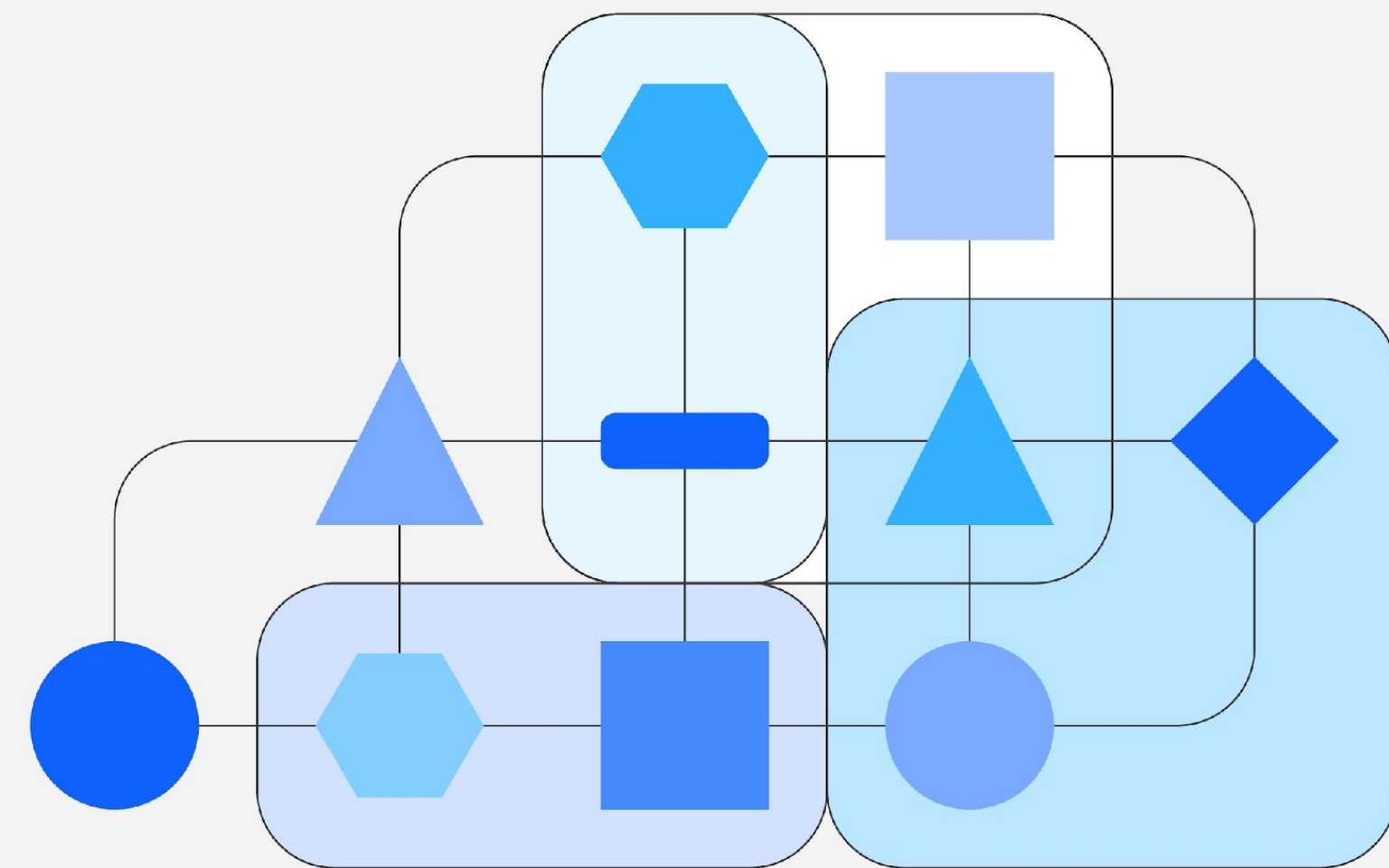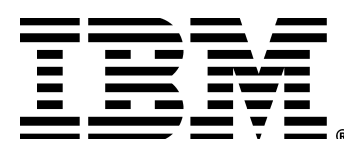# Ensuring Data Quality in Streaming Data Pipelines at Scale

