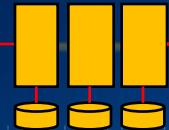
Zookeeper  
Nodes



ClickHouse  
15 nodes

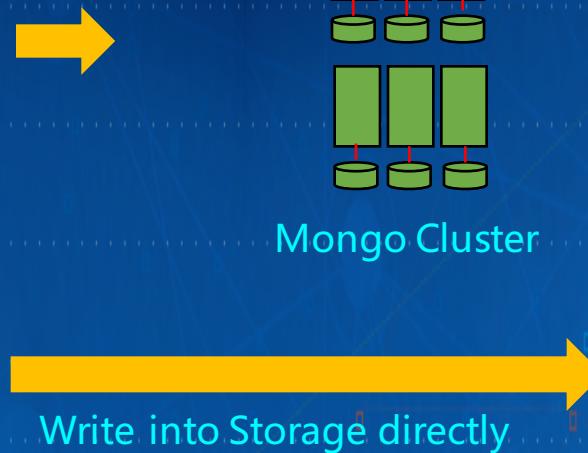
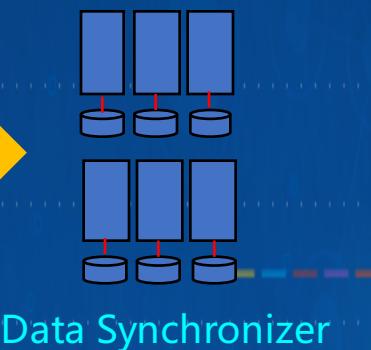
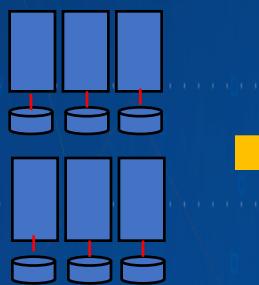


ClickHouse  
15 nodes



# Technology Architecture In the future

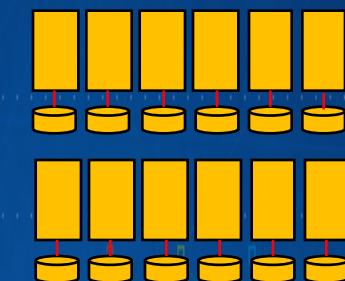
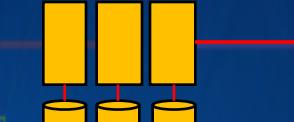
- Memory Table
- ReplicatedMergeTree Table
- Distributed Table



