# Introduction

Exploring the Prezi Data Platform, its architecture as a Data Lakehouse, and the key components.

- What is a data platform and why we need it
- Different Architectures and DataLakehouse
- Components

Data Engineer    @Prezi

Data Infrastructure Team

From Italy🇮🇹
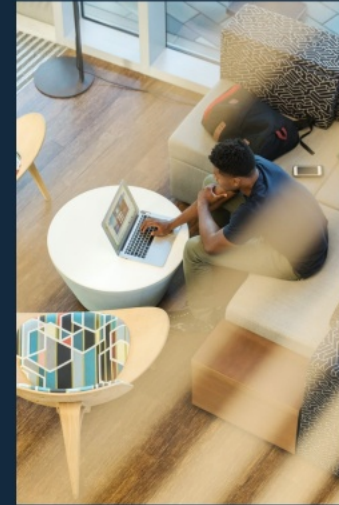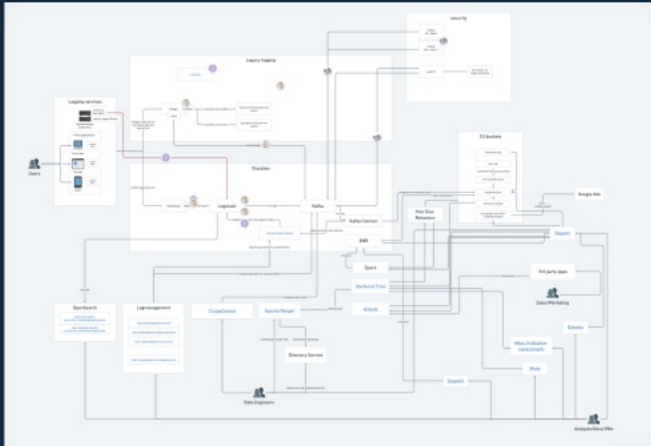
Vincenzo Cassaro

Data Engineer        @Prezi

Data Infrastructure Team

From Italy🇮🇹

# Data Platform

Define, Ingest, Store, Process, Orchestrate
User Events

# Architectures

## Data Warehouse
Table storage + Compute


## Data Lake
Multiple format (messy) storage


## Data LakeHouse
Best of Both World (free file format + Table Api)

# Architectures

Data Warehouse
Table storage + Compute

Data Lake
Multiple format (messy) storage

Data LakeHouse
Best of Both World (free file format + Table Api)

# 00



## Definition

Events Definition

- Logmanagement (Django)
- New schema in Schema Registry
- New topic in Kafka
- Logging Libraries (all languages)

# 00

## Definition

Events Definition

- Logmanagement (Django)
- New schema in Schema Registry
- New topic in Kafka
- Logging Libraries (all languages)

# 01



## Ingestion

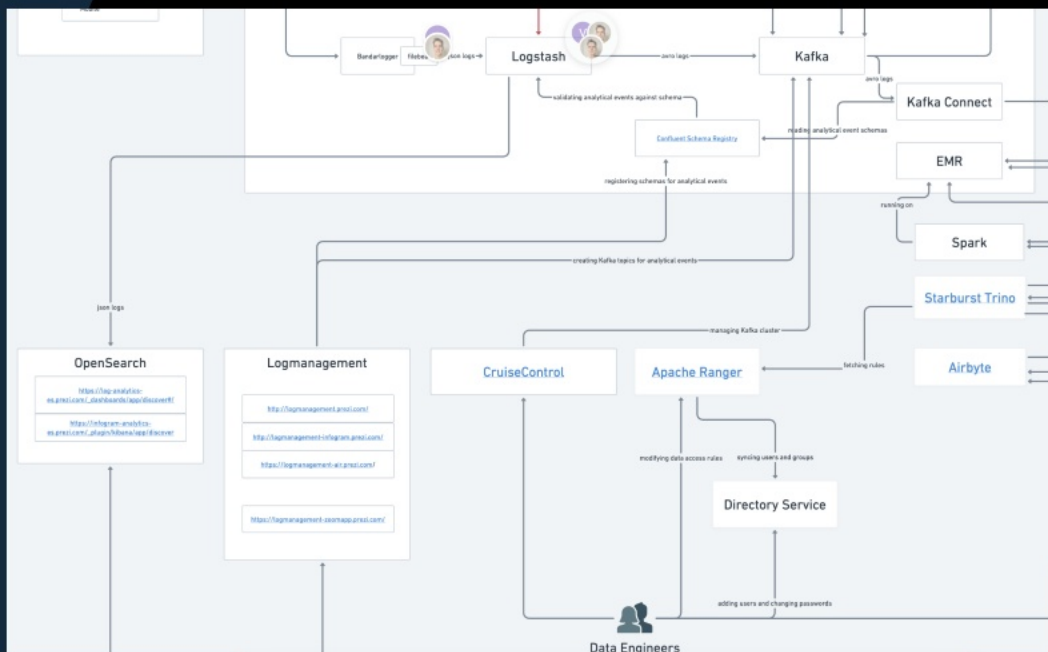Efficient, real-time data ingestion ensures seamless data flow into the Data Lakehouse

