# Harnessing Real-Time Power: Inside Sofascore's Data-Driven Infrastructure
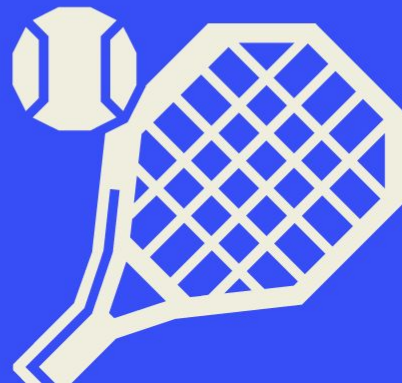
October, 2024: Nog.hr

Sofascore

Sofascore

INTRODUCTION
# Karlo Knežević

Husband

Father

35 years old, from Zagreb

Sofascore, Faculty of Electrical Engineering and Computing, Algebra University

PhD. in Machine Learning and Evolutionary Computation applying to symmetric cryptography

**Head of AI @ Sofascore**

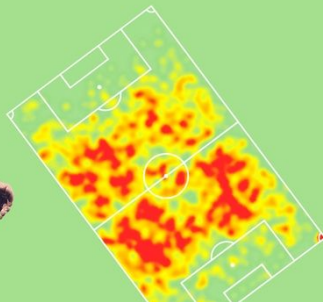**Sofascore**

# What is Sofascore

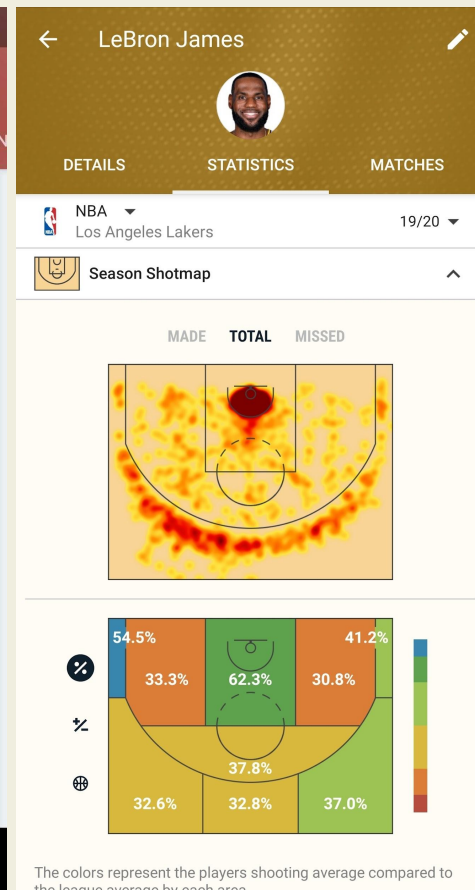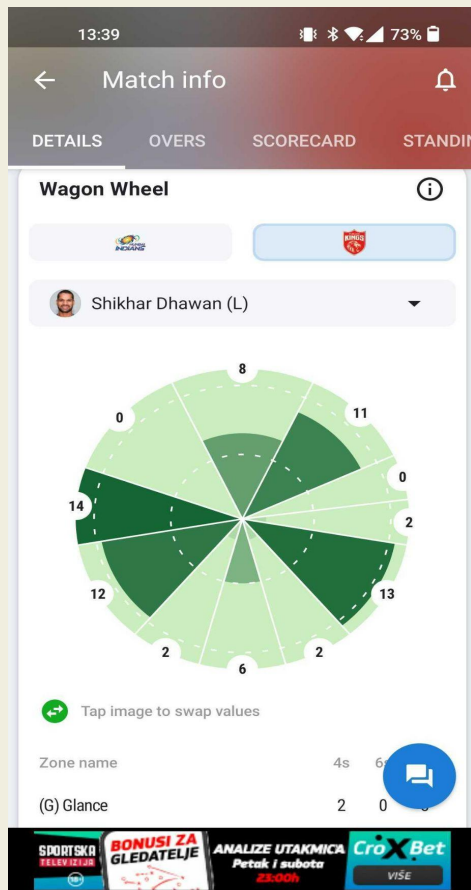# 25

SPORTS
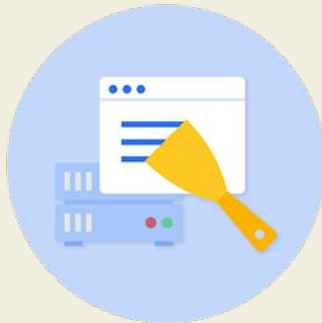
# 28

MILLION USERS

Data in
Sofascore

# SPORTS DATA

# Sports data in Sofascore

- 25 sports (for now)
- Data on players, teams, leagues, events (matches), coaches, referees...
- Insightful statistics
- Player rating
- Graphs and visualizations

# Where does this data come from?

- Buying providers' data (Opta Sports)
- Scraping data (web)
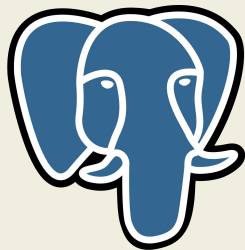- Manual inserting (Sofascore Data Team)
- Crowdsourcing



Sofascore

# How much data are we talking about?

| TABLE | ROWS |
|---|---|
| event | ~6,3M |
| player | ~900k |
| team | ~350k |
| tournament | ~17k |
| … | … |

Sofascore

# USER ACTIONS DATA

# User actions data in Sofascore

**!** **Clickstream** - records of user's actions through their app journey

● +200 actions (events)

