

# Rethinking Data @ SoundCloud

Ana Pereira

Janette Müller-Lehmann

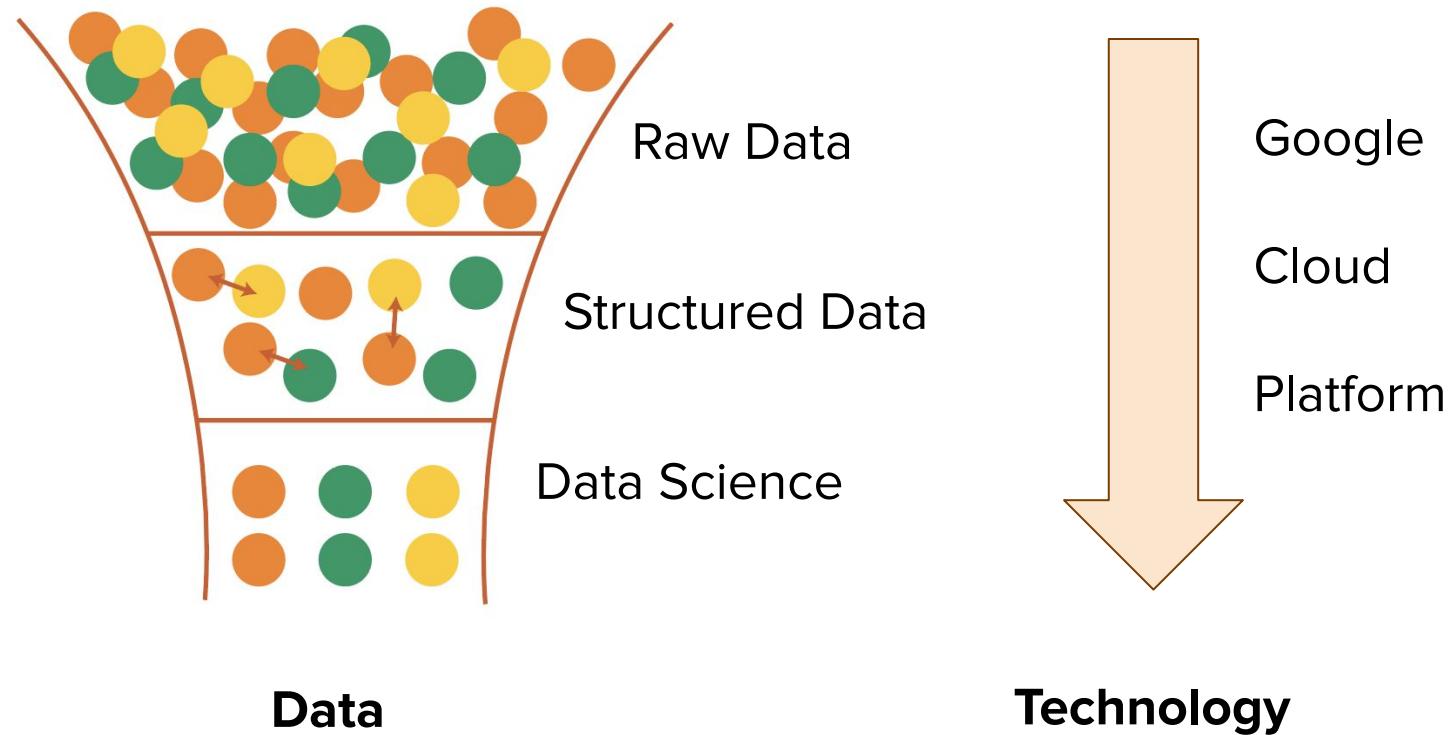
Mahnaz Namazizavareh

Berlin Women in Machine Learning & Data Science Meetup

November 05, 2020



# Agenda



# SoundCloud & Data Science



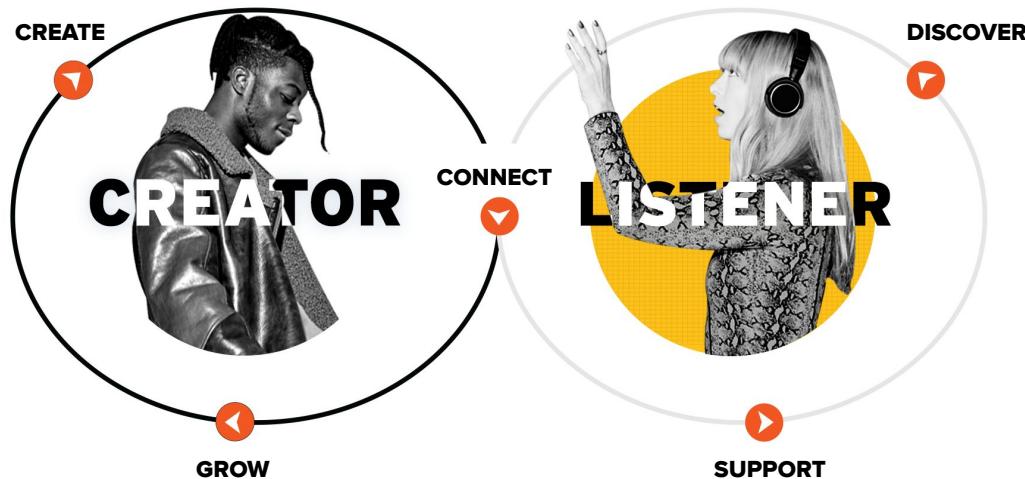
# SOUNDCLOUD

## THE WORLD'S LARGEST OPEN AUDIO PLATFORM

SoundCloud **empowers creators to connect directly** with young, influential listeners driving culture.