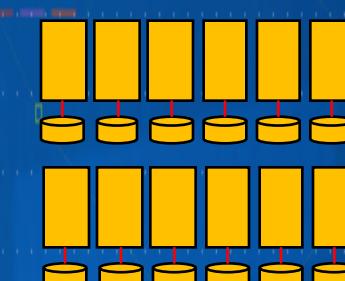
Mongo Cluster

Write into Storage directly

Zookeeper  
Nodes



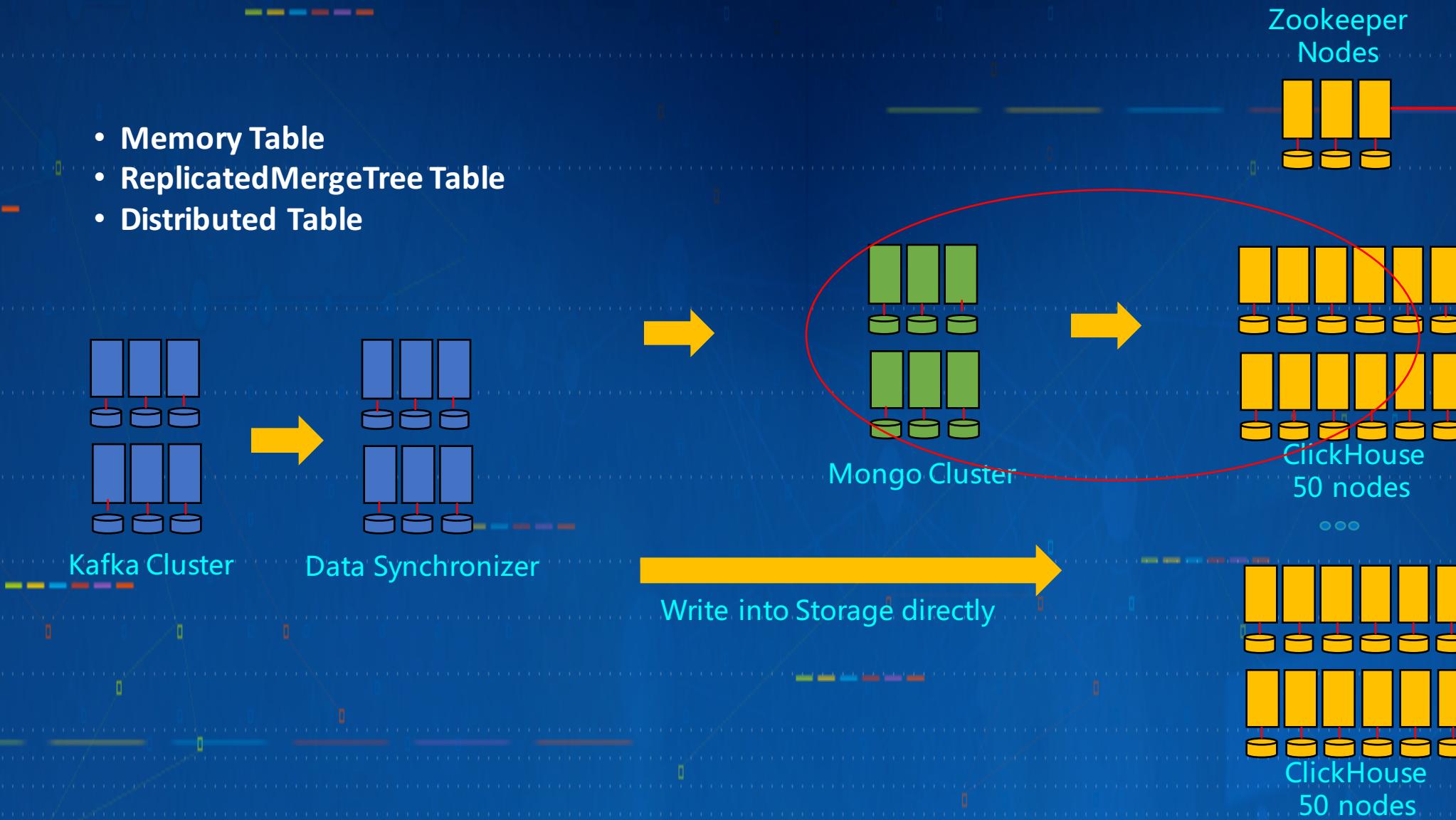
ClickHouse  
50 nodes



ClickHouse  
50 nodes

