# Data Platform Architecture Principles

Julien Le Dem

CTO and co-founder Datakin

@J_

# AGENDA
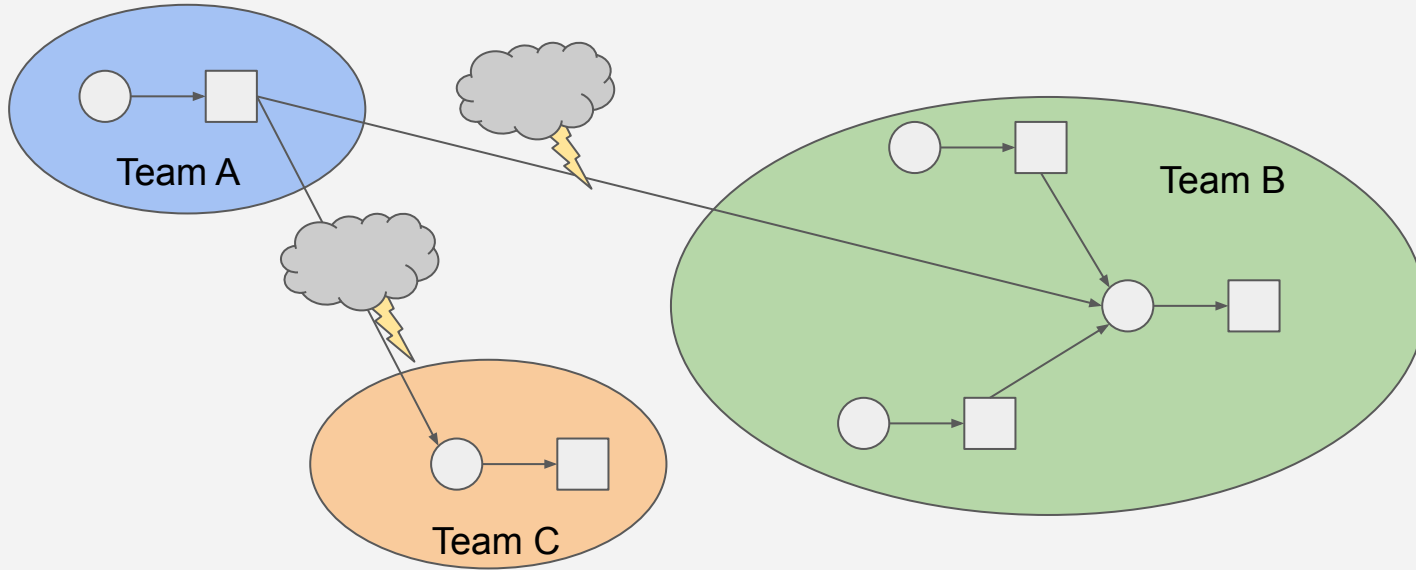
# 01 A Healthy Data Ecosystem

# Team interdependencies

# Explicit contracts

- Schemas

