

ClickHouse at PhysicsWallah

Empowering Real-Time Analytics at scale

Utkarsh G. Srivastava
Software Development Engineer III
PhysicsWallah

Agenda

- Problem statement
- Few solutions
- Implementation
- Challenges
- Impact
- Areas of improvement
- QnA

The problems

Database

Data was stored as aggregates in MongoDB, limiting detailed analytics and product capabilities.

Cost

Need of a scalable, low-cost solution to stay aligned with our cost-sensitive business model serving few 100s of TBs of data.

Time

With timelines already slipping, to build and ship fast—speed was non-negotiable.

Few Solutions

The solutions

Database

- InfluxDB / TimescaleDB
- Pinot
- Druid
- ElasticSearch
- ClickHouse

Cost

- Snowflake/BigQuery
- PostgreSQL with extensions
- ClickHouse

Time

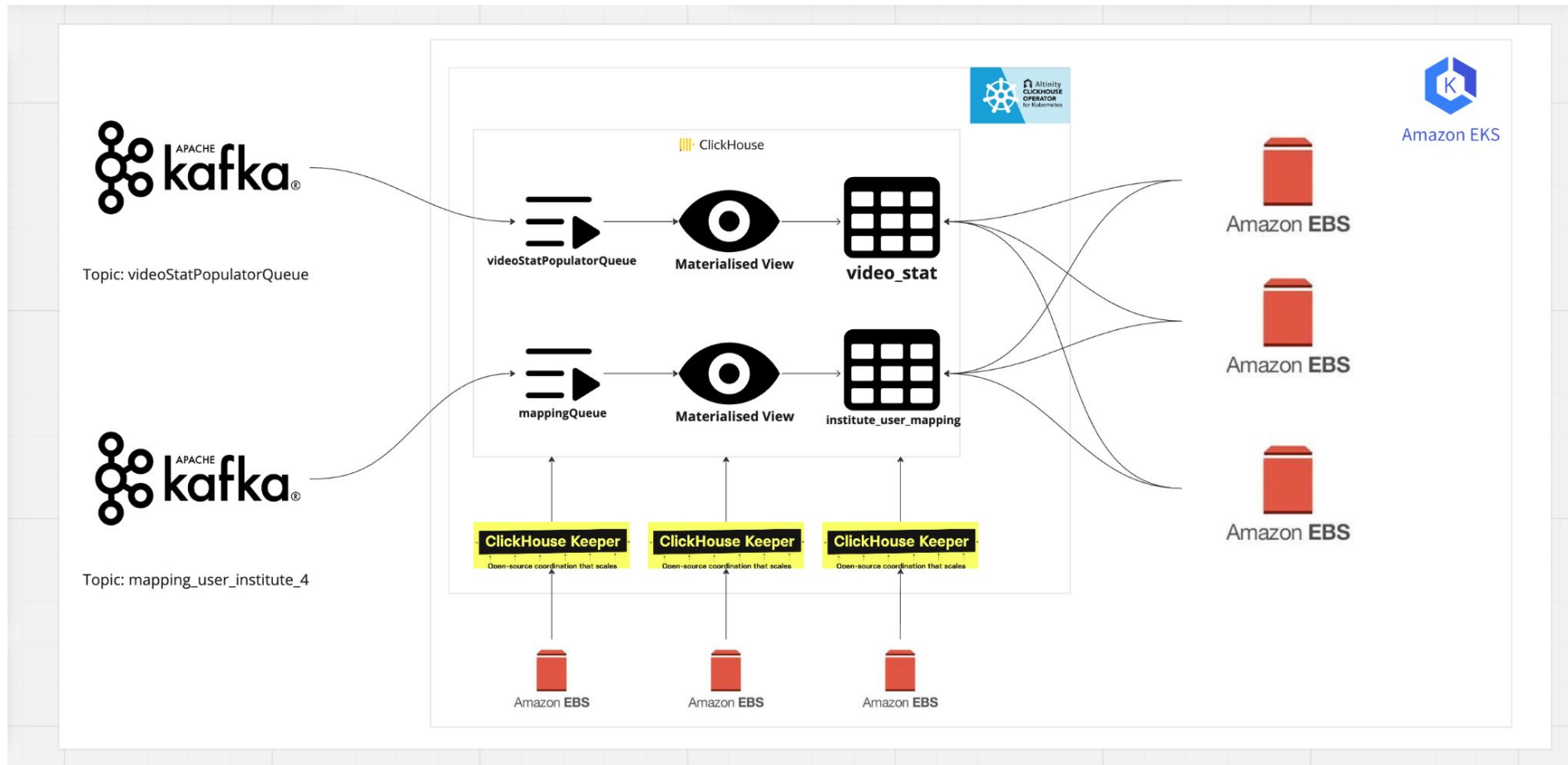
- Custom solution over MongoDB
- Stream processing stack like Flink
- ClickHouse