A rich audio discovery experience, **combining new music, DJ sets, and remixes direct from creators**, feeding vibrant local scenes worldwide.

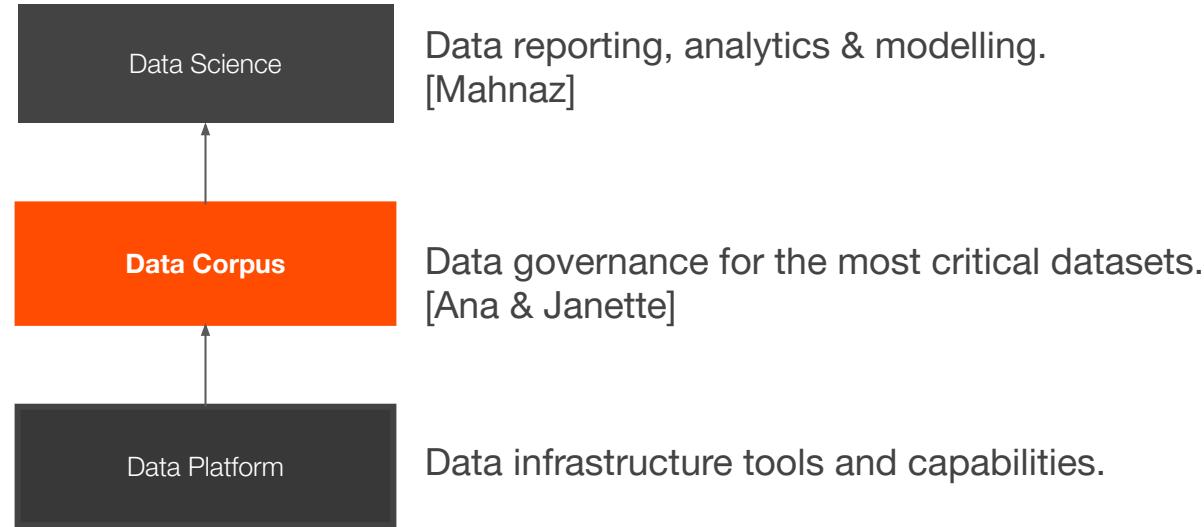
SoundCloud is a global accelerator of audio culture:  
**WHAT'S NEXT IN MUSIC IS FIRST ON SOUNDCLOUD.**



200+ million tracks | 25+ million creators | available in 190 countries and territories



# Data @ SoundCloud



# Data Science @ SoundCloud



## Reporting

How are we performing?

How many plays does my track have (creator)?

...



## Audience Intelligence

How is our influencer community structured?

Who are up-and-coming creators to look out for?

...



## Product Exploration

How is our new feature performing?

Where should we launch next?

...



## Data Products

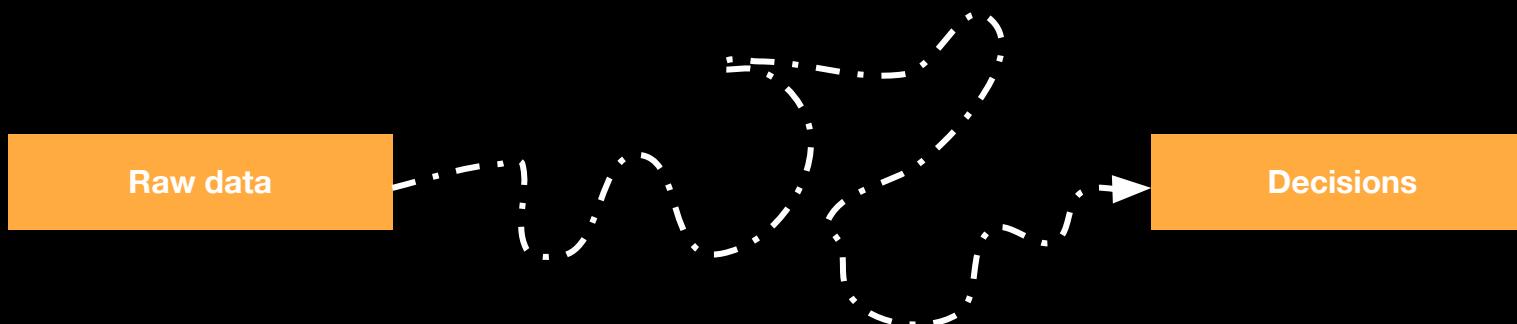
Which are the most interesting tracks for a user?

How do I find relevant hip hop tracks?