**Ryan Yackel**
Ryan.Yackel@ibm.com
GTM Product Manager
IBM Databand

**Marc Sabate**
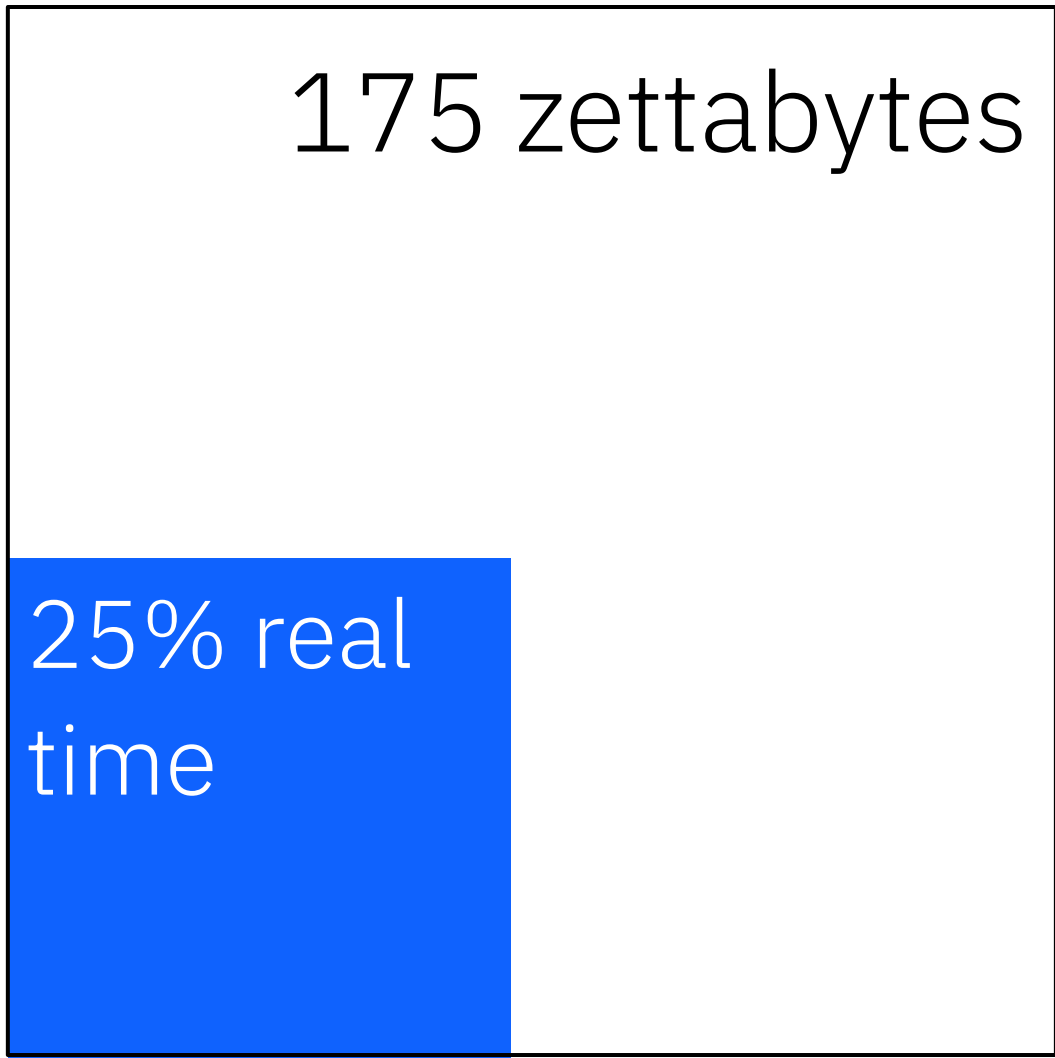Marc.Sabate@ibm.com
Technical Product Manager
IBM StreamSets

IBM

According to the latest estimates,
402.74
million terabytes
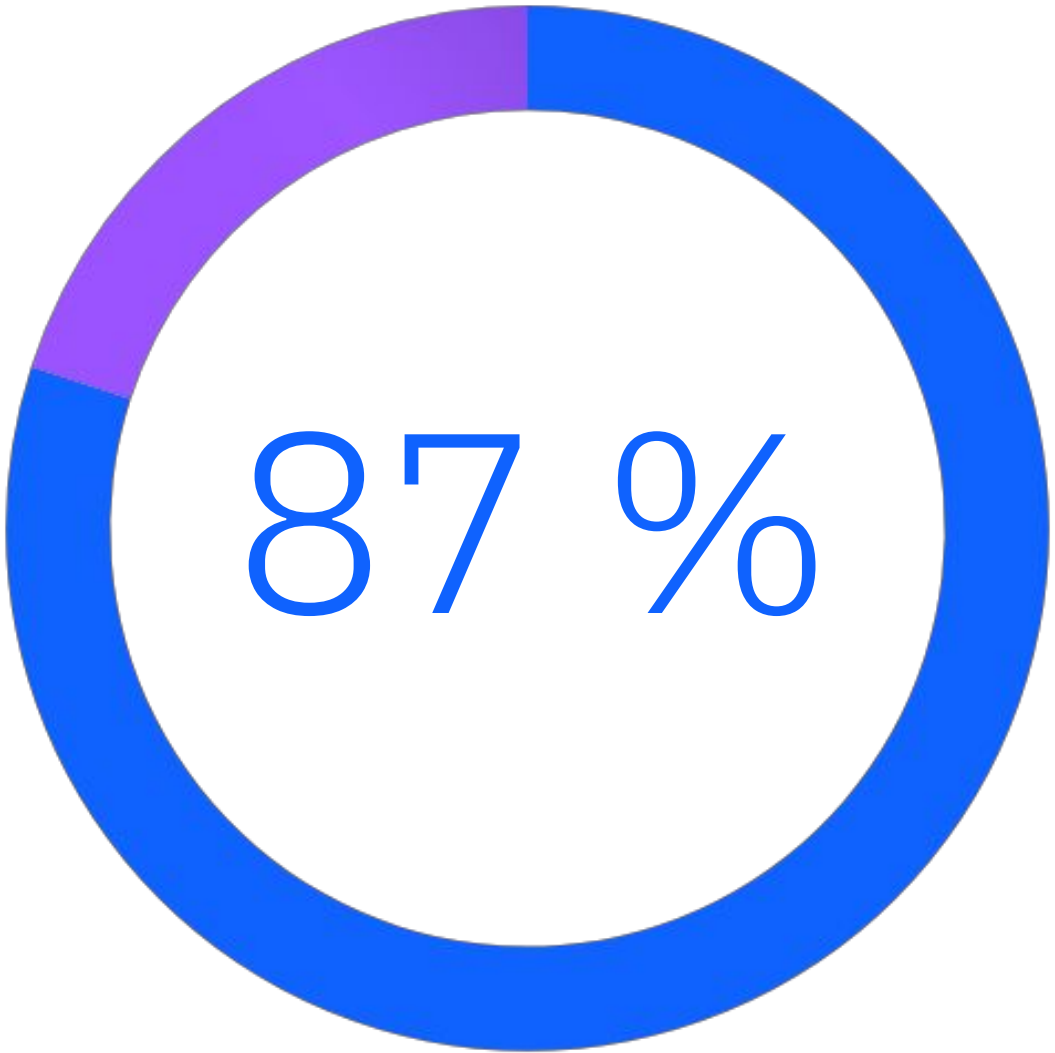of data are created
each day.

Global data volume is projected to grow to 175 zettabytes this year, with nearly 25% of this data being real-time

87% of organizations require data to be ingested and analyzed within one day or faster

175 zettabytes

25% real time

87 %

# While data integration remains a priority for organizations, challenges remain

67% of data leaders believe prioritizing integrating data from various sources is a key priority[1]

## Data access

50% of data analyst time is spent working reactively or trying to find or get access to data[2]