● **Event** - important occurrence in our app that we measure

● e.g. follow_player, open_league, open_player, ad_click, app_remove, ...
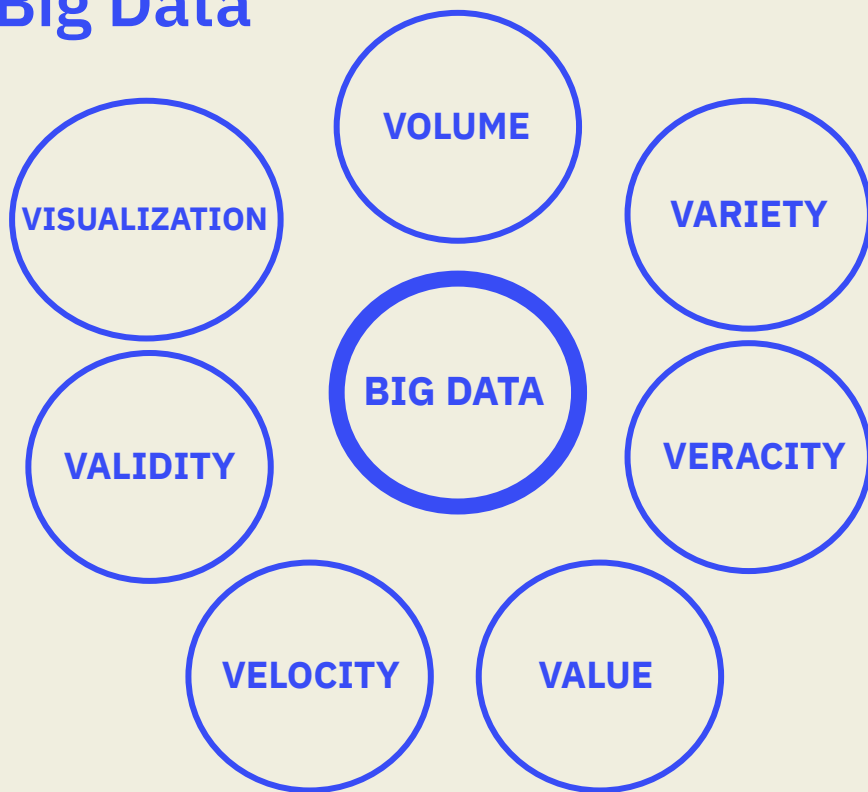
**⌐** **sofascore**

# Big Data

*"extremely large dataset that may be analysed computationally to reveal patterns, trends, and associations, especially relating to human behaviour and interactions"*
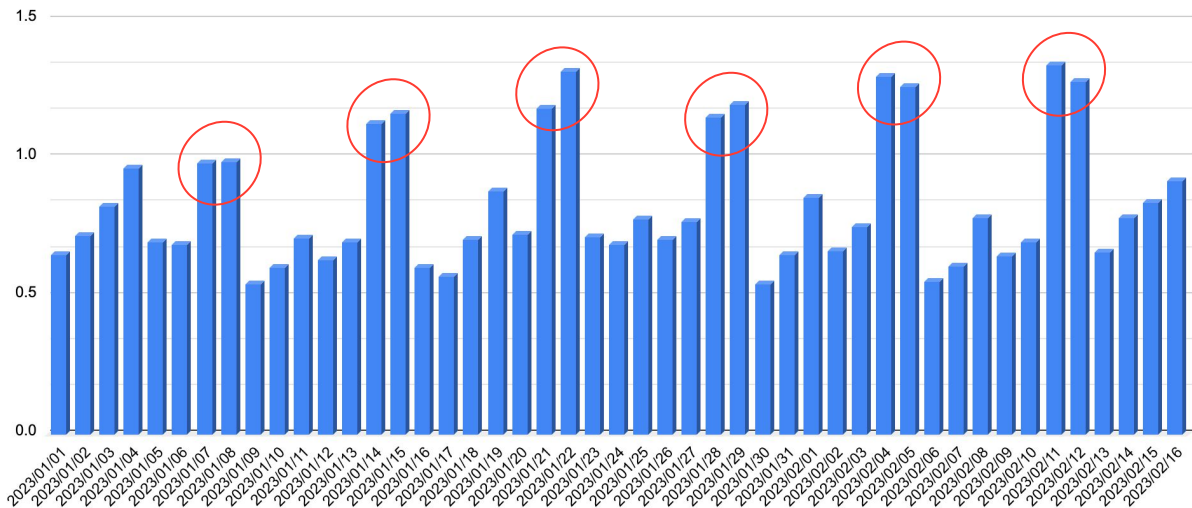
**Sofascore**

# 1PB

# ! The 7 Vs of Big Data



- VOLUME
- VISUALIZATION
- VARIETY
- BIG DATA
- VALIDITY
- VERACITY
- VELOCITY
- VALUE

sofascore

# How much data are we talking about?

VOLUME

- ~1.5PB of user actions data since February 2019
- 1 trillion rows, 268 columns in **bq.events** table in **ClickHouse**
- ~700GB arriving daily