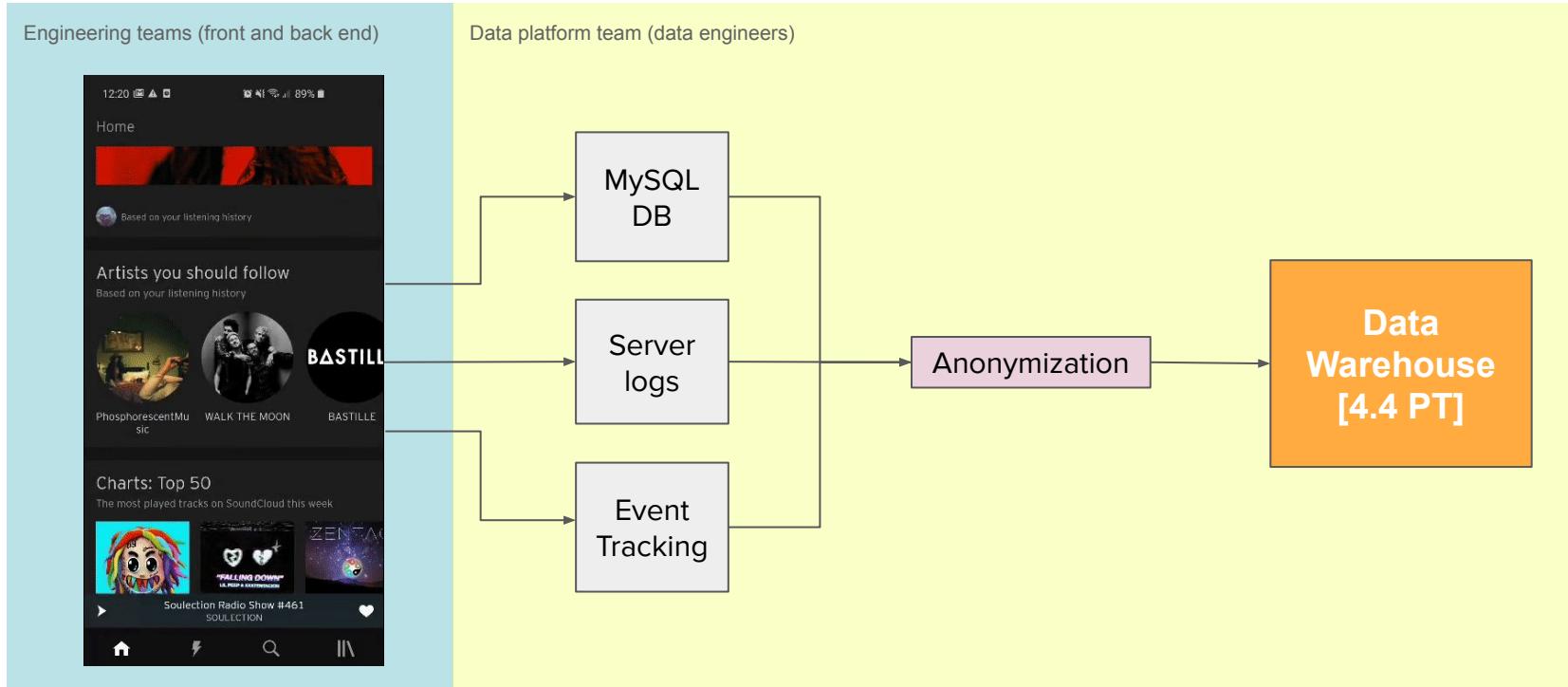
...



