# DataOps as a Service

*"People in both fields operate with beliefs and biases.*
*To the extent you can eliminate both and replace*
*them with **data**, you gain a clear advantage."*

Michael Lewis, Moneyball: The Art of Winning an Unfair Game

**ANTONI IVANOV**

Software Engineer / Versatile Data Kit

aivanov@vmware.com
linkedin.com/in/antoni-ivanov

# Agenda

Data Applications

API for Data

SLO and SLAs for Data

DevOps Cycle for Data

Versatile Data Kit

# Applications

Examples

Application



- E-Commerce application
- Mobile app
- Customer relationship management
- Recommendation system

Data



- Databases
- Log files
- Click streams
- Metrics

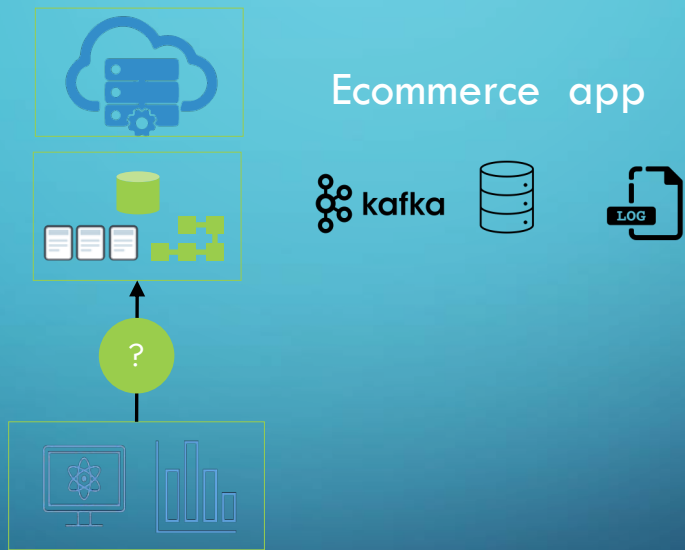# Data Applications

Examples

**Application**

- E-Commerce application
- Mobile app
- Customer relationship management
- Recommendation engine

**Data**

- Databases
- Log files
- Click streams
- Metrics

**Data Applications**

- Usage reporting
- Business intelligence
- Recommendation engine
- Forecasting model

# Data Journey



Ecommerce app

kafka

# The Data Journey

# How do multiple applications and data applications communicate between each other?

# API

**API** is a set of rules, protocols, and tools that allow different software applications to communicate with each other

# API Components

Interface and Contract

Security and Access Control

Usability and Documentation

Monitoring and Operations

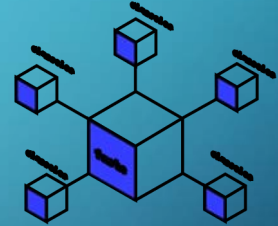# API for Data Example
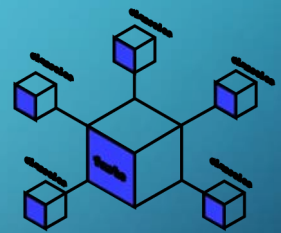
OLTP DB
current products info

S3 Service
legacy products info

Raw Data
(*Data Lake*)

Copy of the data
for further processing

All Product data model
(*Dimensional model*)