• Bandarlogger + Filebeat
• Logstash
• Kafka
• Kafka Connect

# 01



## Ingestion

Efficient, real-time data ingestion ensures seamless data flow into the Data Lakehouse

- Bandarlogger + Filebeat
- Logstash
- Kafka
- Kafka Connect

"On Premise" -> MSK, Gobblin

# 02

## Storage

The robust storage of the Data Lakehouse accommodates structured, unstructured, and semi-structured data, enabling scalable and efficient data management.

S3 + Parquet + Hive

# 02



## Storage

The robust storage of the Data Lakehouse accommodates structured, unstructured, and semi-structured data, enabling scalable and efficient data management.

S3 + Parquet + Hive

Hive Metastore -> Glue

# 02



## Storage

The robust storage of the Data Lakehouse accommodates structured, unstructured, and semi-structured data, enabling scalable and efficient data management.
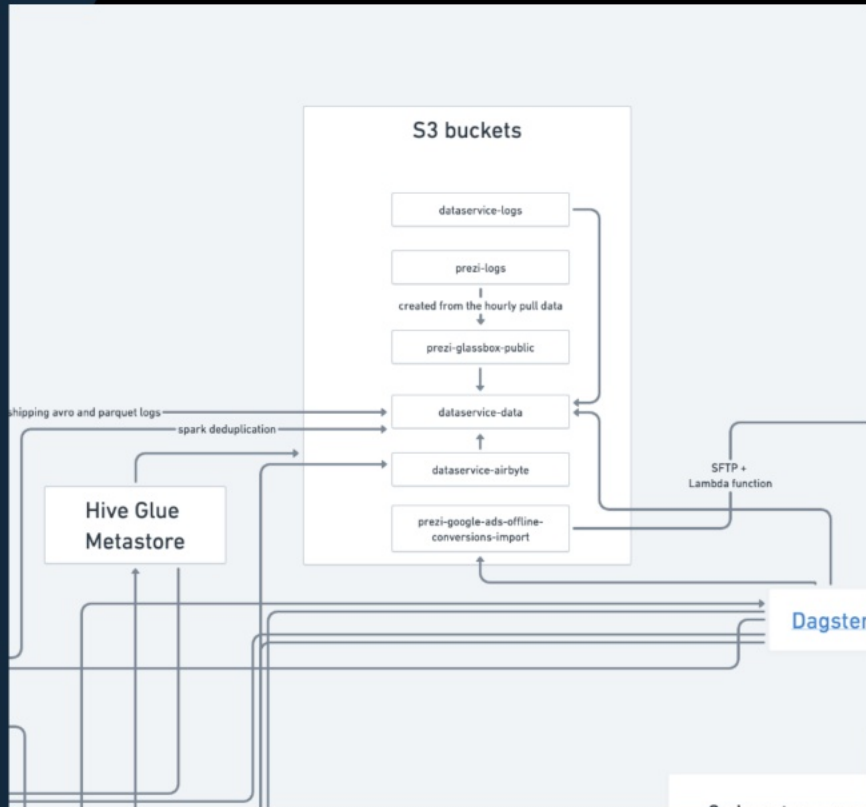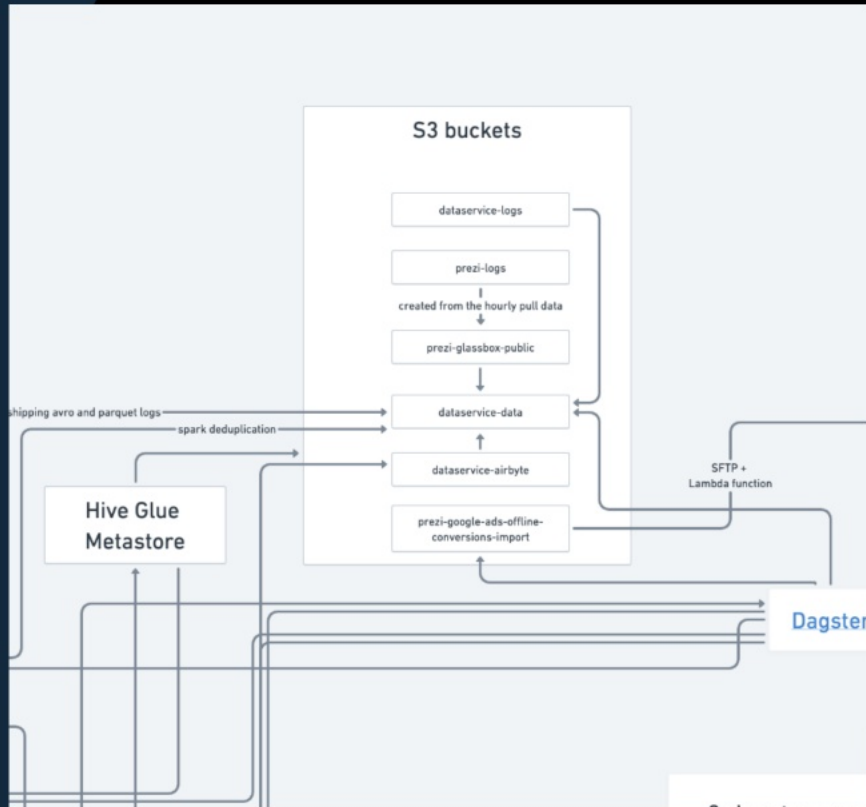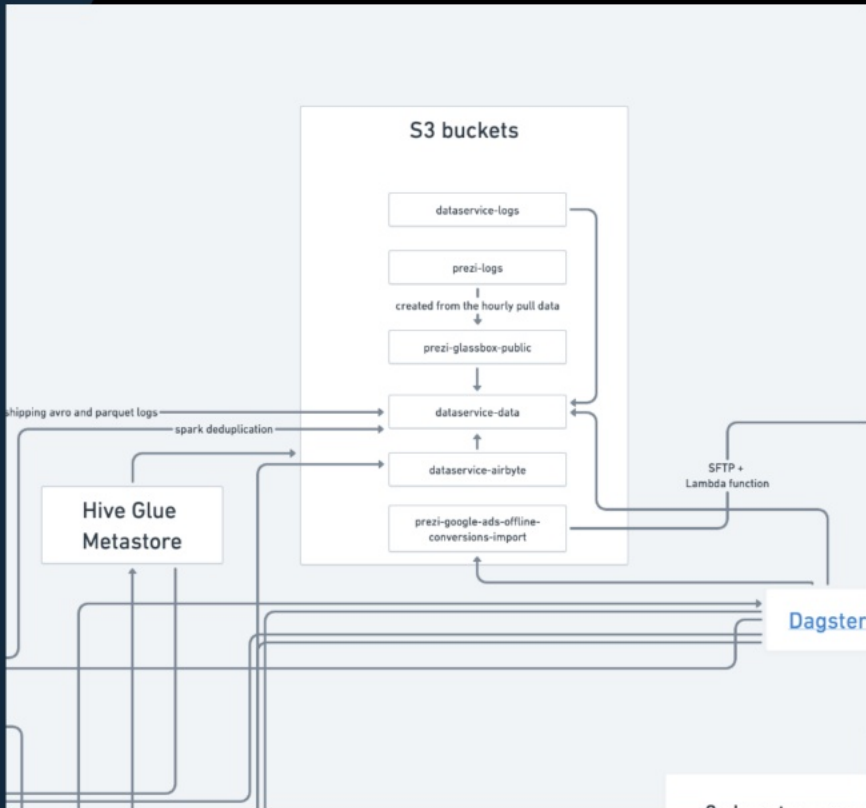
S3 + Parquet + Hive

Hive Metastore -> Glue

Iceberg

# 03



## Processing

I the Data Lakehouse Architecture, processing is offloaded to a separate component. In our infra we use Spark and Trino

7 x r5.2xlarge

# 03



## Processing

I the Data Lakehouse Architecture, processing is offloaded to a separate component. In our infra we use Spark and Trino