- Shared or Private

- SLA: experimental, production ready
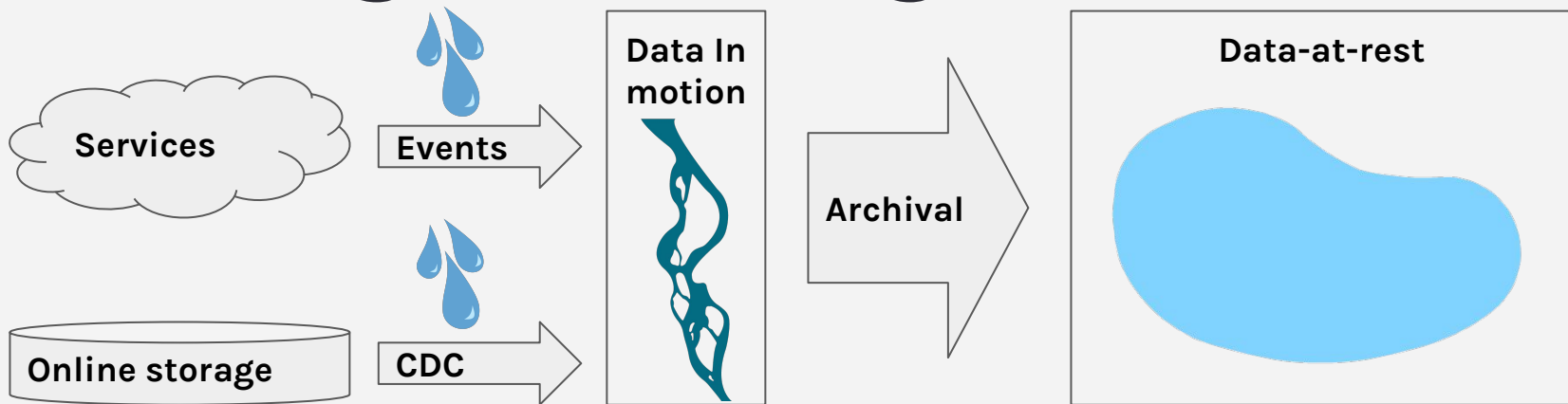
# Understanding dependencies

- **Who do I depend on?**
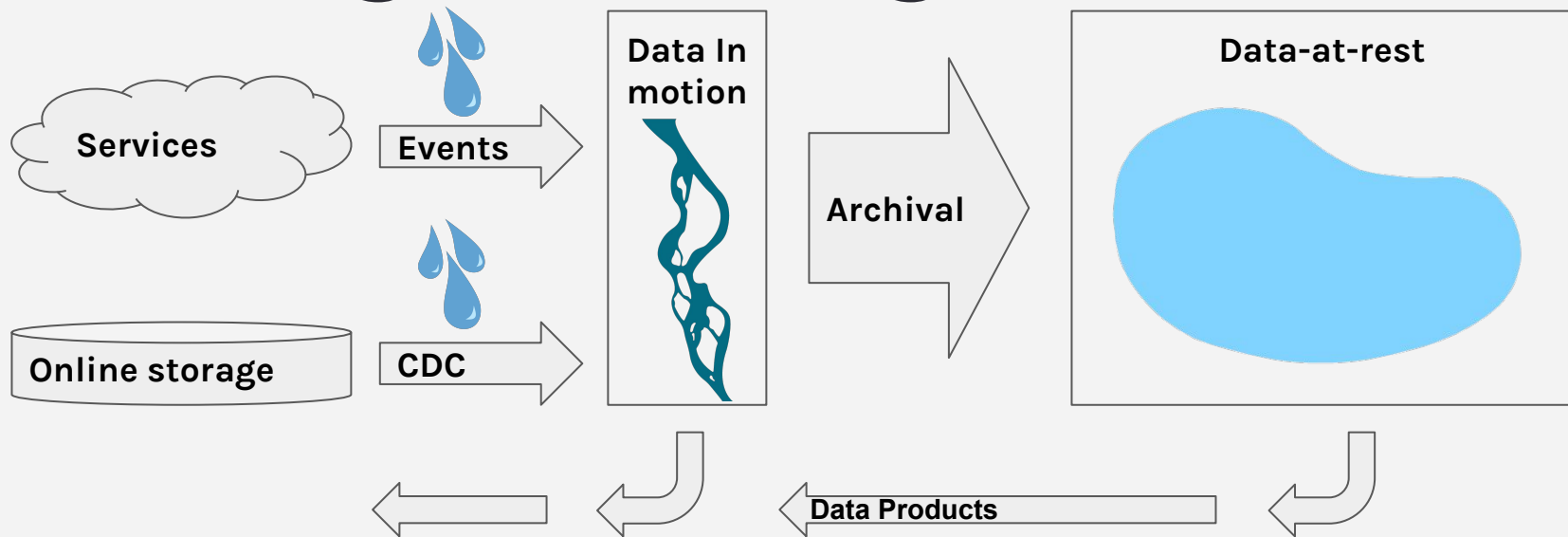- **Who depends on me?**