Billions of rows per day

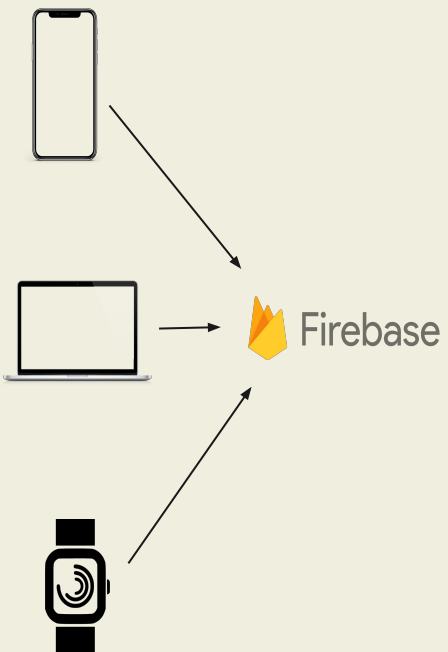**sofascore**

# ! Data Warehouse (DWH)

- A place where all enterprise data from multiple sources is consolidated into a single source of truth
- Finance, marketing, product and users data integrated into one database
- Used by data scientists and BI in order to make data-driven decisions
- Sofascore's DWH - ClickHouse

**Sofascore**

**03**

Data Pipeline

Firebase

Sofascore

**Firebase**

*"Firebase is an app development platform that helps you build and grow apps and games users love. Backed by Google and trusted by millions of businesses around the world."*
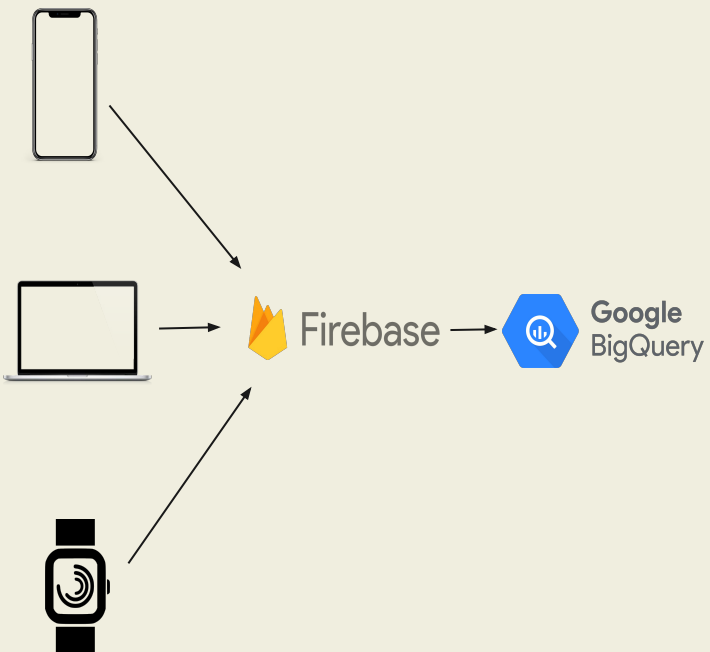
- Offering services like:
  - **Analytics**, Authentication, Databases, File Storage, Push Messages, ...

**Google Analytics**
for Firebase

- Reporting up to 500 different types of events
- Associate up to 25 parameters with each event
- Data is exported daily to **BigQuery**
  - Duplicates occurring (due to client's network issues)
  - < 2% duplicates