7 x r5.2xlarge

Spark 3.5
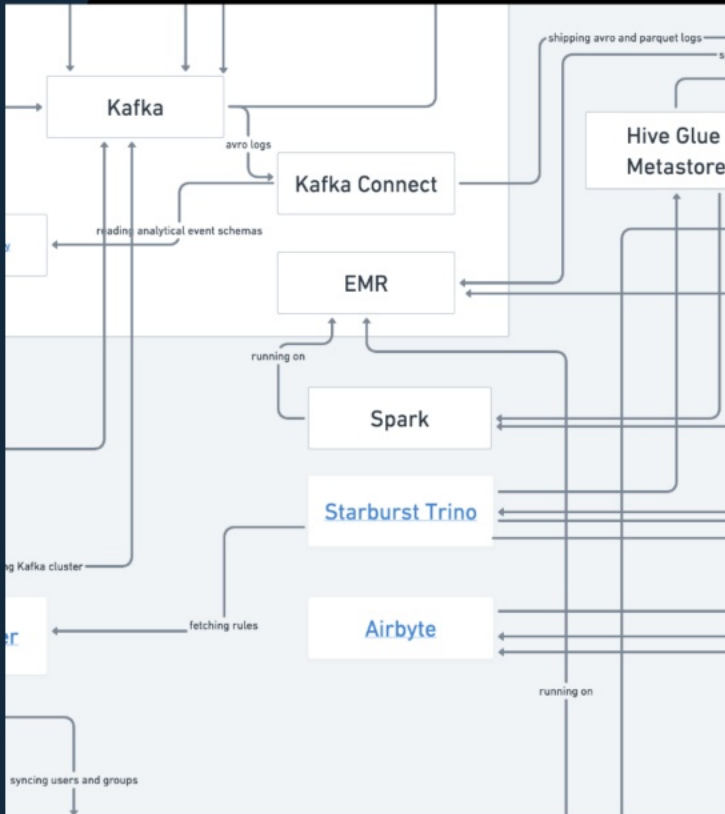
# 03



## Processing

I the Data Lakehouse Architecture, processing is offloaded to a separate component. In our infra we use Spark and Trino

7 x r5.2xlarge

Spark 3.5

Scalability and saas

# 03



## Processing

I the Data Lakehouse Architecture, processing is offloaded to a separate component. In our infra we use Spark and Trino

7 x r5.2xlarge

Spark 3.5

Scalability and saas

HIVE

~25 RDS
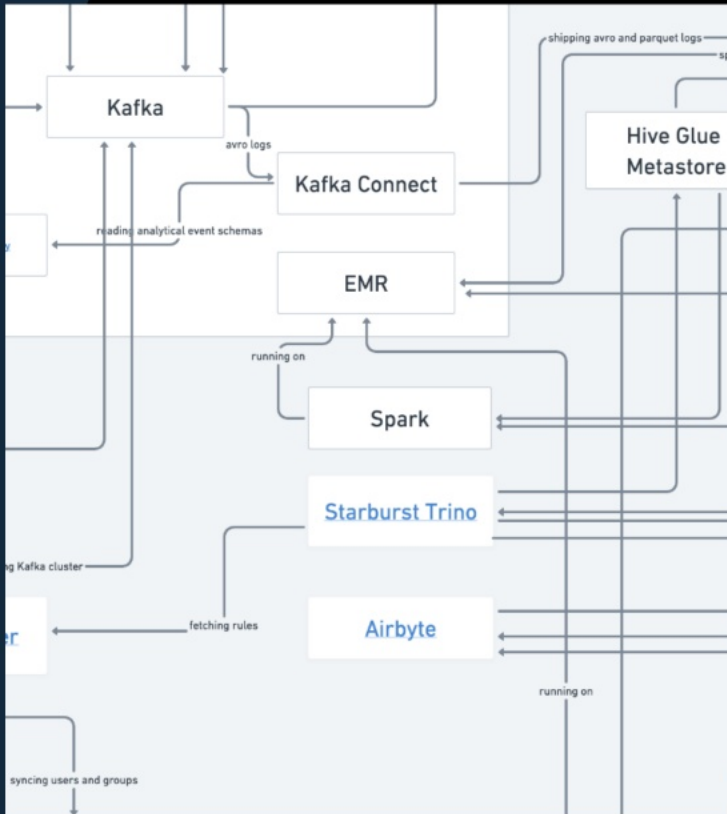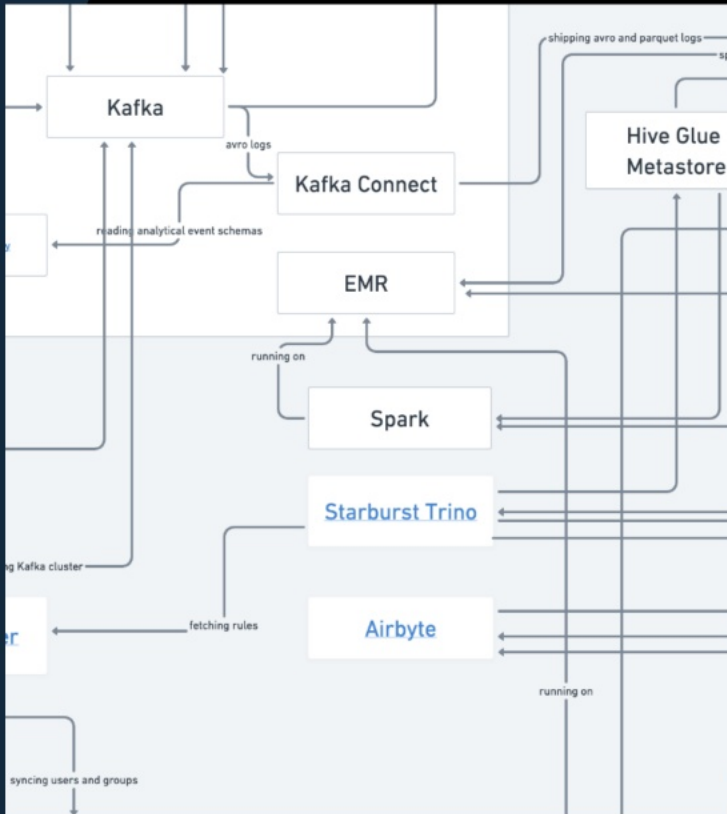
trino