# Quick iterations

- **Fail safe environment: Easy to undo**
- **Quick troubleshooting**
- **Quick feedback**

# 02  Data Platform Abstractions and Services

# Storage and ingestion



Services

Online storage

Events

CDC

Data In motion

Archival

Data-at-rest

# Storage and ingestion

Services

Events

Online storage

CDC

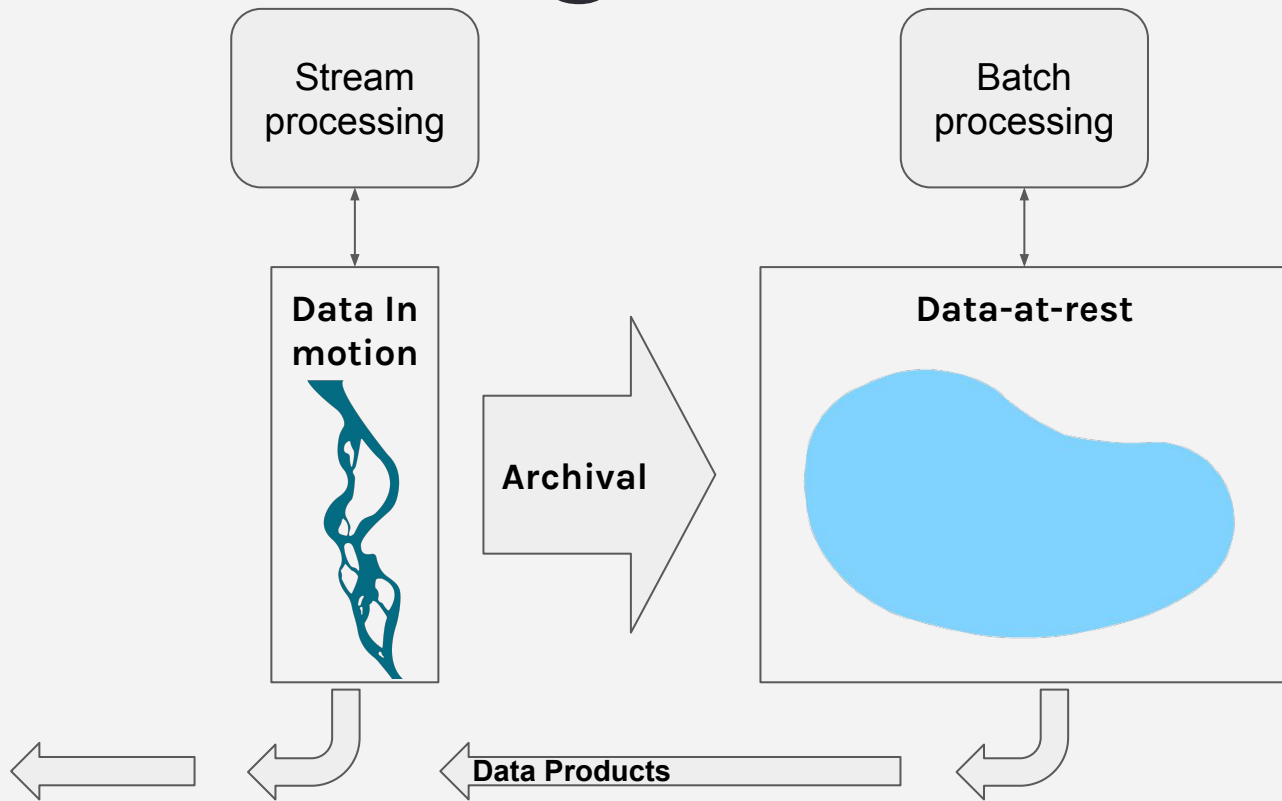Data In motion

Archival

Data-at-rest

Data Products

# Data-in-motion

- Schema registry
- Keyed for CDC
- Horizontally scalable
  - Partitioning
- Candidates: Kafka, Pulsar, …

# Data-at-rest

- Table abstraction:
    - Snapshot Isolation
    - Time travel: can roll back a change
    - Schema evolution
    - Partitioning decoupled from job
- Candidates:
    - Iceberg,
    - Deltalake over cloud blob storage