**VERACITY**

sofascore

**Firebase**
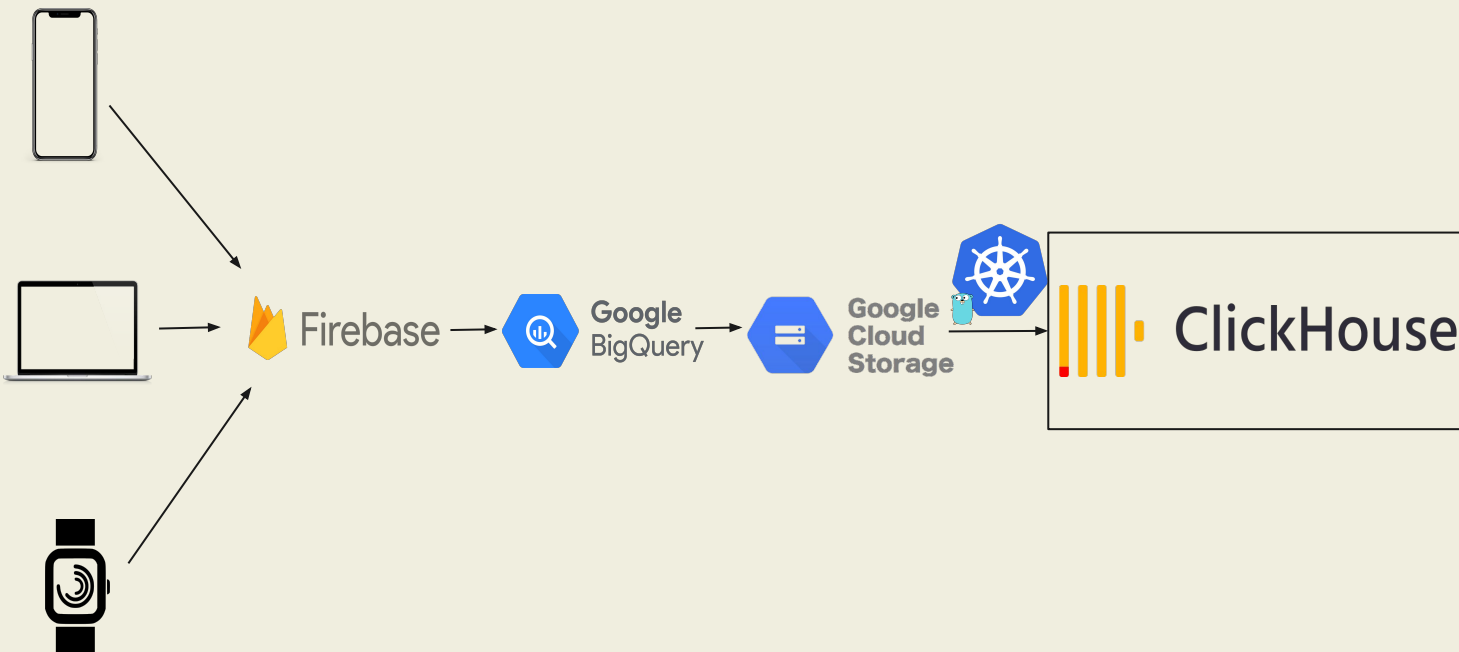
**Google BigQuery**

**Sofascore**

**Google BigQuery**

BigQuery is a completely serverless and cost-effective enterprise data warehouse. It has built-in machine learning and BI that works across clouds, and scales with your data.

- Table named *events_YYYYMMDD* is created each day within BQ dataset
- Each column in that table represents an event-specific parameter
- Apart from event-specific parameters other columns are:
  - Event (event_name, event_date, event_timestamp...)
  - User (user_pseudo_id, user_id, user_first_touch_timestamp)
  - Geographical (continent, country, region...)
  - Device (category, language, operating_system, ...)
  - App info (id, version, install_source, ...)
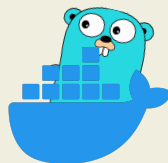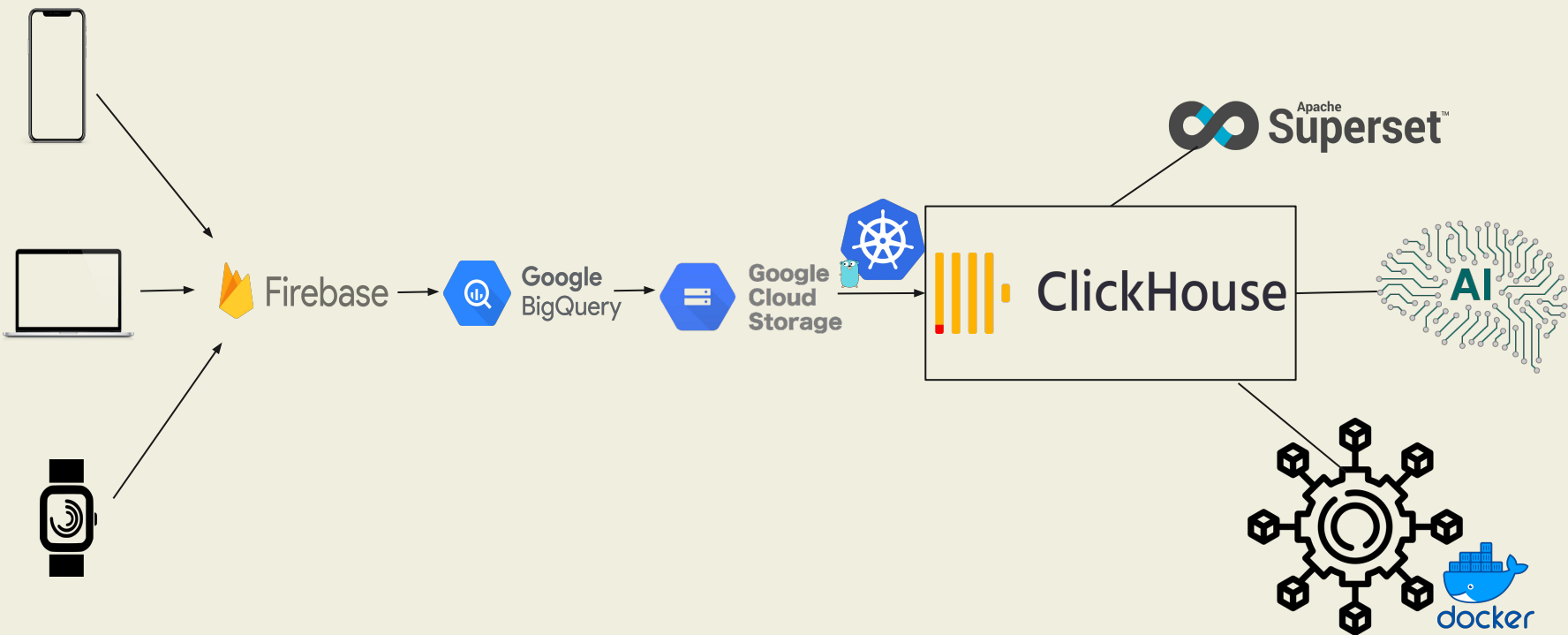  - Traffic Source (name, medium, ...)