DG
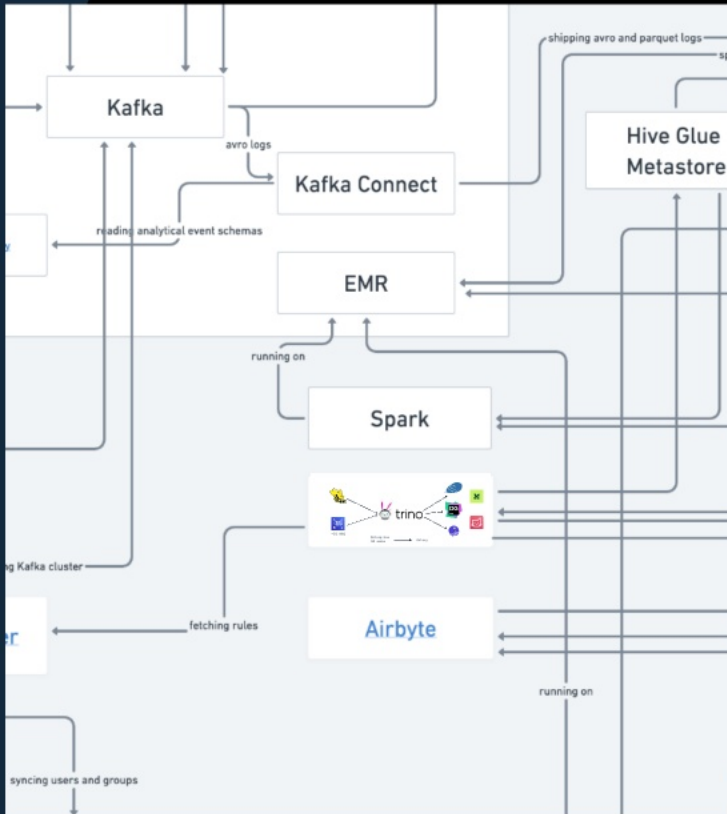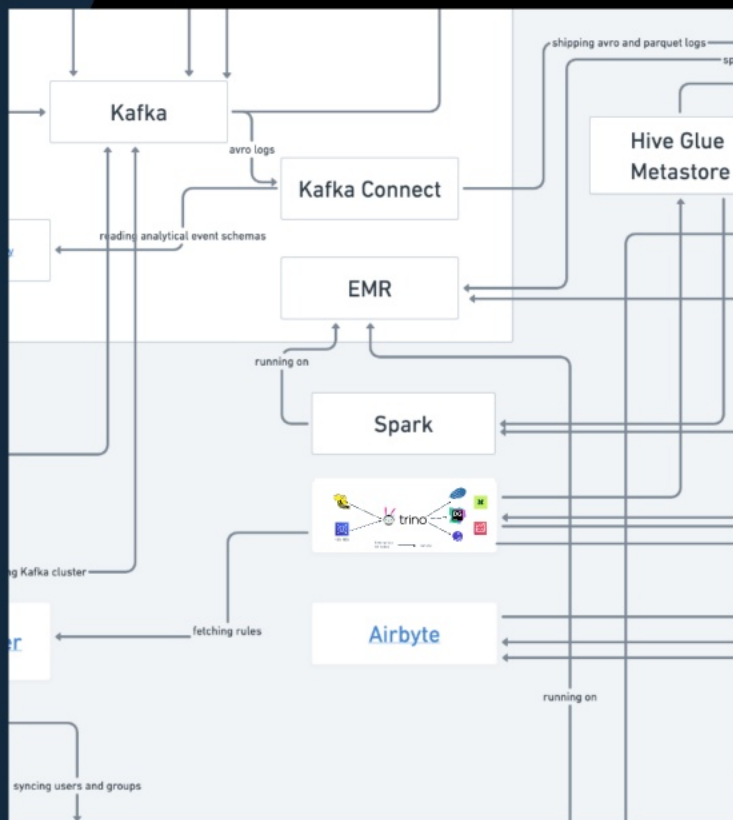
M