# API for Data Example

Table/Entity Name: All Products

Most recent information about each product

**Data Access**

Tables in Database / SQL

Pandas DataFrames

Parquet/Arrow Data Format

## Data Schema

**product_id**: UUID
**name**: String
**category**: String
**price**: Decimal

## Data Semantics

**product_id**: UUID and unique across all records.
**name**: non-empty string representing latest official name
**category**: Must belong to a predefined list of categories.
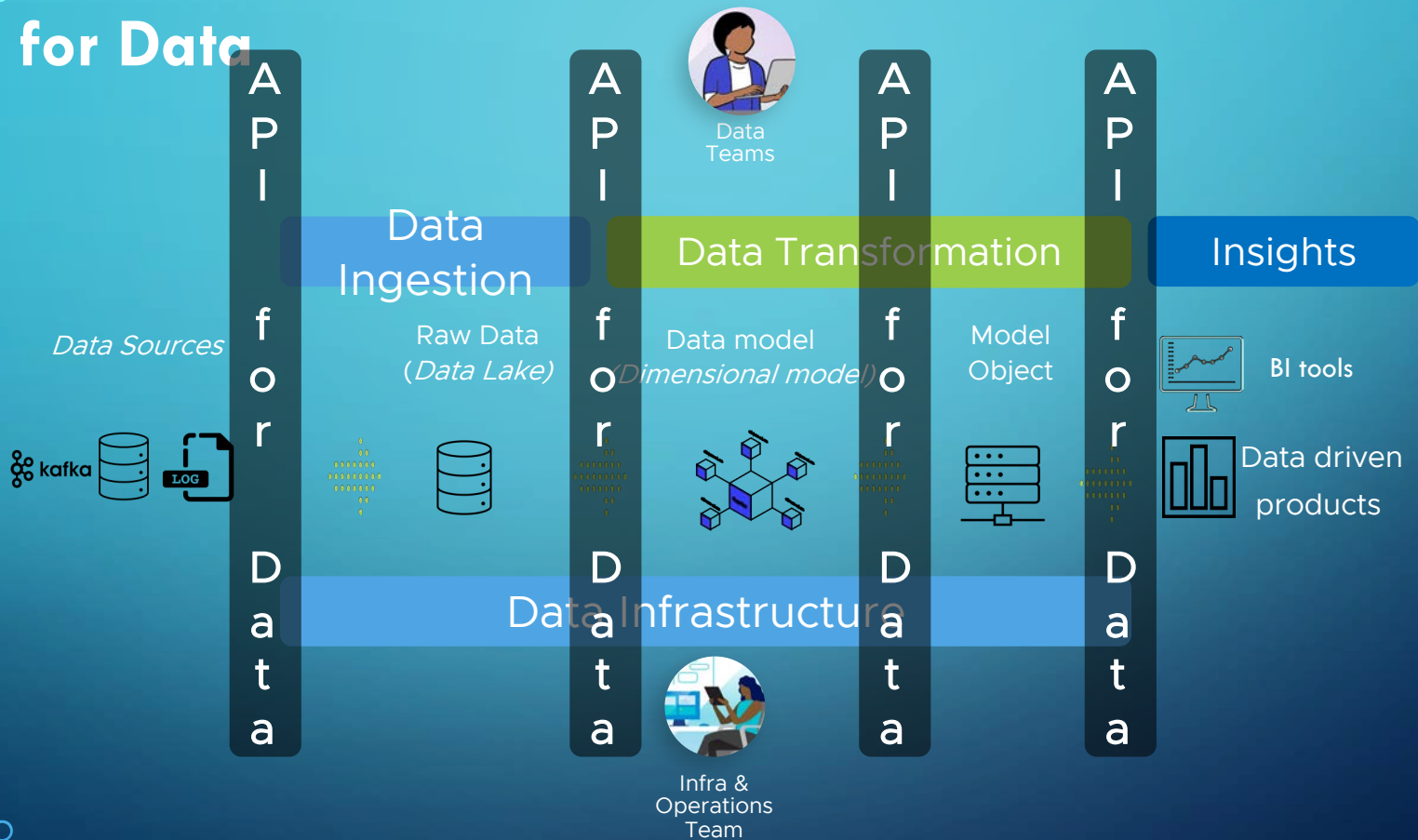**price**: Must be a positive decimal number in currency XXX

# SLOs and SLAs for Data

Data Semantics

Most recent information about each product

**product_id**:  UUID and unique across all records.
**name**: non-empty string representing user facing name
**category**: Must belong to a predefined list of categories.
**price**: Must be a positive decimal number in currency XXX

# SLOs and SLAs for Data

**Data Accuracy SLOs**

**product_id**:  UUID and unique across all records.
**name**: non-empty string representing user facing name
**category**: Must belong to a predefined list of categories.
**price**: Must be a positive decimal number in currency XXX

**Data Availability SLOs**

The "products" table should be queryable 99.9% of the time

**Data Freshness SLOs**

Any changes in the inventory system should be updated in the "products" table within 1 hour
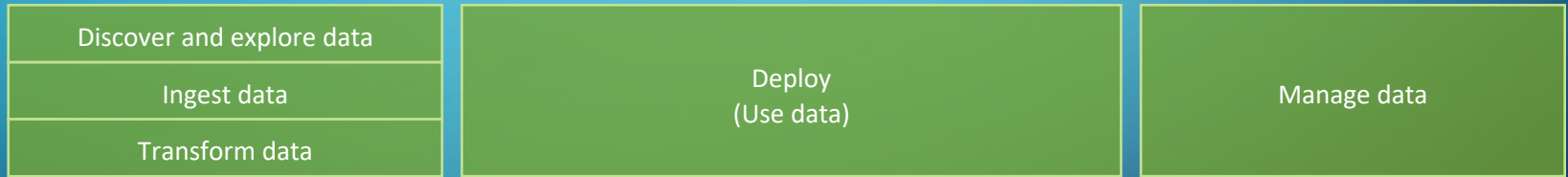
# Data Contracts

## API for Data and SLO/SLAs for data

https://bit.ly/data-contract



Blog posts by Chad Sanderson

# DevOps Cycle for Data

Data Teams

| Plan | → | Code | → | Build | → | Test | → | Release | → | Deploy | → | Operate | → | Monitor |

| Discover and explore data | Deploy (Use data) | Manage data |
| Ingest data | | |
| Transform data | | |