## Data drift

80% of data engineering time is spent on maintaining & fixing existing pipelines due to data drift, creating a costly maintenance burden[3]
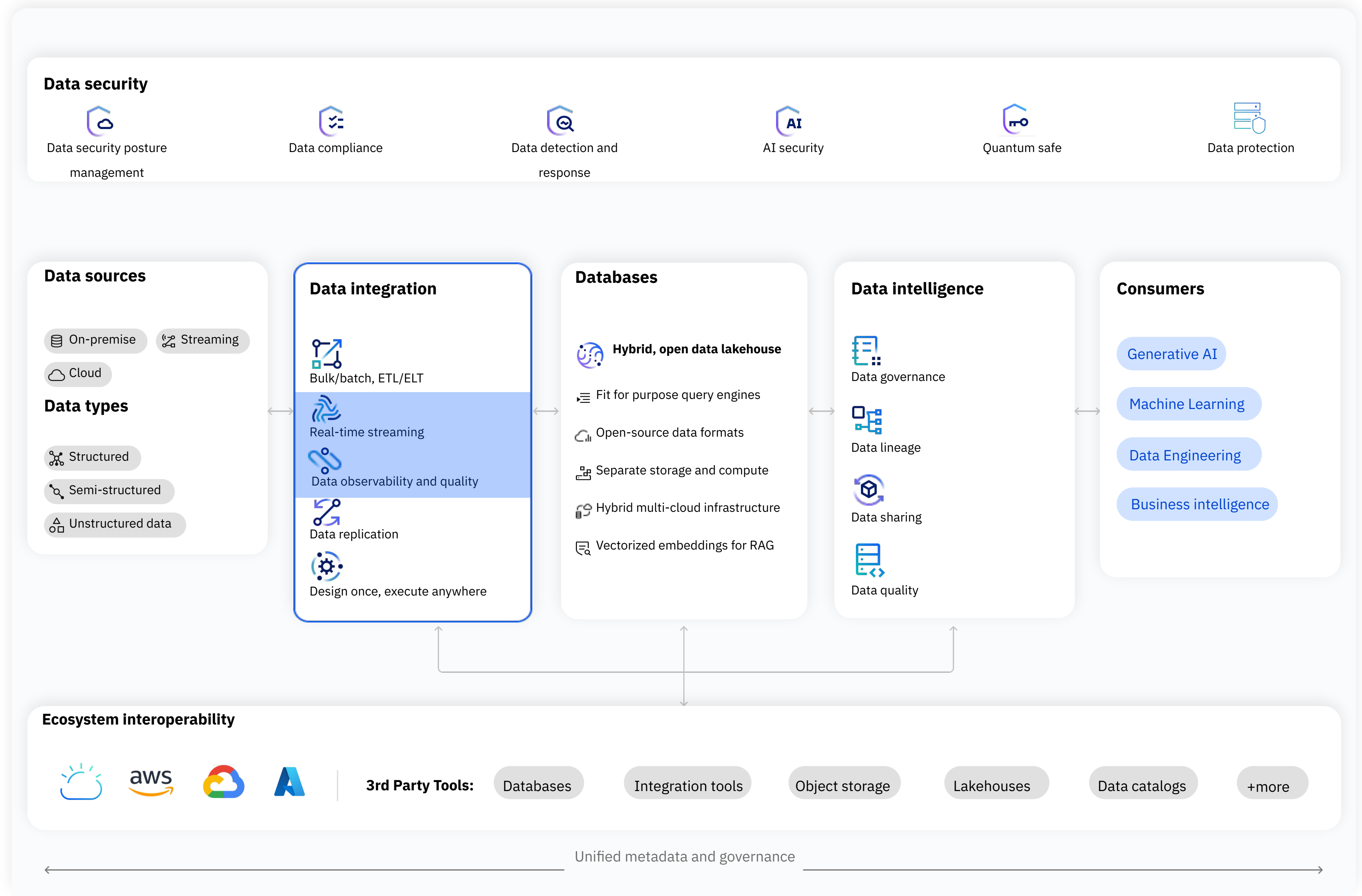
## Data proliferation

With growing data volume and datasets organizations require scalability and faster response times