# Processing

Stream processing

Batch processing

**Data In motion**

**Data-at-rest**

**Archival**

**Data Products**

# Stream processing

- Anti-pattern:
  - Dependencies outside the streaming bubble:
    - Synchronous service calls
    - Database lookup
  - Ingest that data instead (CDC / Domain events)
    - kafka.KTable, flink.DynamicTable
- Candidates:
  - Flink, Spark Streaming, Kafka Streams

# Batch processing

- **Your job as a function: inputs and outputs are parameters.**
  - **Testable transformation:**
  - **Multiple instances in parallel**
- **Atomic runs:**
  - **output is complete or not visible**
- **Understand dependencies**
  - **Jobs depend on their inputs**

# Interactive

- Notebooks:
  - Source control for saving state
  - Repeatable environments: docker images
- Warehouse technology:
  - Decoupled storage and compute
  - Interconnection with data storage
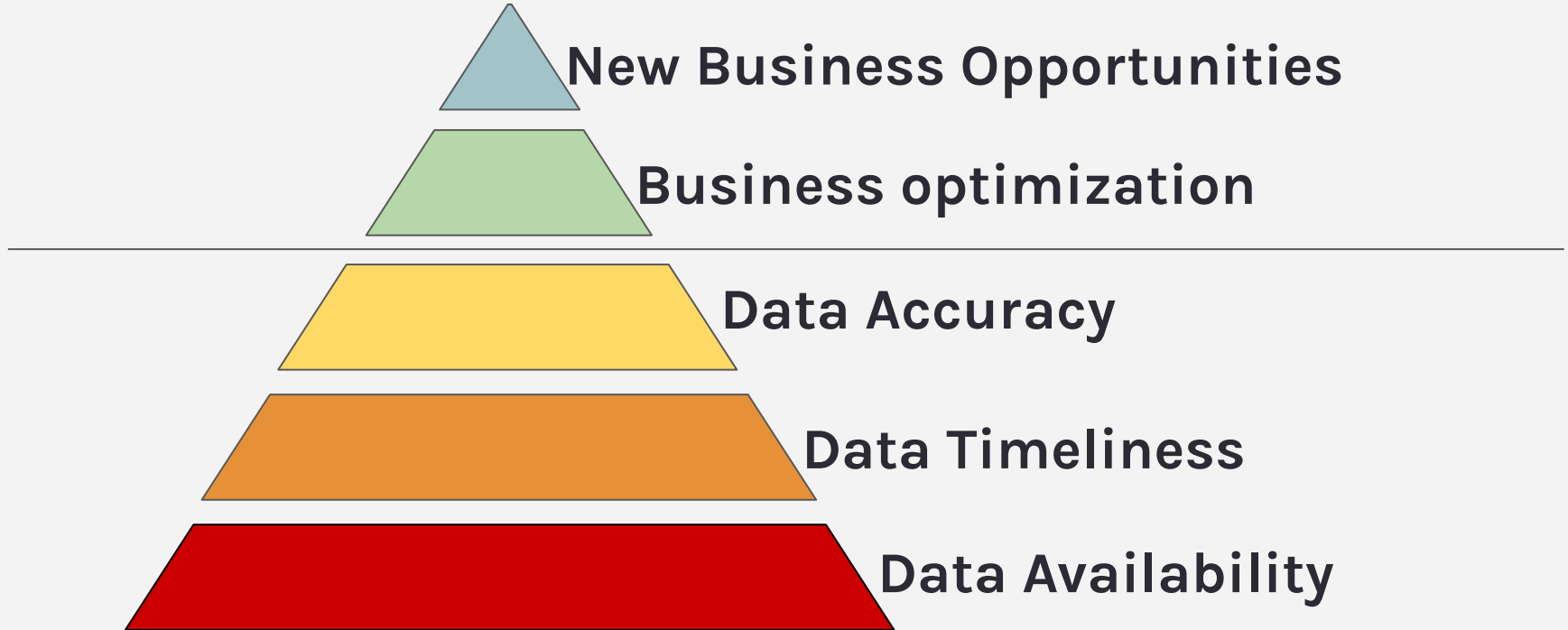
# 03 Observability for data pipelines

# Today: Limited context

DATA

- What is the data source?
- What is the schema?
- Who is the owner?
- How often is it updated?
- Where is it coming from?
- Who is using the data?
- What has changed?