# The Journey of Data



# Raw Data



Users interact  
with SoundCloud

Ingestion

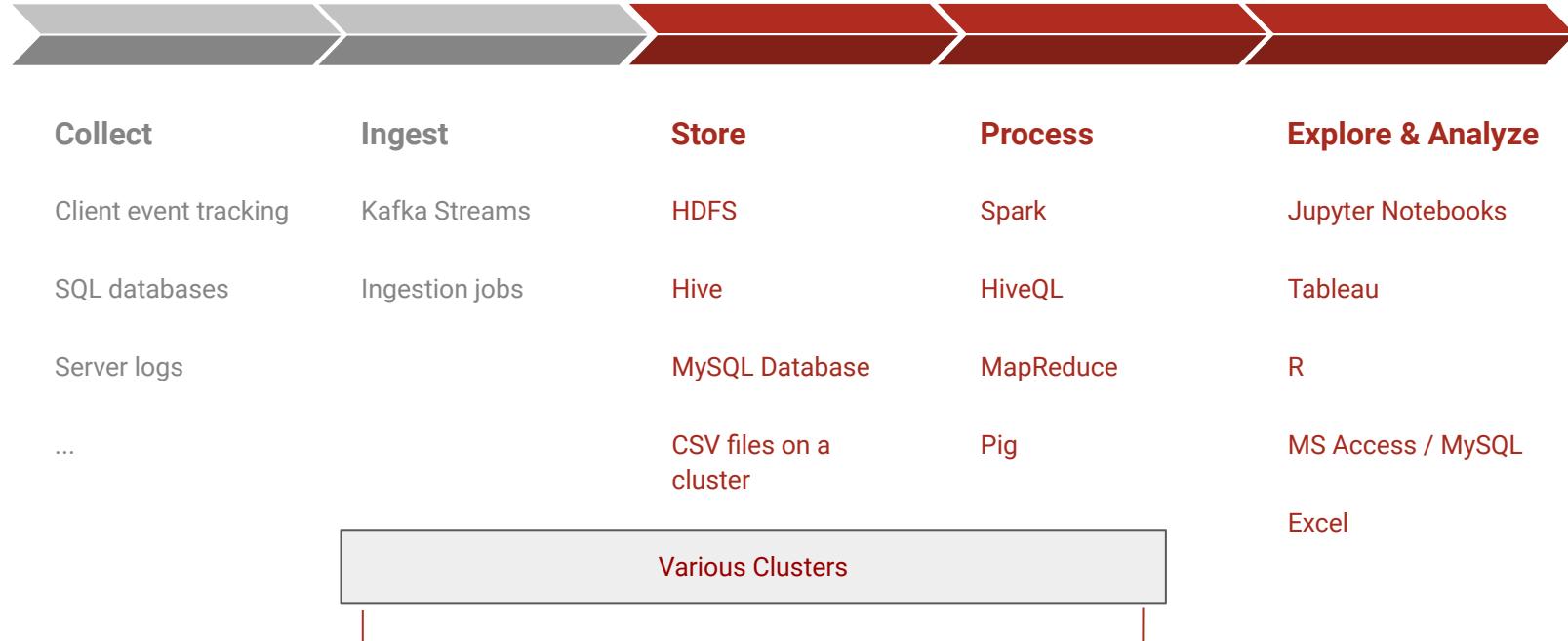
Storage



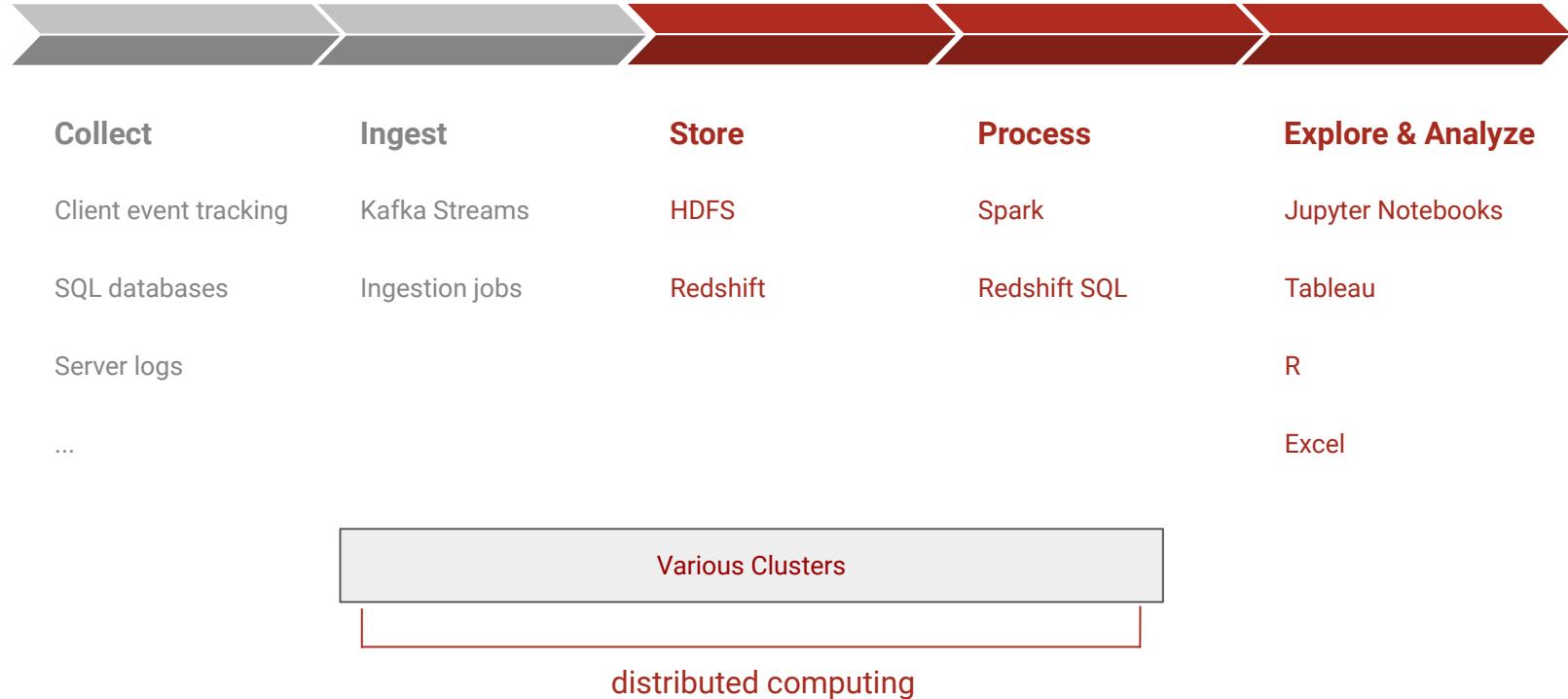
# **Problem 1**

## **The Infrastructure**

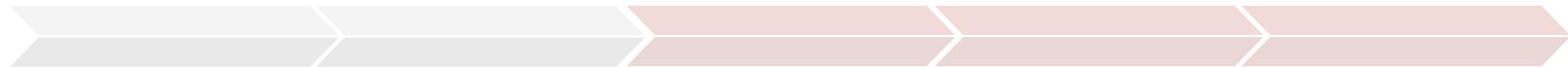
# Data Science in the Past - In General



# Data Science in the Past @ SoundCloud



# Data Science in the Past



Collect

Client event tracking

SQL databases

Server logs

...

High **configuration and maintenance costs** of various clusters, systems and applications.

Batch processing frameworks are **not suitable for ad hoc and trial-and-error data analysis**.

A lot of **copy & paste** between different systems and services.

Not an option for people **without coding experiences**.