## Data latency

Need for processing data as it arrives to reduce delays in business decision and operations
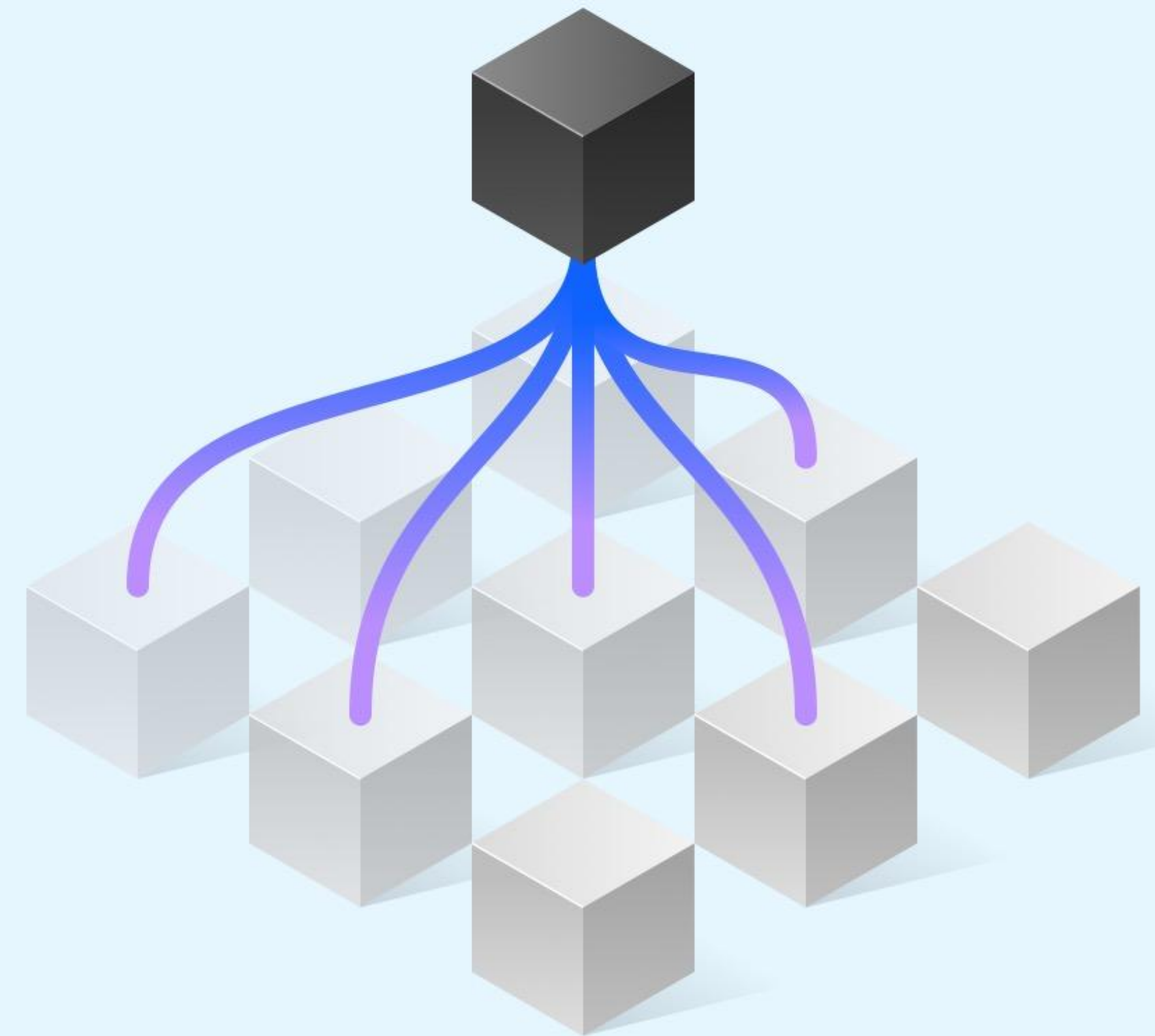
A reimagined data strategy is imperative to meet the demands of contemporary use cases and democratize data access across the enterprise at scale

# Integrate, access, govern, and secure all data types with an open and hybrid data architecture

## Data security

- Data security posture management
- Data compliance
- Data detection and response
- AI security
- Quantum safe
- Data protection

## Data sources

- On-premise
- Streaming
- Cloud

### Data types

- Structured
- Semi-structured
- Unstructured data

## Data integration

- Bulk/batch, ETL/ELT
- Real-time streaming
- Data observability and quality
- Data replication
- Design once, execute anywhere

## Databases

**Hybrid, open data lakehouse**

- Fit for purpose query engines
- Open-source data formats
- Separate storage and compute
- Hybrid multi-cloud infrastructure
- Vectorized embeddings for RAG

## Data intelligence

- Data governance
- Data lineage
- Data sharing
- Data quality

## Consumers

- Generative AI
- Machine Learning
- Data Engineering
- Business intelligence

## Ecosystem interoperability

**3rd Party Tools:** Databases | Integration tools | Object storage | Lakehouses | Data catalogs | +more

Unified metadata and governance

What is Real-Time Data Integration?