Enterprise
10 nodes ⟶ Galaxy

# 03



## Processing

I the Data Lakehouse Architecture, processing is offloaded to a separate component. In our infra we use Spark and Trino
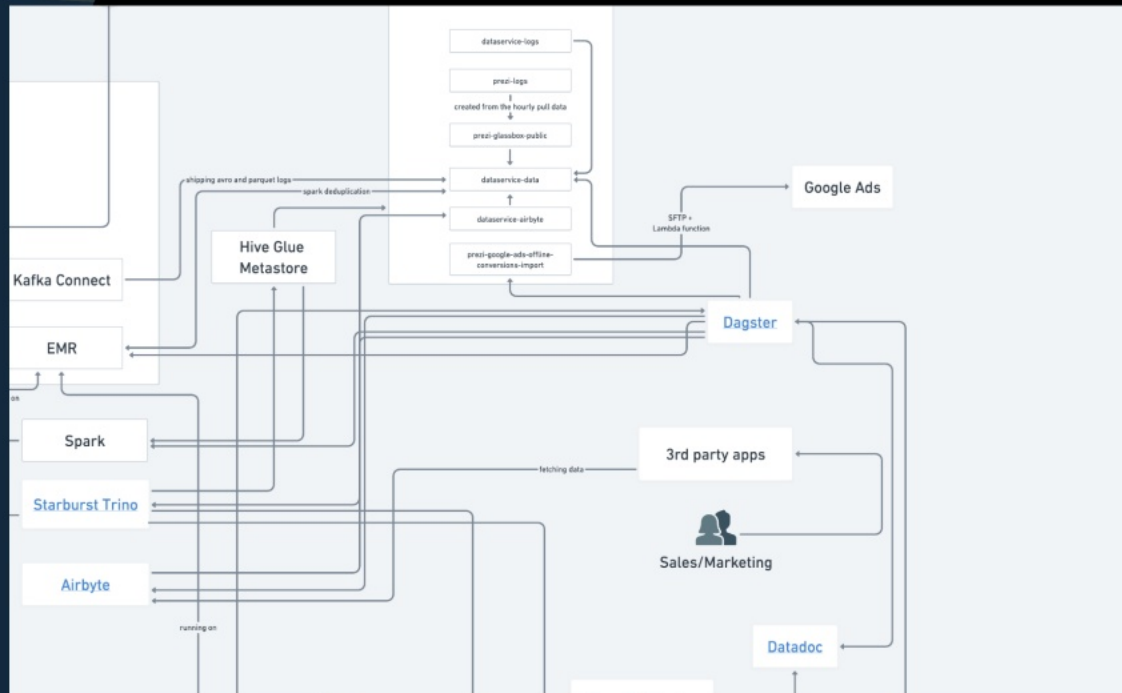
7 x r5.2xlarge

Spark 3.5

Scalability and saas

# 04

## Orchestration

The numerous jobs we have, are coordinated/orchestrated by a component called Dagster. It's responsibility is to run the correct job at the correct time

# 04

## Orchestration

The numerous jobs we have, are coordinated/orchestrated by a component called Dagster. It's responsibility is to run the correct job at the correct time