**Sofascore**

**sofascore**

- **Kubernetes** (**K8s**) - open-source system for automating deployment, scaling, and management of containerized applications
- Dockerized Go script that imports data into Clickhouse from files exported to GCS from BigQuery
- Parallel import on 15 servers in Sofascore Kubernetes cluster
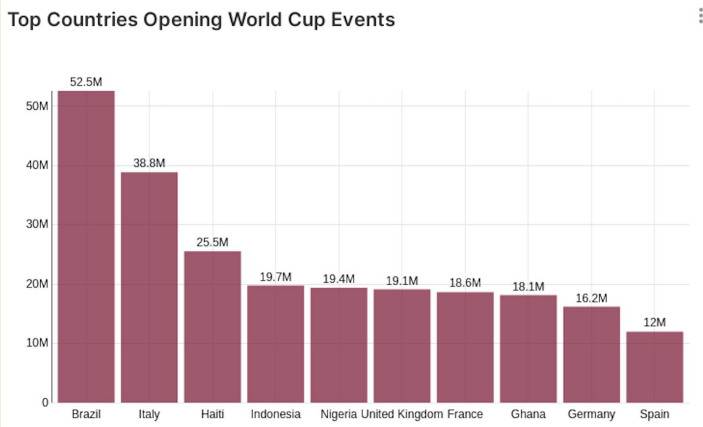- Import using K8s cluster speeds up the import ~15 times 🚀

VELOCITY

**Sofascore**

DATA PIPELINE

Firebase → Google BigQuery → Google Cloud Storage → ClickHouse → AI

Apache Superset

docker

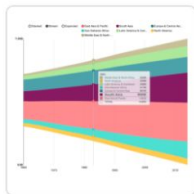sofascore

Apache **Superset**™

- Most popular (by Github ⭐) **open-source** BI and analytics visualization platform in the world
- Software application for data exploration that handles data at petabyte scale
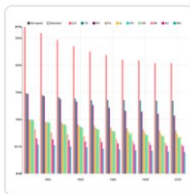- Wide range of database support

**Most Followed Players**

| Player | Nationality | Follows |
|---|---|---|
| Lionel Messi | | 53,855 |
| Cristiano Ronaldo | | 40,067 |
| Neymar | | 26,416 |
| Kylian Mbappé | | 20,752 |
| Robert Lewandowski | | 12,033 |
| Vinícius Júnior | | 11,071 |
| Luka Modrić | | 8,402 |
| Ángel Di María | | 8,030 |
| Kevin De Bruyne | | 7,929 |
| Julián Álvarez | | 7,533 |

**Top Countries Opening World Cup Events**



| Country | Value |
|---|---|
| Brazil | 52.5M |
| Italy | 38.8M |
| Haiti | 25.5M |
| Indonesia | 19.7M |
| Nigeria | 19.4M |
| United Kingdom | 19.1M |
| France | 18.6M |
| Ghana | 18.1M |
| Germany | 16.2M |
| Spain | 12M |

sofascore

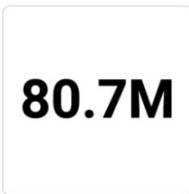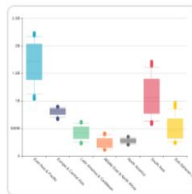# Apache Superset™

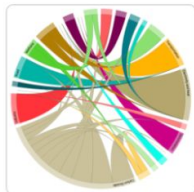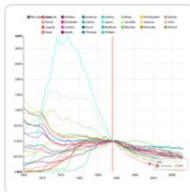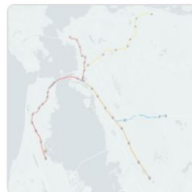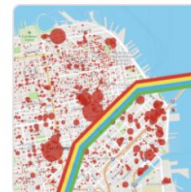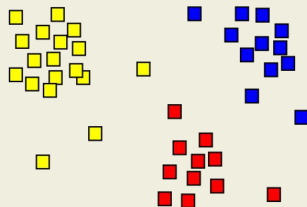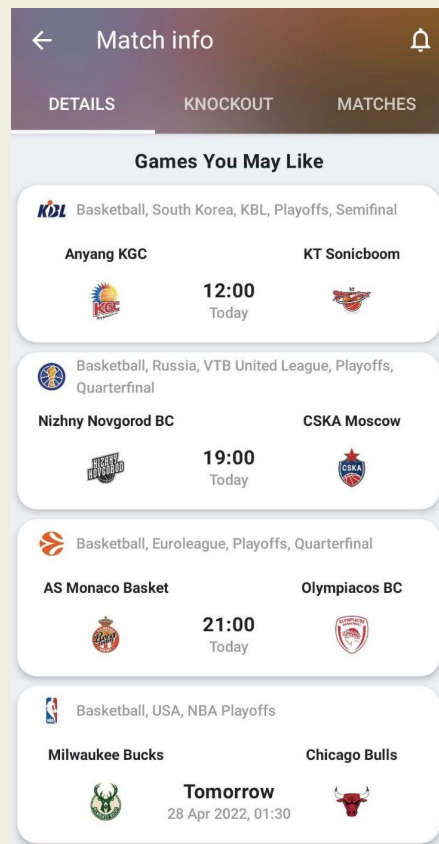| | | |
|---|---|---|
| Area Chart | Time-series Bar Chart | Big Number with Trendline |
| Big Number | Box Plot | Bubble Chart |
| Bullet Chart | Calendar Heatmap | Chord Diagram |
| Time-series Percent Change | Country Map | deck.gl Arc |
| deck.gl Geojson | deck.gl Grid | deck.gl 3D Hexagon |
| deck.gl Multiple Layers | | |