IBM defines real-time data integration as the ability to ingest, process, and write data as soon as it's available instead of on an intermittent or scheduled basis.

IBM StreamSets

A no-code/low-code streaming data integration tool
for engineers who *hate* no-code/low code tools
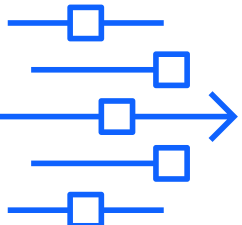
Best-in-class
developer
experience

Capabilities without
compromise

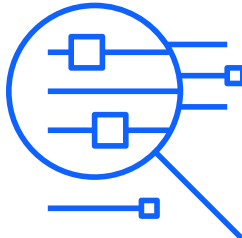Built for enterprise
scale

# IBM StreamSets

## Benefits

↓

Real-time data integration solution for building streaming data pipelines to enhance real-time decision-making and mitigate risks

**Enable real-time data ingestion at scale**

**Reduce data drift with intelligent streaming data pipelines**

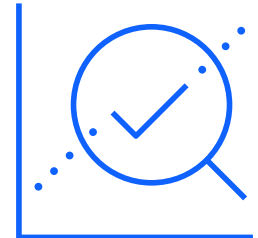**Stream any data from any source**

# What about data quality?
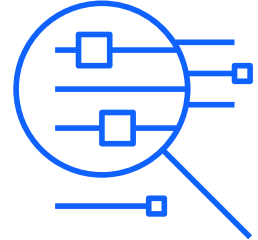
# IBM Databand

## Benefits

↓

Databand empowers data platform teams to deliver reliable and trustworthy data with continuous data observability.

**Detect issues earlier**

**Resolve issues faster**

**Improve data reliability**

# How Databand works

## 1. Collect
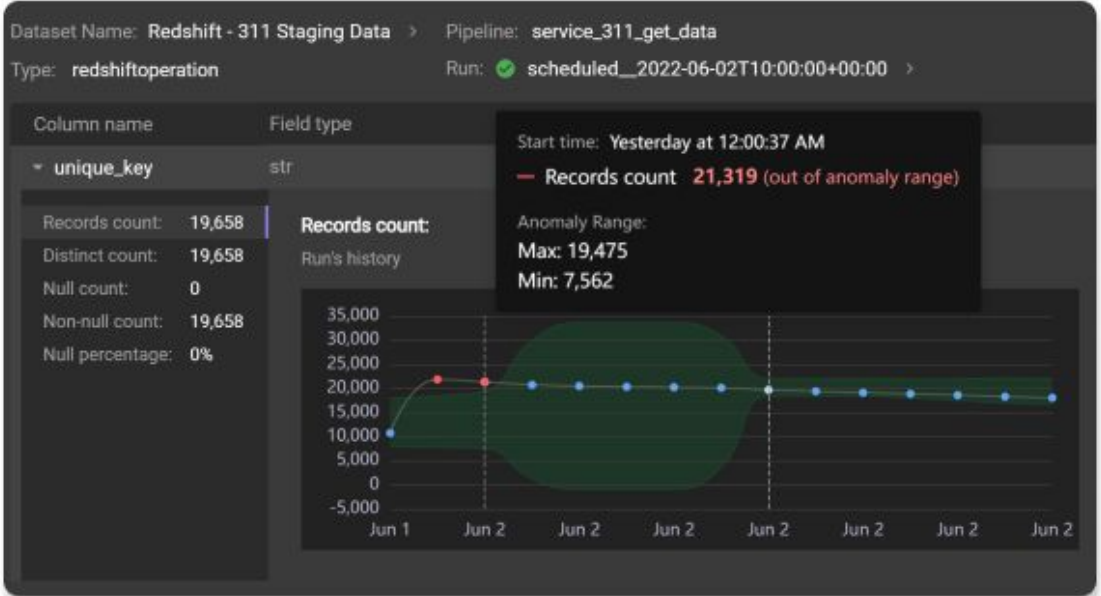**Automatically collect metadata.**
From all key solutions in the modern data stack.