# Technology Architecture In the future

- Memory Table
- ReplicatedMergeTree Table
- Distributed Table

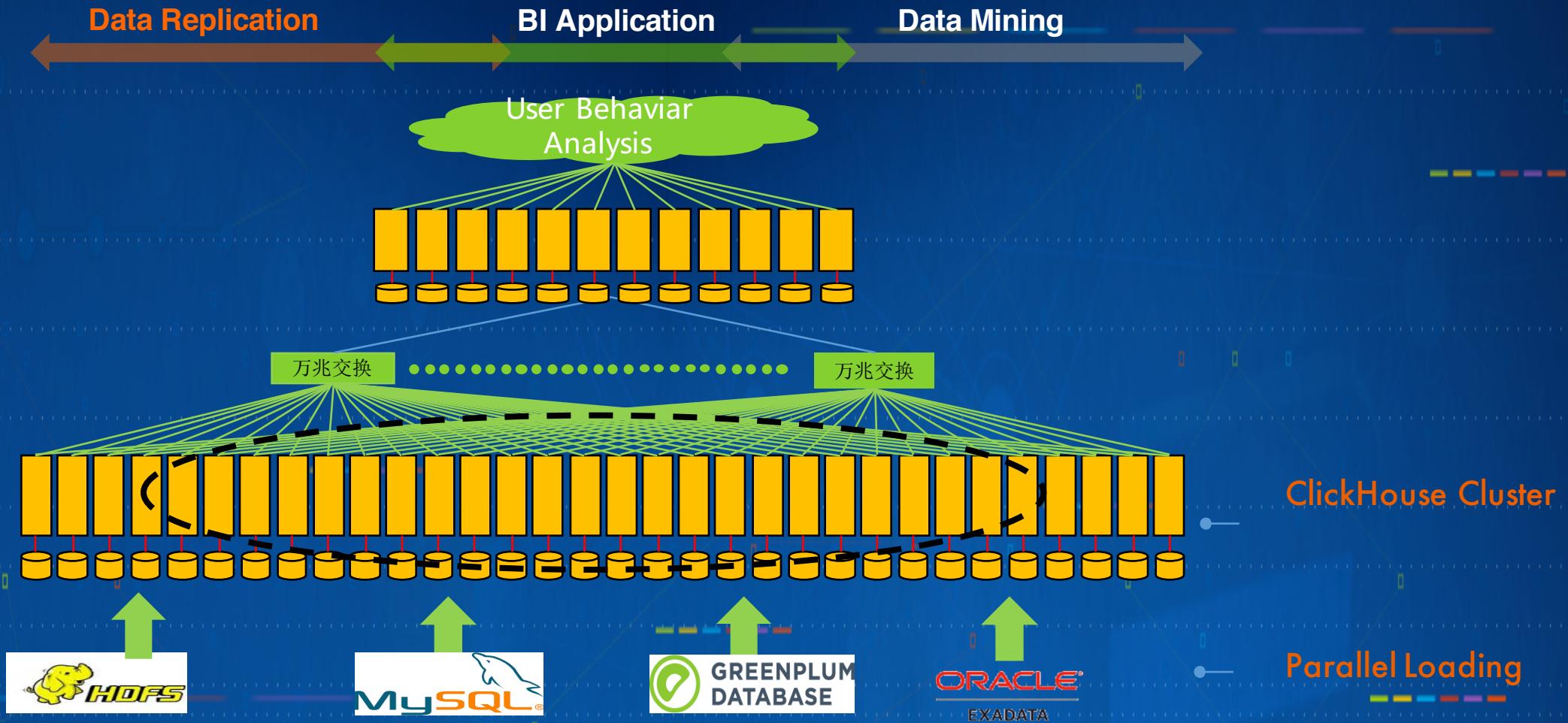


# Technology Architecture In the future

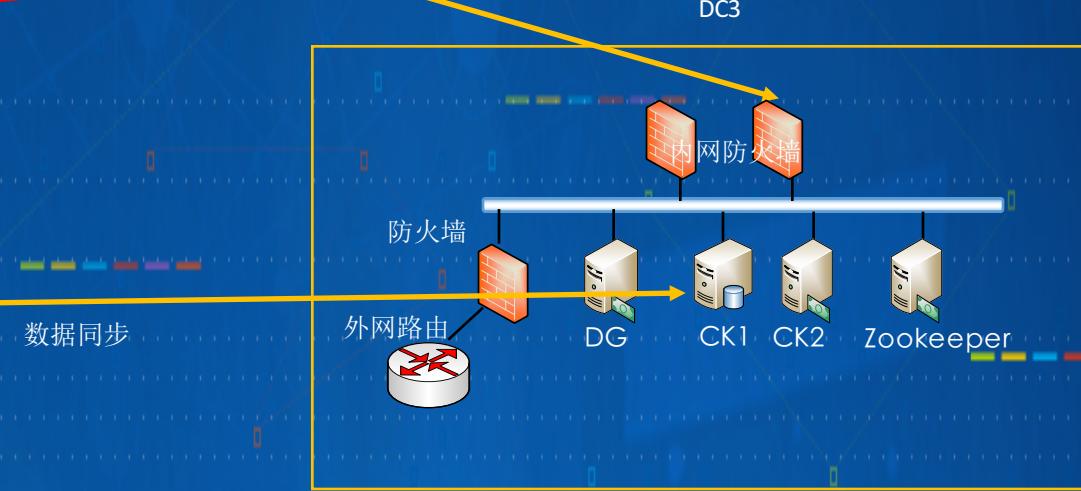