# ~~Maslow's~~ Data hierarchy of needs



New Business Opportunities

Business optimization

Data Accuracy

Data Timeliness

Data Availability

# Observability for data

- Dependencies: Lineage
- availability, timeliness, accuracy
- Change management
    - Schema
    - Code
    - Size
    - Duration

# Observability for data

- **Dependencies: Lineage**
- **availability, timeliness, accuracy**
- **Change management**
  - **Schema**
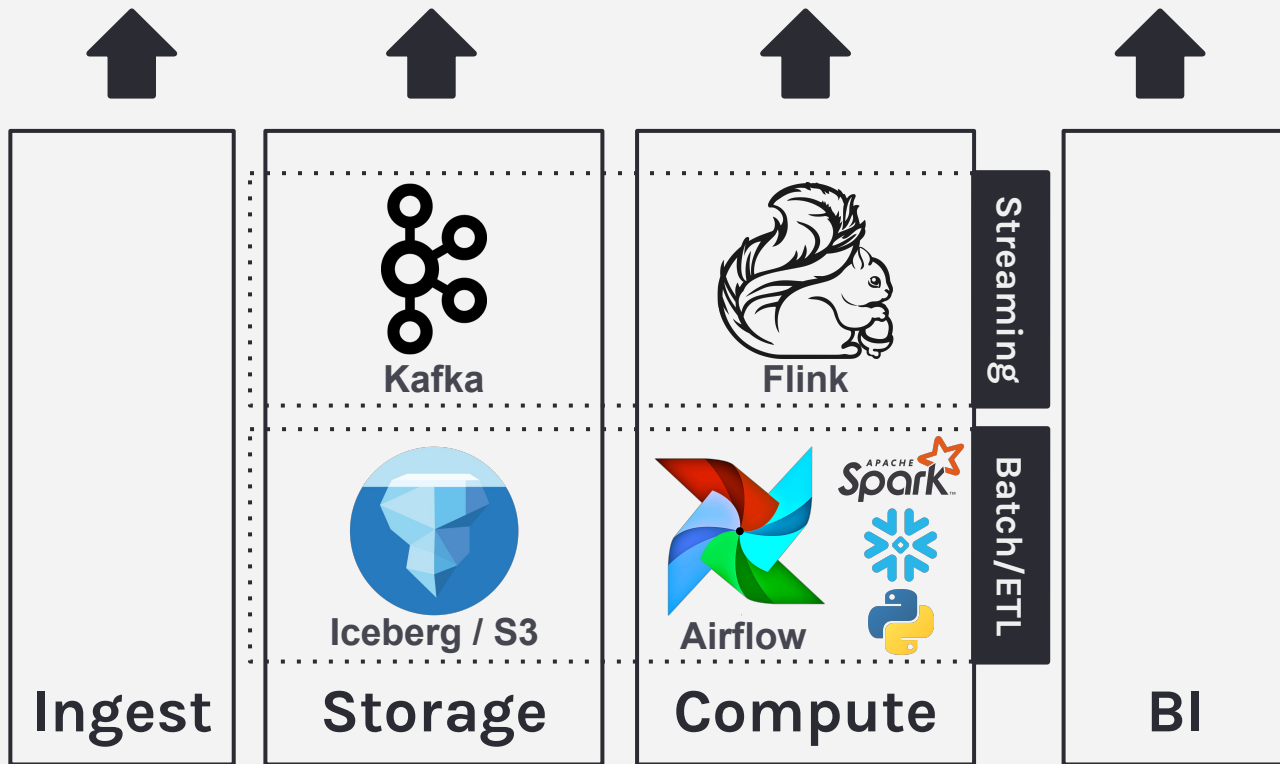  - **Code**
  - **Size**
  - **Duration**
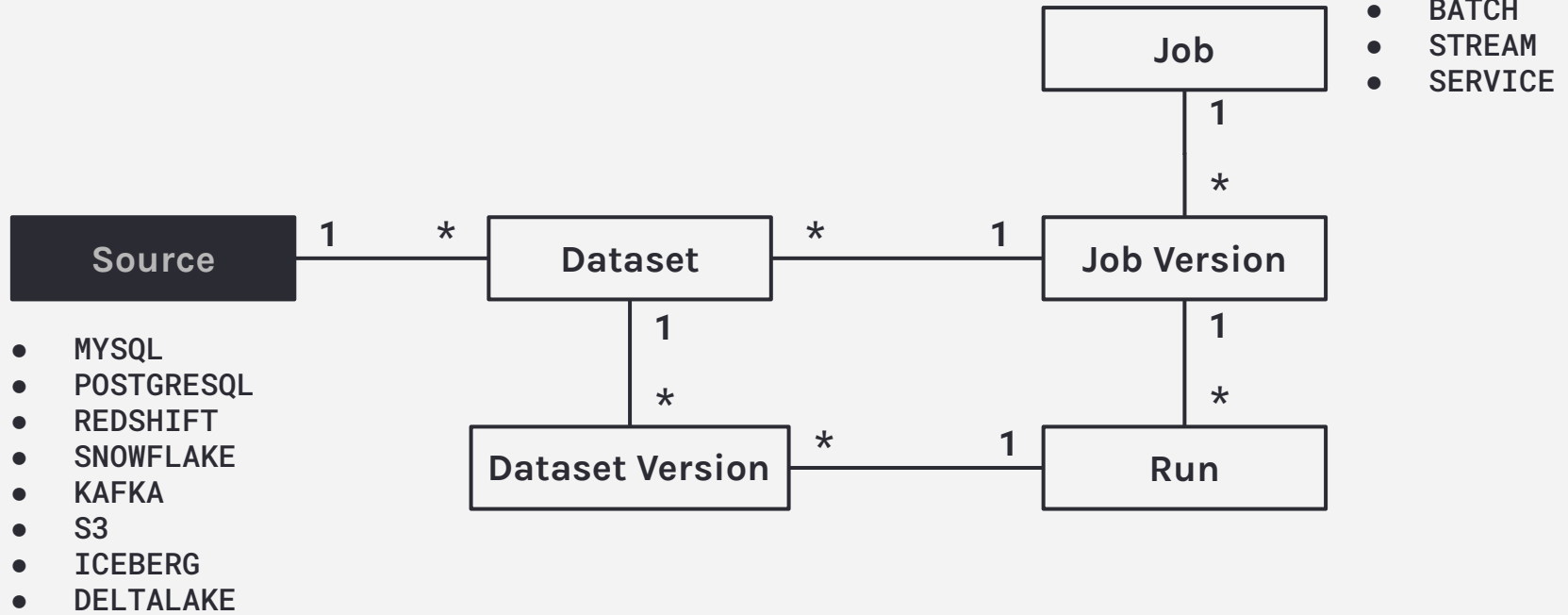
In the **services** world it's called **traces**

# Marquez: Data model

**Job**

- BATCH
- STREAM
- SERVICE

**Source**

- MYSQL
- POSTGRESQL
- REDSHIFT
- SNOWFLAKE
- KAFKA
- S3
- ICEBERG
- DELTALAKE

**Dataset**

**Job Version**

1

*

1 * *  1

**Dataset Version**

1

*

* 1

1

*

**Run**

# Datakin leverages Marquez metadata



- **Marquez standardizes metadata collection**
  - Job runs
  - parameters
  - version
  - inputs / outputs

- **Datakin enables**
  - Understanding operational dependencies
  - Impact analysis
  - Troubleshooting: What has changed since the last time it worked?

Thanks! <o/

# Questions?