## 2. Profile
**Build historical baseline.**
Based on common data pipeline behavior.



## 4. Resolve
**Resolve the root cause.**
Create smart communication workflows to resolve data quality issues & meet SLAs.



## 3. Alert
**Alert on anomalies and rules.**
Based on deviations or breaches.

# DEMO

## Overview

1. Streaming pipeline ingestion with masked data with AWS S3 destination
2. Streaming pipeline ingestion with product data with Snowflake destination
3. Data quality alerts based on Json format errors
4. Data quality alerts based on infrastructure errors

Marc Sabate
Marc.Sabate@ibm.com
Technical Product Manager
IBM StreamSets

# Q&A Time



**Ryan Yackel**
Ryan.Yackel@ibm.com
GTM Product Manager
IBM Databand



**Marc Sabate**
Marc.Sabate@ibm.com
Technical Product Manager
IBM StreamSets

# Three ways to get started with IBM data integration today

Want to read more about IBM data integration?

[Visit the IBM data integration website →](#)

See how it works for you.

[Start a data integration trial →](#)

Kickstart your project with the IBM technical experts

[Book here →](#)

Autodesk and IBM

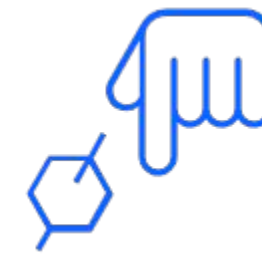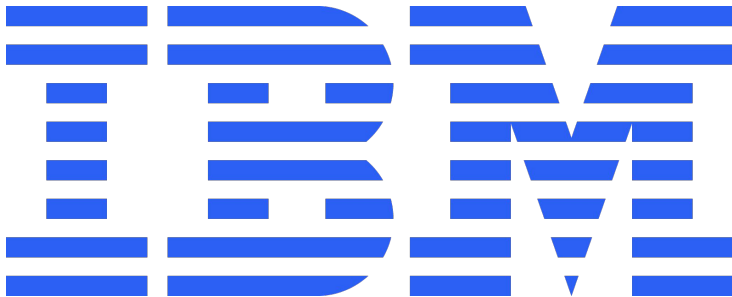# Moving from reactive to proactive data quality

"We got tired of being caught off guard by repeated types of data incidents with no owner to tackle these incidents. With Databand, we've been able to reduce our mean time to detection down to almost zero.

At Autodesk, we encourage innovation, so we saw this as an internal opportunity to bring Databand's data observability to the business."