# 04

## Orchestration

The numerous jobs we have, are coordinated/orchestrated by a component called Dagster. It's responsibility is to run the correct job at the correct time
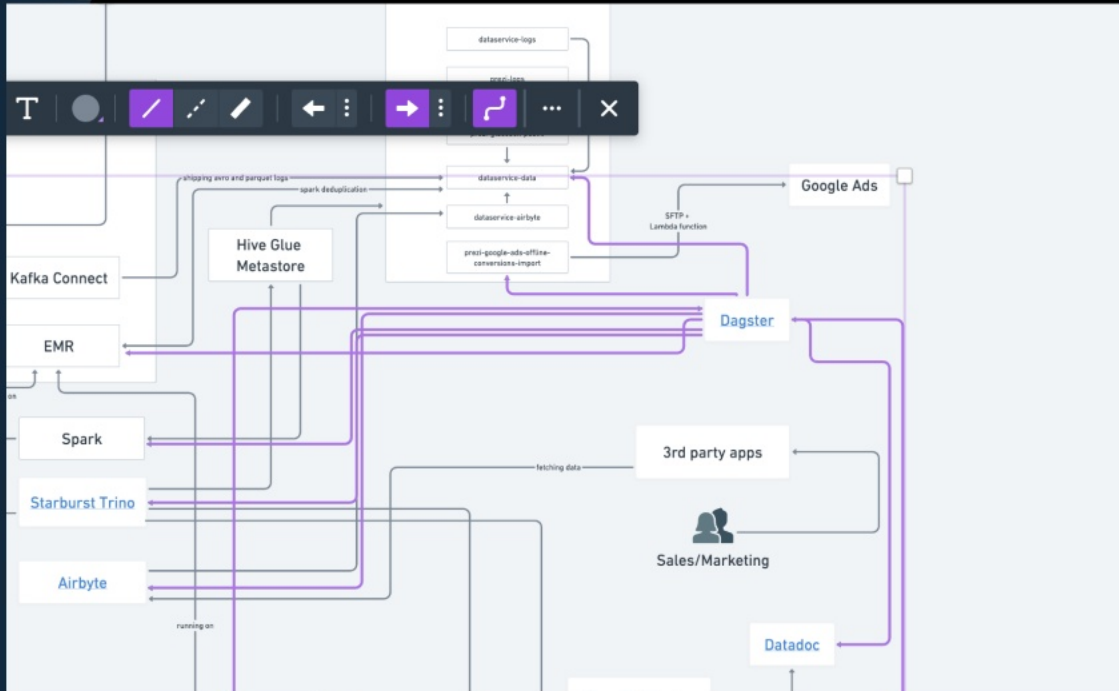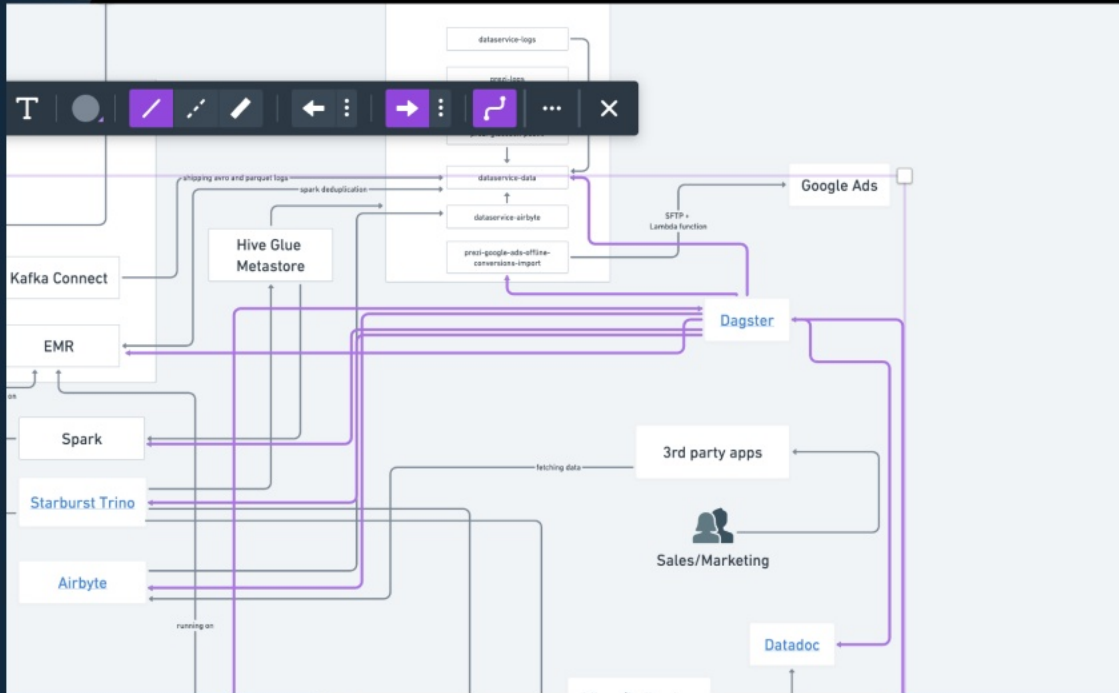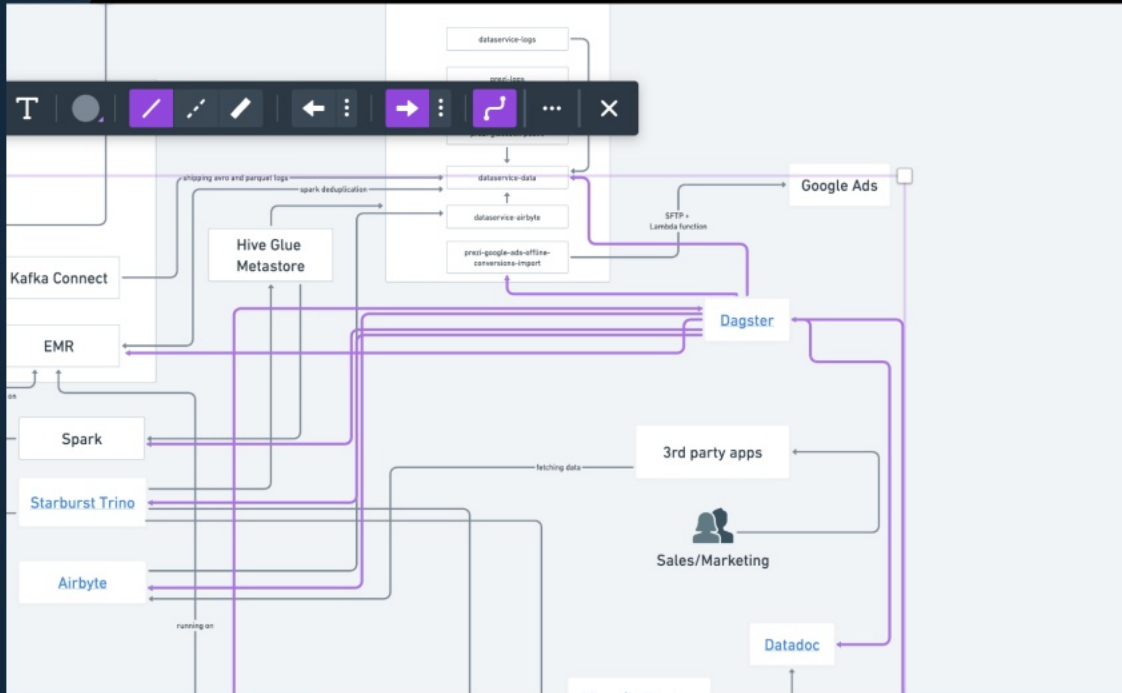
**More Flowkeeperless**

# 04



## Orchestration

The numerous jobs we have, are coordinated/orchestrated by a component called Dagster. It's responsibility is to run the correct job at the correct time

**More Flowkeeperless**

**Hybrid Saas**

# Conclusion

To sum it up:
- Prezi Data Platform is founded on stable/state-of-the-art techs and architecture decision
- It allows the introduction of cutting edge technology, as it is pluggable and modular
- It's constantly evolving

Trino Summit 2024

Thanks!

@vincenzocassaro

# State of
# Prezi Data Platform

## Introduction

Exploring the Prezi Data Platform, its architecture as a Data Lakehouse, and the key components.

- What is a data platform and why we need it
- Different Architectures and DataLakehouse
- Components

## Conclusion

To sum it up:

- Prezi Data Platform is founded on stable/state-of-the-art techs and architecture decision
- It allows the introduction of cutting edge technology, as it is pluggable and modular
- It's constantly evolving

Trino Summit 2024

@vincenzocassaro

Thanks!

**00** Definition

**01** Ingestion

**02** Storage

**03** Processing

**04** Orchestration