Sofascore

AI

- **User clustering**

- **Predictions (anticipate change in market)**

- **Content recommendation (user retention)**

← Match info 🔔

DETAILS    KNOCKOUT    MATCHES

**Games You May Like**

KBL  Basketball, South Korea, KBL, Playoffs, Semifinal

Anyang KGC                    KT Sonicboom

**12:00**
Today

Basketball, Russia, VTB United League, Playoffs, Quarterfinal

Nizhny Novgorod BC            CSKA Moscow

**19:00**
Today

Basketball, Euroleague, Playoffs, Quarterfinal

AS Monaco Basket             Olympiacos BC

**21:00**
Today

Basketball, USA, NBA Playoffs

Milwaukee Bucks              Chicago Bulls

**Tomorrow**
28 Apr 2022, 01:30

VALUE

Sofascore

- Dockerized microservices written in Go and Python
- Defined as different resources inside the kubernetes cluster (cronjobs, jobs, deployments…)
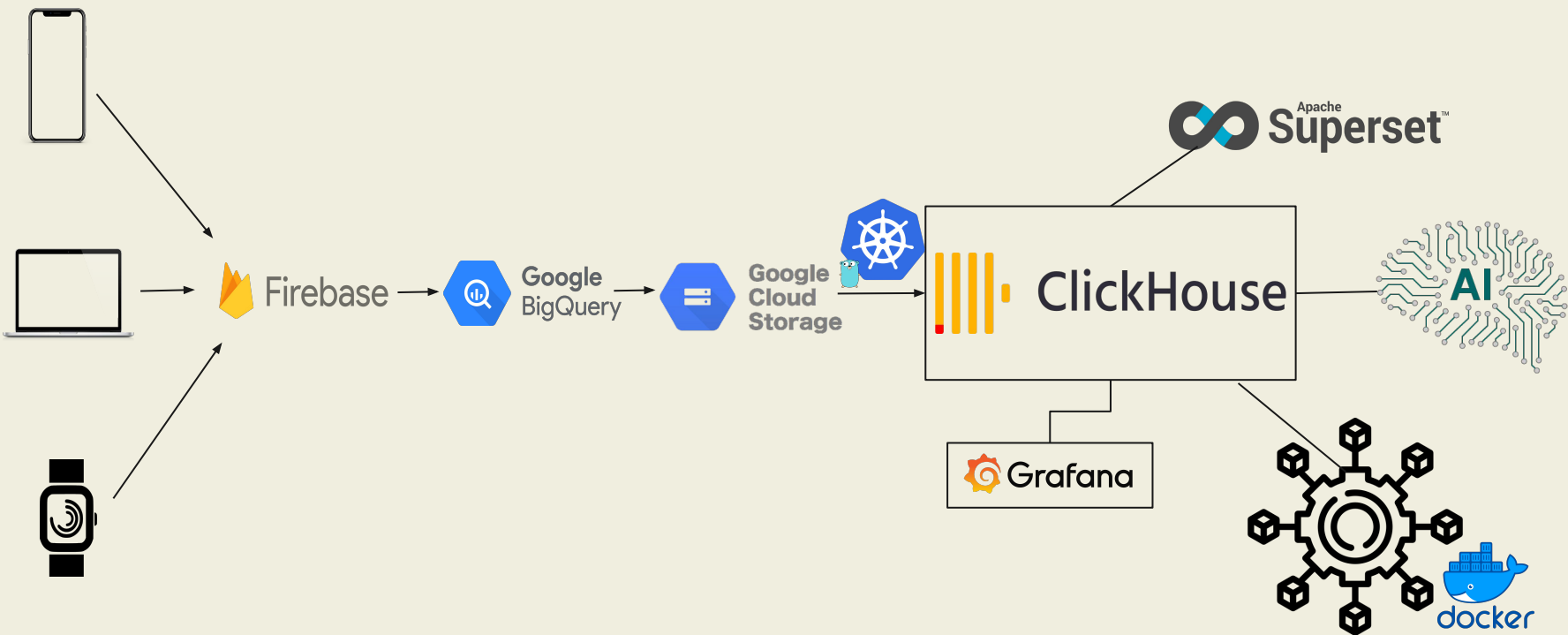- Importing sales and financial data from Google Sheets and Smartsheets
- Importing data from APIs