Explore & Analyze

Jupyter Notebooks

Tableau

R

MS Access / MySQL

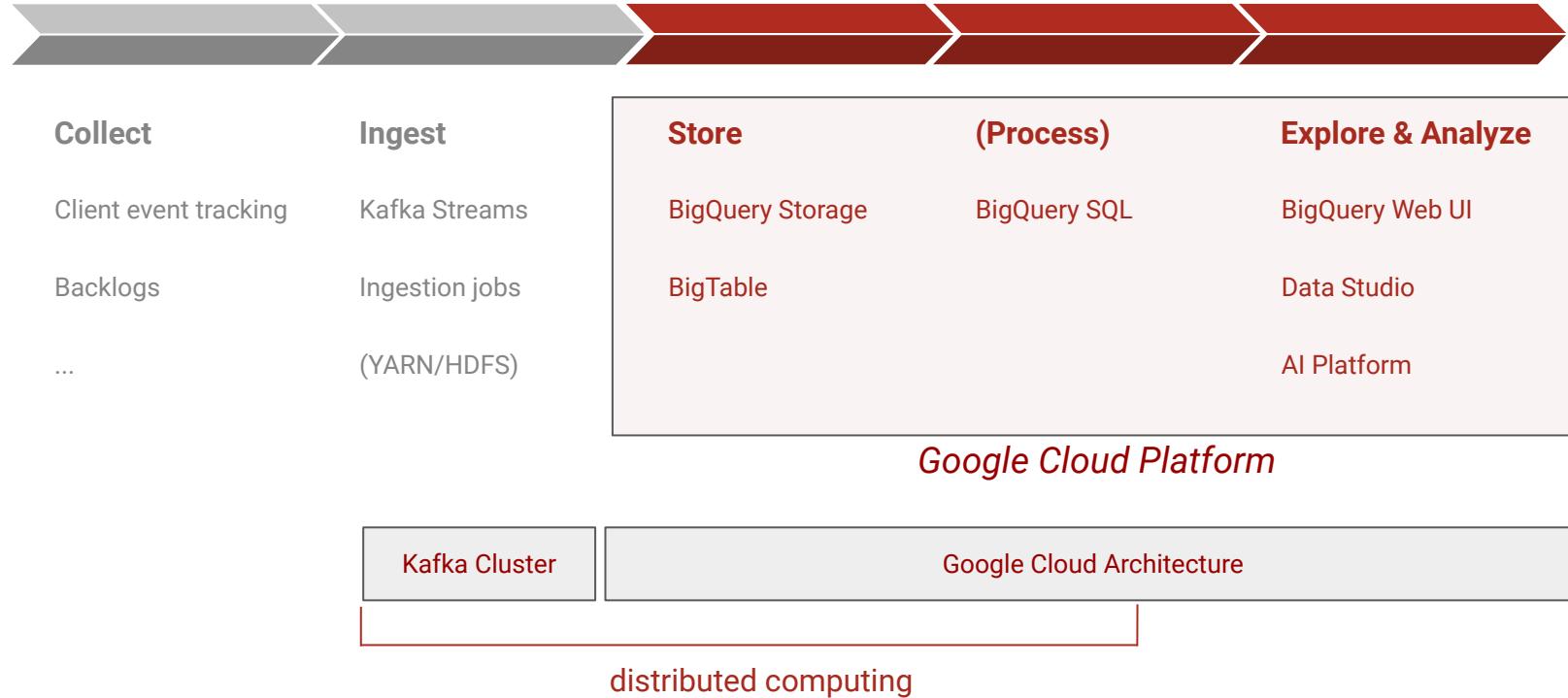
Excel

Various Clusters

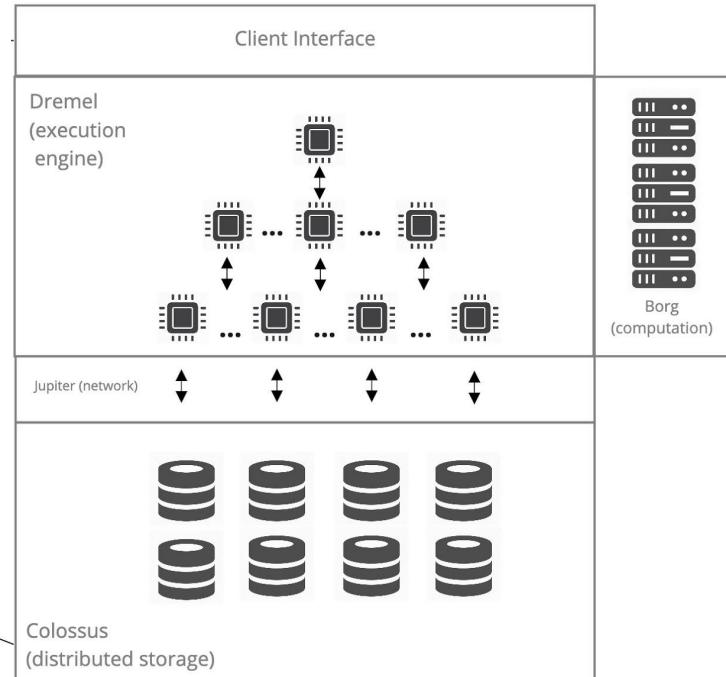
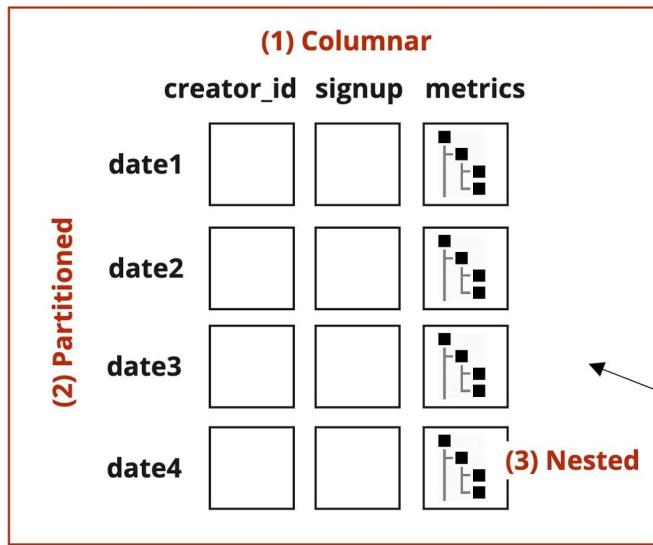
distributed computing