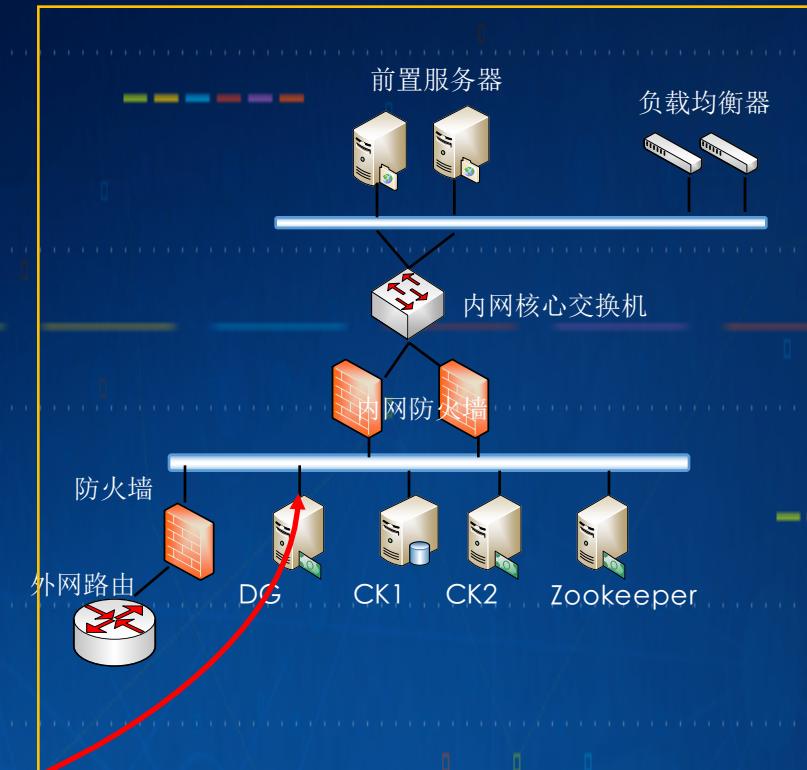
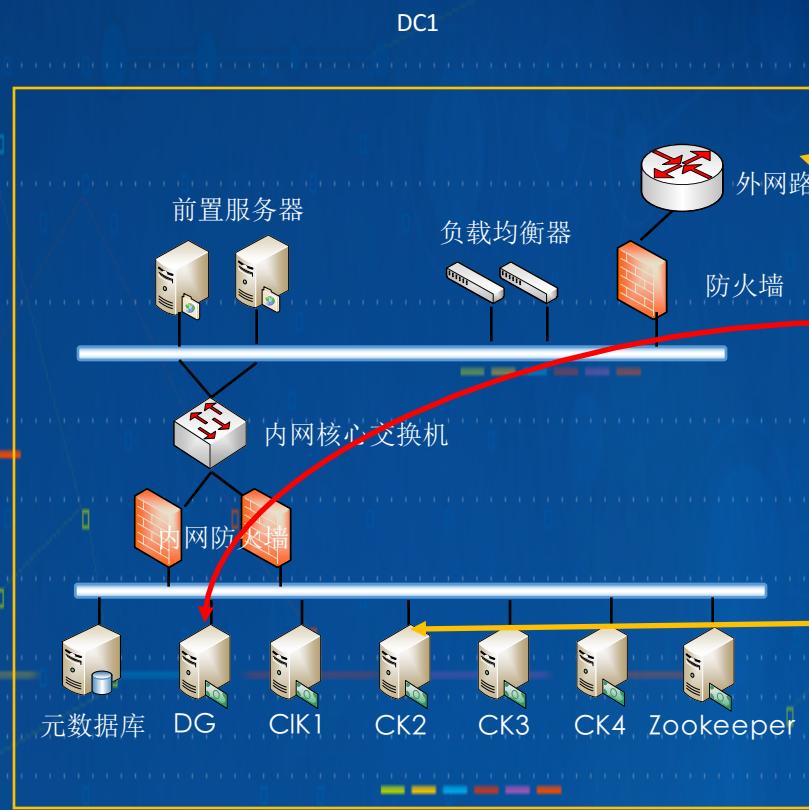
- Create a realtime sync tool between Mongo and Clickhouse
- Create a realtime sync tool between MySQL and Clickhouse