# Versatile Data Kit

Develop Data Jobs

VDK SDK

Deploy and Monitor

Control Plane and Operations UI



**https://github.com/vmware/versatile-data-kit**

# WHAT ARE WE GOING TO DO?



```
1    INSERT INTO tableName (sddc_sk,active_from,active_to,sddc_id,updated_by_user_id,s
•    '500'),(sddc_sk,active_from,active_to,sddc_id,updated_by_user_id,state,is_nsxt,cl
2    ....
3

•    '2', 'RUNNING', 'TRUE', 'AWS', '497'),('sddc03-v01', '2.01.19', '3.01.19', '3',
5208 ,('sddc01-v01', '1.01.19', '2.01.19', '1', '9', 'STOPPED', 'FALSE', 'AWS', '500
•    '2', 'RUNNING', 'TRUE', 'AWS', '497'),('sddc03-v01', '2.01.19', '3.01.19', '3',
5209 ,('sddc01-v01', '1.01.19', '2.01.19', '1', '9', 'STOPPED', 'FALSE', 'AWS', '500
•    '2', 'RUNNING', 'TRUE', 'AWS', '497'),('sddc03-v01', '2.01.19', '3.01.19', '3',
```

😭 ????

# Versatile Data Kit Control Service

# Versatile Data Kit Control Service
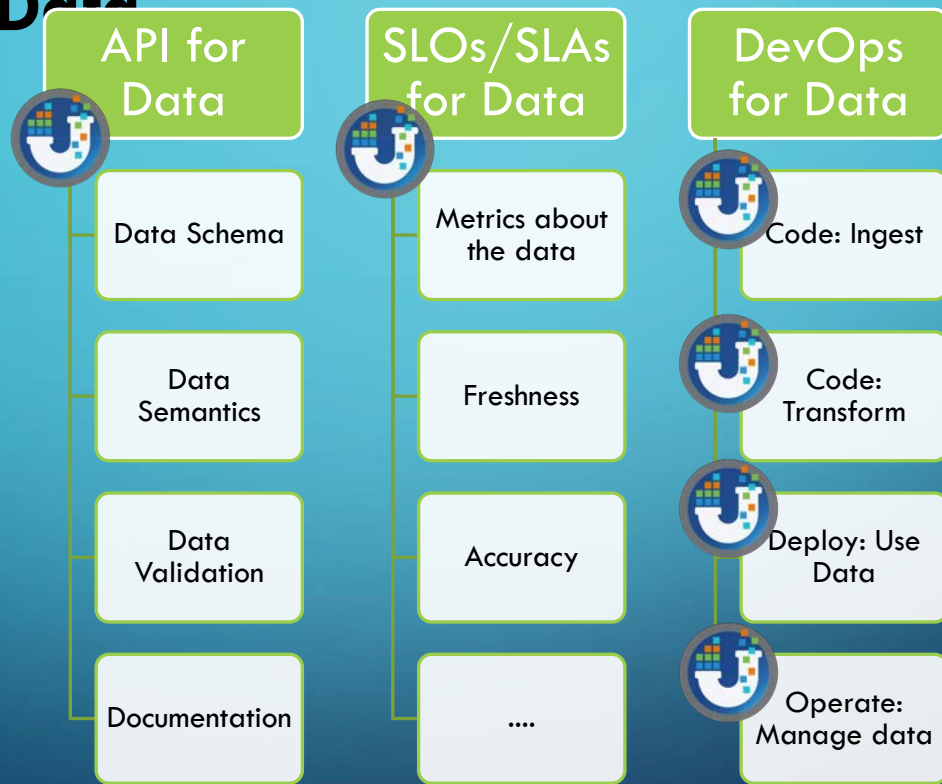
Plan → Code → Build → Test → Release → Deploy → Operate → Monitor

```
2  ≫  FROM versatiledatakit/job-builder
3     💡
4     # Run system test before accepting the new job code
5     RUN pytest system_test.py || die 'Failed system test'
6
7     #   Remove execution privileges from files during container build
8     RUN chmod -R -x $job_name/
9
```

Establish standard system tests and security hardening

# DataOps for Data

| API for Data | SLOs/SLAs for Data | DevOps for Data |
|---|---|---|
| Data Schema | Metrics about the data | Code: Ingest |
| Data Semantics | Freshness | Code: Transform |
| Data Validation | Accuracy | Deploy: Use Data |
| Documentation | .... | Operate: Manage data |

# Thank you



https://github.com/vmware/versatile-data-kit

# DataOps: DevOps for Data

*Ineffcient Operations*      *Stalled development.*

Domain knowledge

Implement business logic

Optimizes for agility and speed

DevOps & Infrastructure knowledge

Maintain infrastructure

Optimizes reliability, availability and security

Wall of conflict

Data Team

Data Team

Data Team

Data Team

Data Team

Infra & Operations Team

*Blurred lines of responsibility*

# Building a data app is hard

# Complexity increases