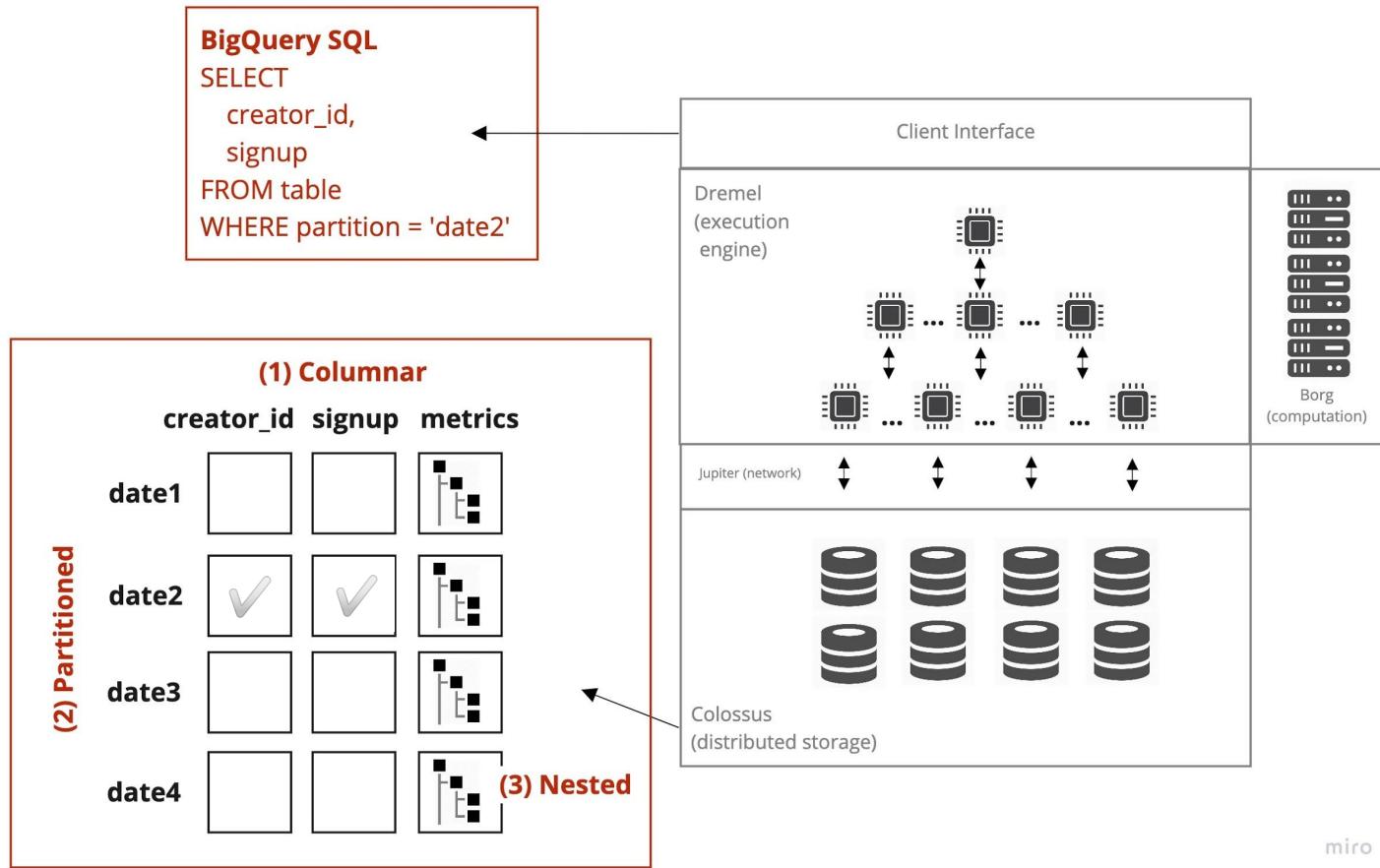
# Data Science in the Cloud @ SoundCloud



# BigQuery Architecture



# BigQuery Architecture



## **Problem 2**

## **The Data Itself**

# The Lake of Raw Data

Raw data has **different schemas**, analysing the data is challenging and error-prone.

Raw data from different systems can be **inconsistent**.

Raw data is **complex** and has a high level of detail.

Raw data is large, even with distributed computing it **takes some time** to get results.