VARIETY

Sofascore

Firebase → Google BigQuery → Google Cloud Storage → ClickHouse

Apache Superset
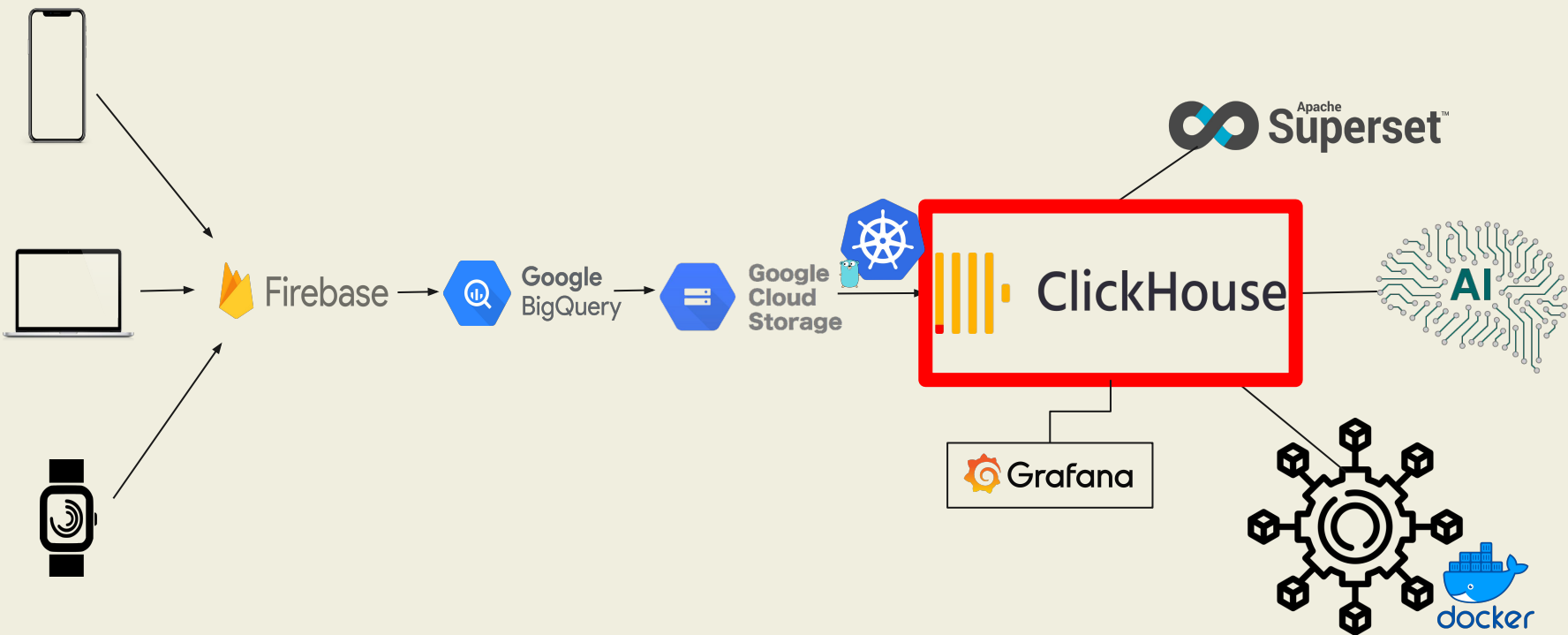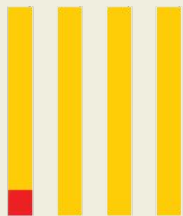
AI

Grafana

docker

sofascore

**Grafana**

- Open source application for multi-platform analytics and interactive visualization
- Allows you to query, visualize, alert on and understand your metrics
- Well suited for time series data visualization
- **Monitoring**:
  - Services
  - Clickhouse performance
  - Server load

VALIDITY



sofascore

sofascore

# ClickHouse

- **Click**stream + Data Ware**house** = **ClickHouse (CH)**
- Open source **column-oriented** SQL database management system (DBMS) for online analytical processing (**OLAP**)
- Fist developed at Yandex and launched in production in 2012 to power **Yandex.Metrica**

**Yandex** Metrica

# Row oriented DBs

- Data associated with a **record** next to each other in memory
- Common row oriented DBs (Postgres, MySQL, MariaDB...)

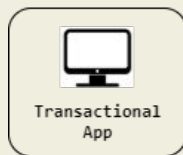| Row | WatchID | JavaEnable | Title | GoodEvent | EventTime |
|-----|---------|------------|-------|-----------|-----------|
| #0 | 89354350662 | 1 | Investor Relations | 1 | 2016-05-18 05:19:20 |
| #1 | 90329509958 | 0 | Contact us | 1 | 2016-05-18 08:10:20 |
| #2 | 89953706054 | 1 | Mission | 1 | 2016-05-18 07:38:00 |
| #N | ... | ... | ... | ... | ... |

# Column oriented DBs

- Data associated with a **field** next to each other in memory
- Common column oriented DBs (Redshift, BigQuery, ClickHouse...)

| Row: | #0 | #1 | #2 | #N |
|------|-----|-----|-----|-----|
| WatchID: | 89354350662 | 90329509958 | 89953706054 | ... |
| JavaEnable: | 1 | 0 | 1 | ... |
| Title: | Investor Relations | Contact us | Mission | ... |
| GoodEvent: | 1 | 1 | 1 | ... |
| EventTime: | 2016-05-18 05:19:20 | 2016-05-18 08:10:20 | 2016-05-18 07:38:00 | ... |

Sofascore

# Row oriented DBs

# Column oriented DBs

## OLTP

- Online **Transaction** Processing
- Most row oriented DBs
- Many users performing varied queries and updates
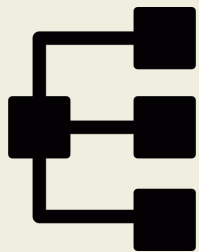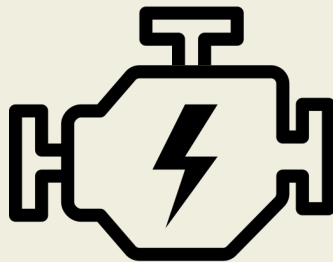- SQL primary language for interaction

## OLAP

- Online **Analytical** Processing
- Most column oriented DBs
- Fewer users performing deep data analysis
- Utilizes particular query language other than SQL