# Technology Architecture In the future



# Disaster Recovery in 3 DCs



数据同步

数据同步

# Need to work

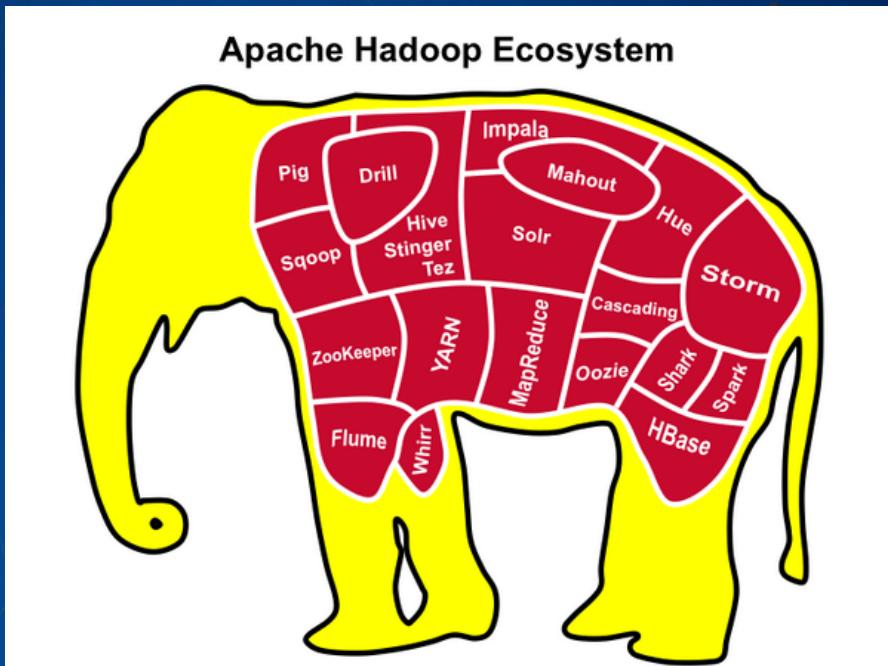
- DML SQL(Update , Delete)
- SQL99/2003
- Map-Reduce
- Automated Operation Tools
- Clickhouse on HDFS(like Hawq)

# Shortage Use Case

Lack point: no server node state information

cluster	shard_num	shard_weight	replica_num	host_name	host_address	port	is_local	user	default_database
bip_ck_cluster	1	1	1	cloudera3	x.x.x.x	9000	0	default	
bip_ck_cluster	2	1	1	cloudera2	x.x.x.x	9000	1	default	
bip_ck_cluster	3	1	1	cloudera4	x.x.x.x	9000	0	default	
test_shard_localhost	1	1	1	localhost	127.0.0.1	9000	1	default	