Implementation

Things at work

- AWS Cloud
- EKS
- EBS
- Altinity ClickHouse Operator
- Altinity ClickHouse CRDs

Typical Flow



Few Challenges

- Choosing appropriate infra
- ClickHouse Cluster Configuration
- Data Ingestion
- Data Manipulation
- Cluster Performance and Visualisation
- Backup and DR

And their solutions

- Choosing appropriate infra - using on-demand instances and node tainting and graviton nodes
- ClickHouse Cluster Configuration - adapting to altinity's way of supplying CH configs in the CRD
- Data Ingestion - using kafka tables to ingest data from kafka topics to CH
- Data Manipulation - using mat views to manipulate data and using running aggregates
- Cluster Performance and Visualisation - using Altinity Grafana Dashboard
- Backup and DR - a blend of full and incremental backups

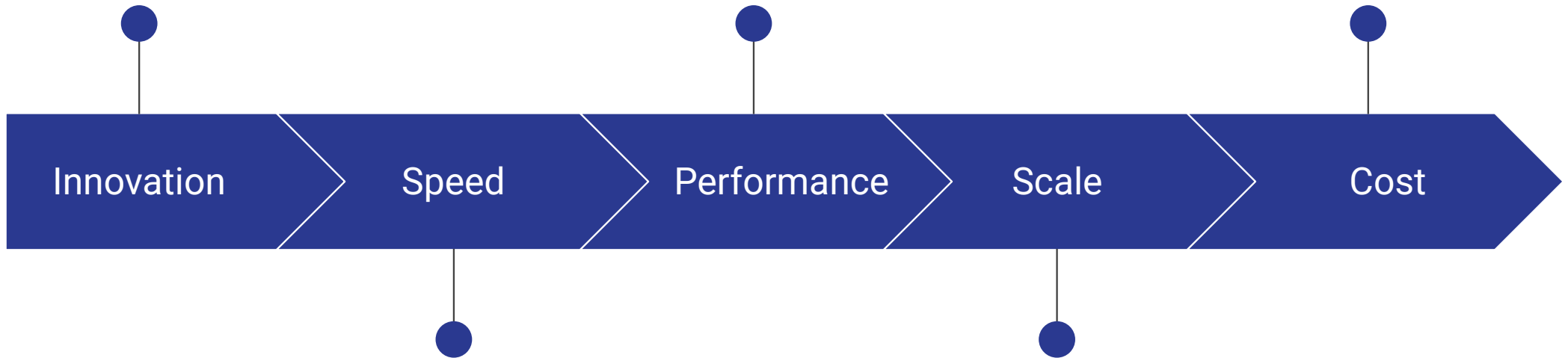
Impact



currently, more than 5 critical customer facing products are now powered by ClickHouse

Able to ingest 50 million records per day with a speed of 20k records per second

Able to create a raw layer of data and reduce overall cost by deduplication



Able to serve 3K QPS with 190ms P99 without any caching on a table containing more than 5 billion rows

Serving more than 10 million users with 1.5 million DAU and reduce overall cost by self-hosting and horizontal scaling

Areas of improvement

- Cost
- Multi-tenancy
- Schema Evolution
- Backup and DR

Special Thanks

- ClickHouse Cloud & Open Source Community
- Altinity
- PhysicsWallah

QnA