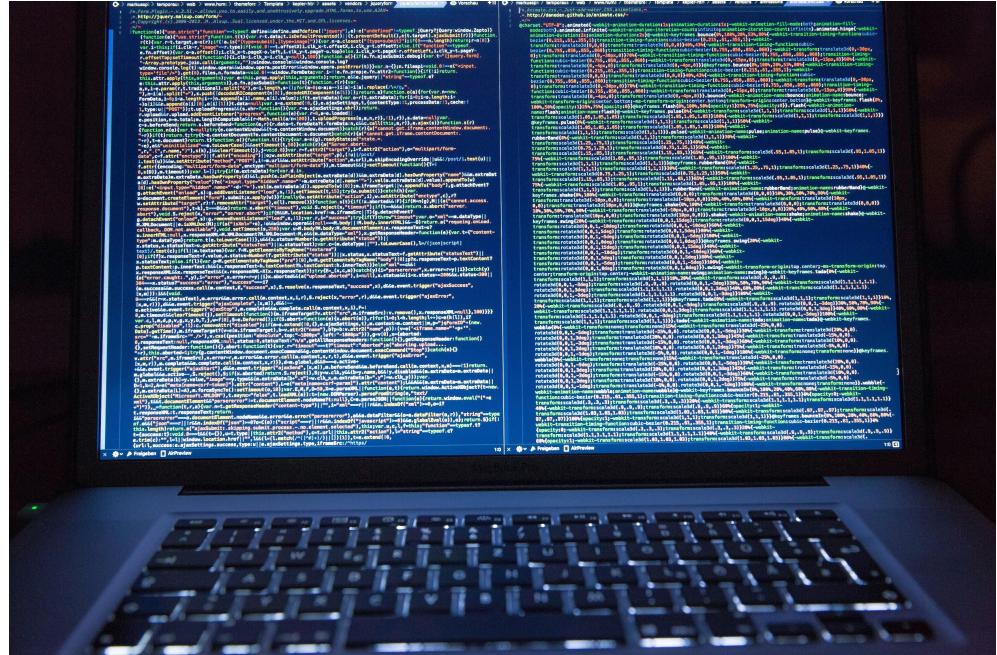


Photo by [Markus Spiske](#) from Pexels

Raw data

# Data Corpus Team

Decisions



# Mission Statement

Implement **data governance** on the datasets that are **most critical** for product and business decision making.

[Definition of data governance](#)



# What characterizes high data quality?

1. Integrity
2. Compliance
3. Timeliness
4. Usability
5. Efficiency
6. Maintainability



# What characterizes high data quality?

1. **Integrity**: Providing an unambiguous and validated source of truth.

- ✓ Manual validations (when building something new)
- ✓ Data quality checks that run daily
- ✓ Anomaly detection
- ✓ Unit testing



# What characterizes high data quality?

2. **Compliance:** Ensuring compliance of data with GDPR.



# What characterizes high data quality?

3. **Timeliness:** Daily updates are ready by 10am CET.

- ✓ Efficiency of ETLs
- ✓ Timeliness of dependencies

# What characterizes high data quality?

4. **Usability:** Providing easily accessible and consumable data.

- ✓ BQ table and field descriptions
- ✓ Design documents
- ✓ Ease of use of resources
- ✓ Providing Training

# What characterizes high data quality?

5. **Efficiency**: Optimize execution time and storage regarding performance and costs.

- ✓ Storage
- ✓ Execution time
- ✓ Make the most of BQ features (sketches, nesting, UDFs)

# What characterizes high data quality?

6. **Maintainability**: The code should be easy to read, extend and maintain.

- ✓ 80/20 principle
- ✓ Best practices
- ✓ Scalability

# Our tech stack

- BigQuery - data warehouse
  - BigQuery SQL - querying language
  - Airflow - scheduler
  - terraform - resource management
  - github - version control
- 
- Ready to use out of the box!

# Learning to work with BQ (some examples)

- The storage location of the data is relevant- we use EU, not US
- Make sure to partition your data, if possible
- SELECT only the columns you need
- Nested and repeated fields reduce storage and computational time
- UDFs (User Defined Functions) allow you to abstract away common coding patterns
- Wide tables perform better