# Kill Hadoop using clickhouse

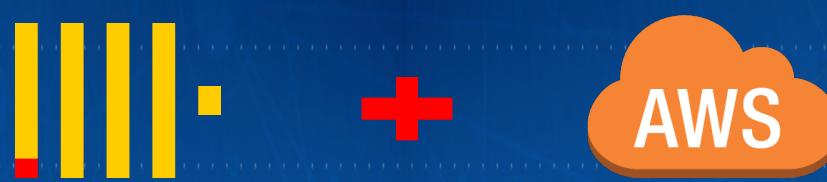


一只大象拆分后，有价值的东西所剩不多

- Kafka
- HDFS
- Spark

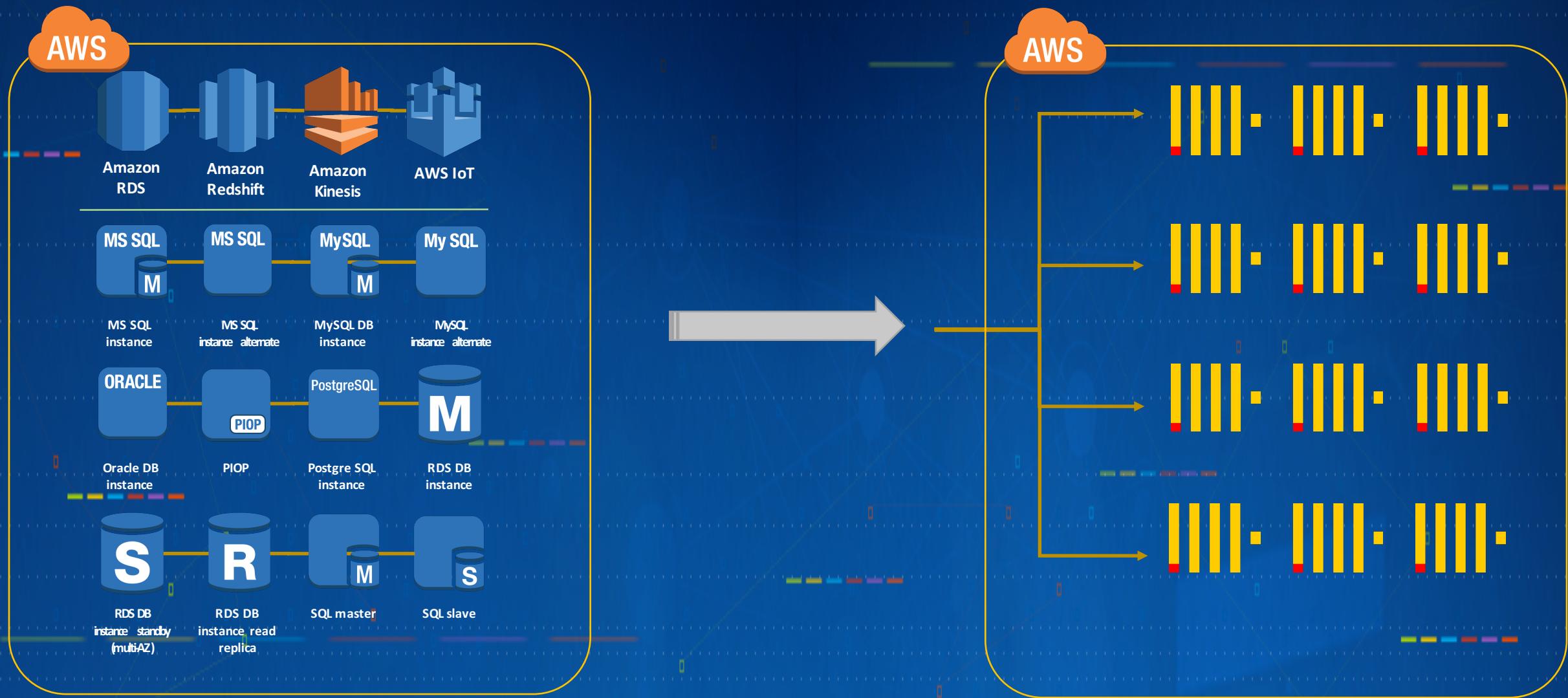


# ClickHouse on AWS

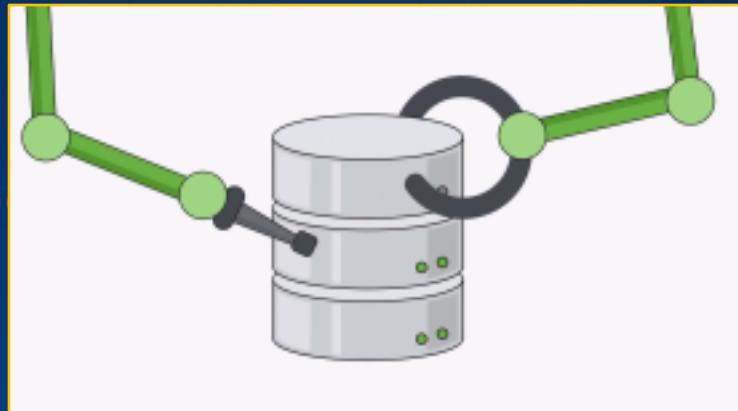


提供 ClickHouse Cloud 云服务

# Clickhouse On AWS



# Clickhouse as a service



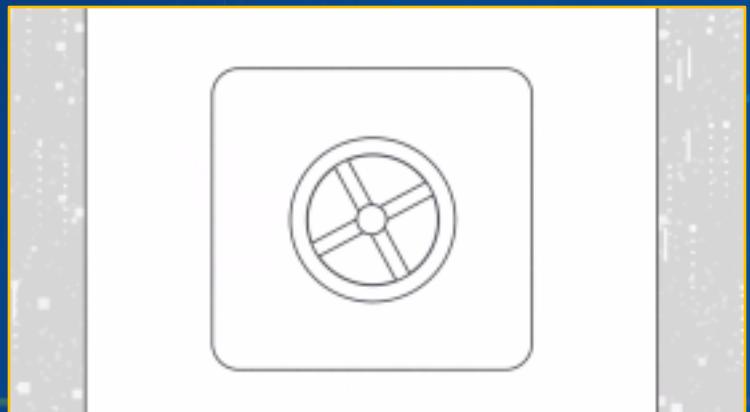
自动化



按需使用



水平扩展



安全可靠



高可用



自动备份

公司招聘：

Database kernel developer

Clickhouse integration developer

欢迎加入我们