Transactional App

Analytics, Reporting

OLTP

OLAP

Sofascore

# Clickhouse features



Parallel query processing



Multiple table engines


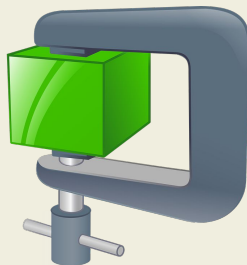
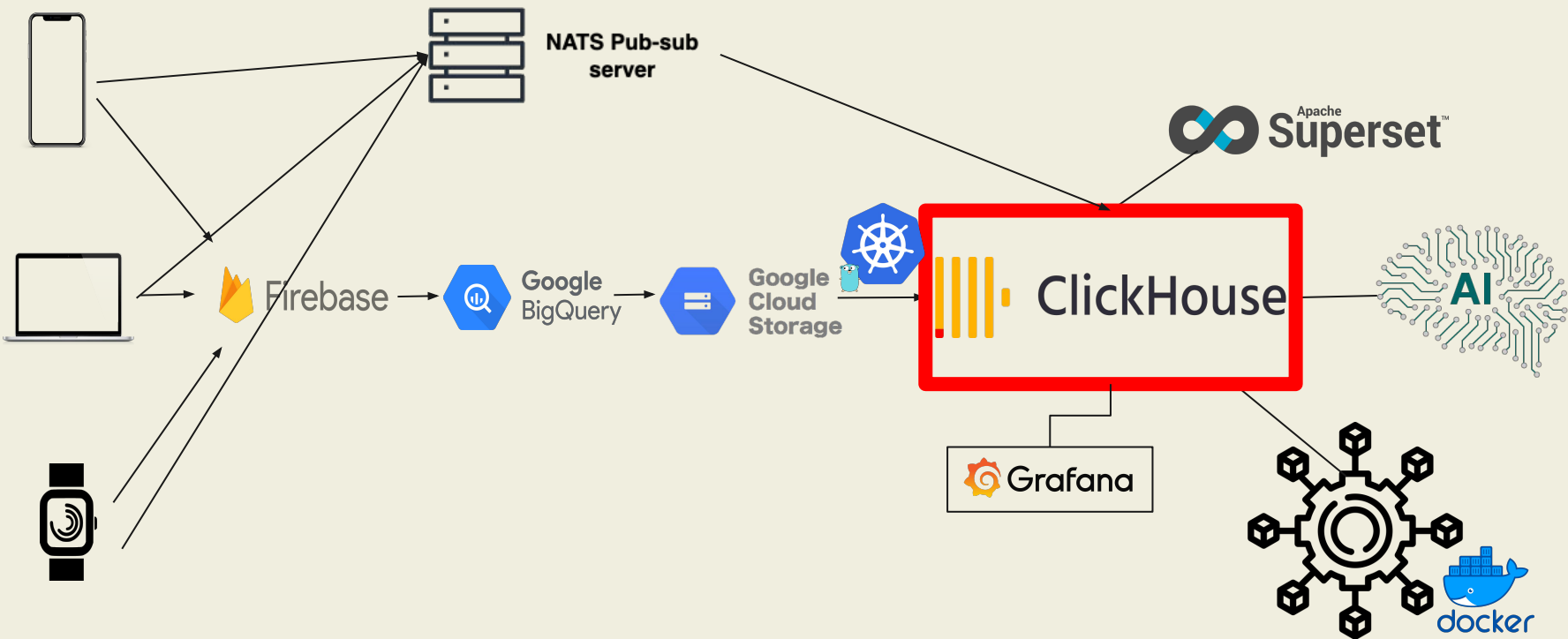OPEN SOURCE

Cost effective

$f(x)$

HUGE number of functions



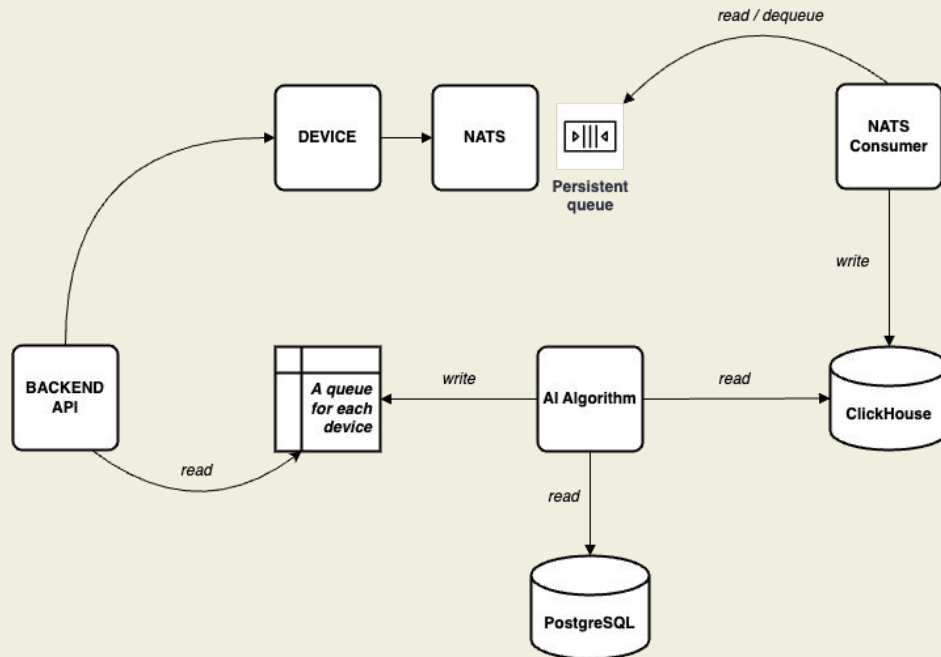Compression



Index support

sofascore

NATS & REAL-TIME ANALYTICS AND AI

NATS Pub-sub server

Apache Superset

Firebase

Google BigQuery

Google Cloud Storage

ClickHouse

AI

Grafana

docker

sofascore

# NATS

Summary

# Summary

- ➔ No. 1 sports platform in the world
- ➔ Big Data: challengeable
- ➔ User action data > sports data
- ➔ Data latency
- ➔ Real-time data pipeline with NATS

**Sofascore**

# Thank you!

karlo.knezevic@sofascore.com