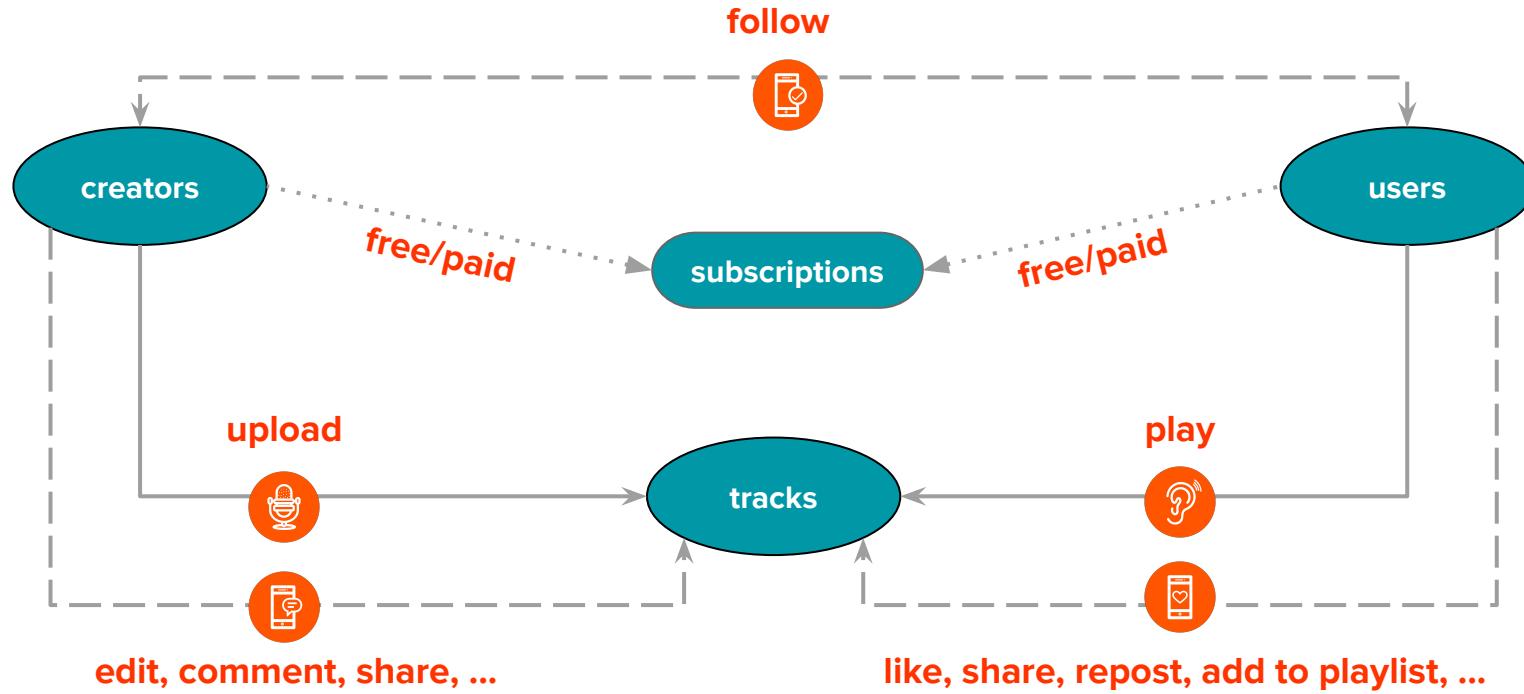
**Data Corpus**

size: ~ 456 TB

← 90% less storage than raw data!



# Data Modelling - Our Main Entities



# daily.creator\_metrics table

Each entity has a table with its main metrics.

partition column

nested and repeated fields

date	creator_id	has_record_deal	actions_received. is_logged_in	actions_received. metric_name	actions_received. metric_value
2020-05-20	soundcloud:creator:2	TRUE	FALSE	n_plays	10
			TRUE	n_plays	484
			TRUE	n_shares	1
			FALSE	listening_time_ms	75829
2020-05-20	soundcloud:creator:3	FALSE	TRUE	n_plays	199
			TRUE	listening_time_ms	29733
2020-05-20	soundcloud:creator:4	FALSE	FALSE	n_plays	9
			FALSE	listening_time_ms	28485

one row per creator (on each day, if they received an action)



Raw data

# Data Corpus Use Cases

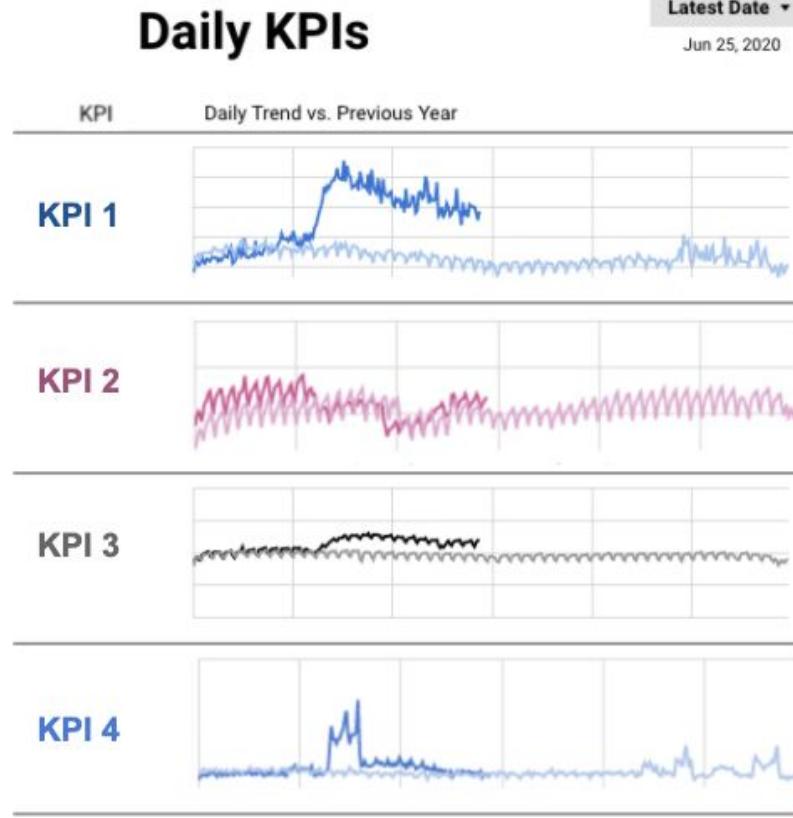
Reporting

Self-serve

Analysis & Modelling



# Reporting → Data Studio



- High level reporting
- Everyone in the company can view
- Connects directly to BQ
- KPI- Key Performance Indicator



# Self-Serve → BigQuery Web UI

What were the **top 10 creators** for the first half of 2020 (in terms of number of plays)?

The screenshot shows the BigQuery Web UI interface. At the top, there is a navigation bar with a dropdown menu "Corpus - Centralized SC Data", a search bar "Search products and resources", and a user profile icon. Below the navigation bar, the main area has tabs for "Unsaved query" (which is currently selected) and "Edited". There are buttons for "+ COMPOSE NEW QUERY", "HIDE EDITOR", and "FULL SCREEN". The main content area displays a SQL query:

```
1 SELECT
2   creator_id,
3   SUM(metric_value) AS n_plays
4 FROM
5   `sc-corpus.daily.creator_metrics`, UNNEST(actions_received)
6 WHERE
7   date BETWEEN '2020-01-01' AND '2020-06-30'
8   AND metric_name = 'n_plays'
9 GROUP BY creator_id
10 ORDER BY n_plays DESC
11 LIMIT 10
```

At the bottom, there are several action buttons: "Run", "Save query", "Save view", "Schedule query", and "More". A "SOUNDCLOUD" logo is visible on the far left.

# Analysis & Modelling → AI Platform Notebooks



- Managed [JupyterLab](#) notebook instances
- For the more advanced users
- You can use whatever environment you'd like
- You can also define how powerful your machine should be machine
- Connects seamlessly to BQ

# Want to join?

- For this and more job opportunities, check out our [jobs page](#).
- **Engineering trainee program:** DevBridge ([Blog post](#) about current edition).

# Questions?