Senior Manager for Data Engineering
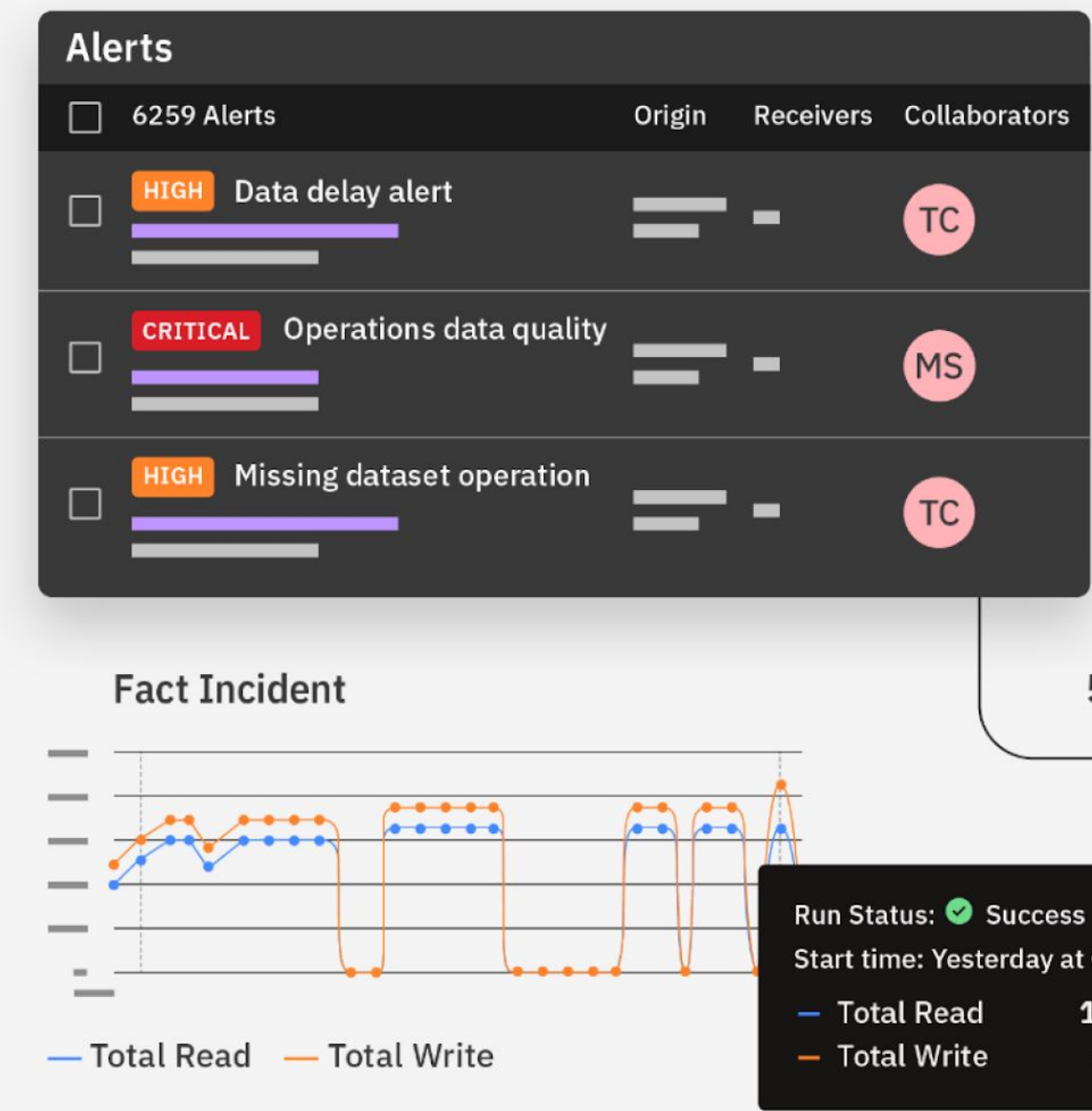and Visualization
Autodesk

## Technologies
– IBM Databand
– Apache Airflow
– Apache Spark
– dbt Core
– Snowflake

## Use cases
– Batch processing monitoring
  Databand extensively monitors production batch processing > 1,000 DAGs for ML and AI teams.
– Inline testing
  The team uses Databand's inline testing capabilities to detect data quality issues in real time.
– Data products support
  Databand supports pipelines that deliver insights and in-product messaging for Autodesk's customers.

## Results
– Reduction in MTTD
  Databand reduced the time to detect data quality issues from days to minutes. This immediate detection allowed the team to address problems before they could cause major disruptions.
– Reduction in MTTR
  The mean time to resolving data issues dropped from weeks to days. Detecting incidents like late-arriving data, schema changes and pipeline failures helps maintain trust and efficiency.
– Cost savings
  Autodesk saw a decrease in cloud consumption costs by detecting issues early and avoiding reruns.

**AUTODESK**

NatWest Group and IBM

## Modernize and compete with challenger banks using innovation and smarter data

NatWest Group PLC (formerly The Royal Bank of Scotland Group PLC), is a majority state-owned British banking and insurance holding company based in Edinburgh, Scotland. The group operates a wide variety of banking brands offering personal and business banking, private banking, insurance and corporate finance, and offers its services to over 19 million retail customers across the UK and Ireland. It also provides business banking services for around 1 in 4 businesses across the UK and Ireland, from startups to multinationals.

Technologies
– IBM StreamSets
– Kafka
–  Hadoop, AWS (S3, EMR)
–  Snowflake, MongoDB

Use cases
– **Risk management and compliance:** Comply with regulations, reduce operational and credit risk exposure, reduce fines
– **Improve customer interactions:** Provide class-leading notification services
– **Optimize customer service:** Understand customer interactions and lifecycles, and offer an industry-leading experience

Results
– **Cost savings:** Reduce customer messaging application cost by £400k
– **Regulatory compliance:** PS2 regulation compliance achieved, delivering 18 million alerts daily. Financial crime compliance achieved
– **Single view** of customer achieved