dojo®

**Building a business-critical data platform to process over £50bn in card transactions**
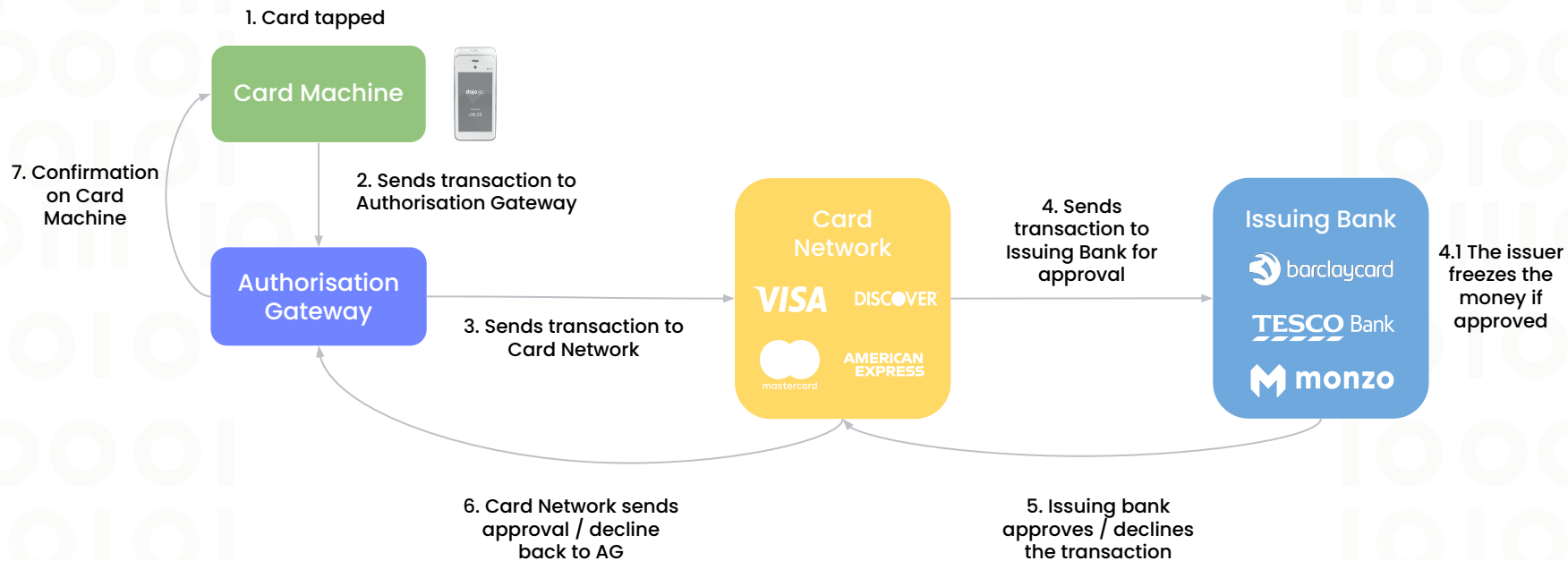
# dojo.

www.dojo.tech

# One of the largest Fintechs in Europe

Enabling 130,000 businesses to take card payments from 4 million consumers per day



dojo.

# Taking a card payment is complicated and highly regulated



1. Card tapped

**Card Machine**

7. Confirmation on Card Machine

2. Sends transaction to Authorisation Gateway

**Authorisation Gateway**

3. Sends transaction to Card Network

**Card Network**

VISA  DISCOVER

mastercard  AMERICAN EXPRESS

4. Sends transaction to Issuing Bank for approval

**Issuing Bank**

barclaycard

TESCO Bank

monzo

4.1 The issuer freezes the money if approved

6. Card Network sends approval / decline back to AG

5. Issuing bank approves / declines the transaction

dojo.

# Challenges

Building a nuclear power station - it cannot fail

| Regulatory | Complexity | Scalability |
|:---:|:---:|:---:|

**FCA** FINANCIAL CONDUCT AUTHORITY

Safeguard customer funds from working capital at all times

**PCI** Security Standards Council ®
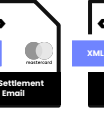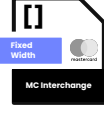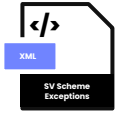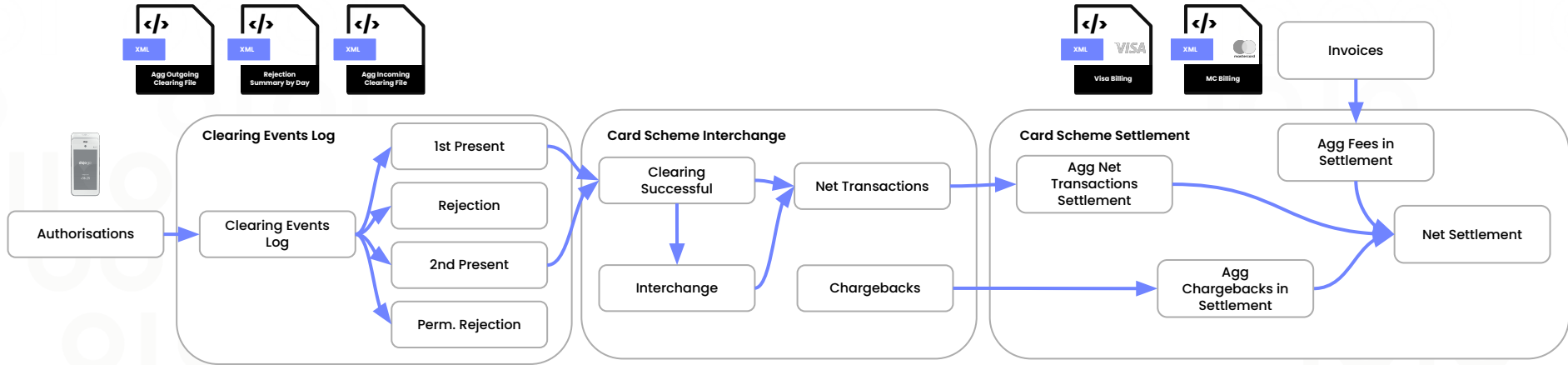
Data contains raw card numbers (PANs)

**Over 1000+ schemas**
Vast number of different proprietary file formats that are fragile and frequently change

**Multiple file formats and sizes**
Files vary from few megabyte XML all the way through to multi-gigabyte proprietary files with tens of millions of rows.
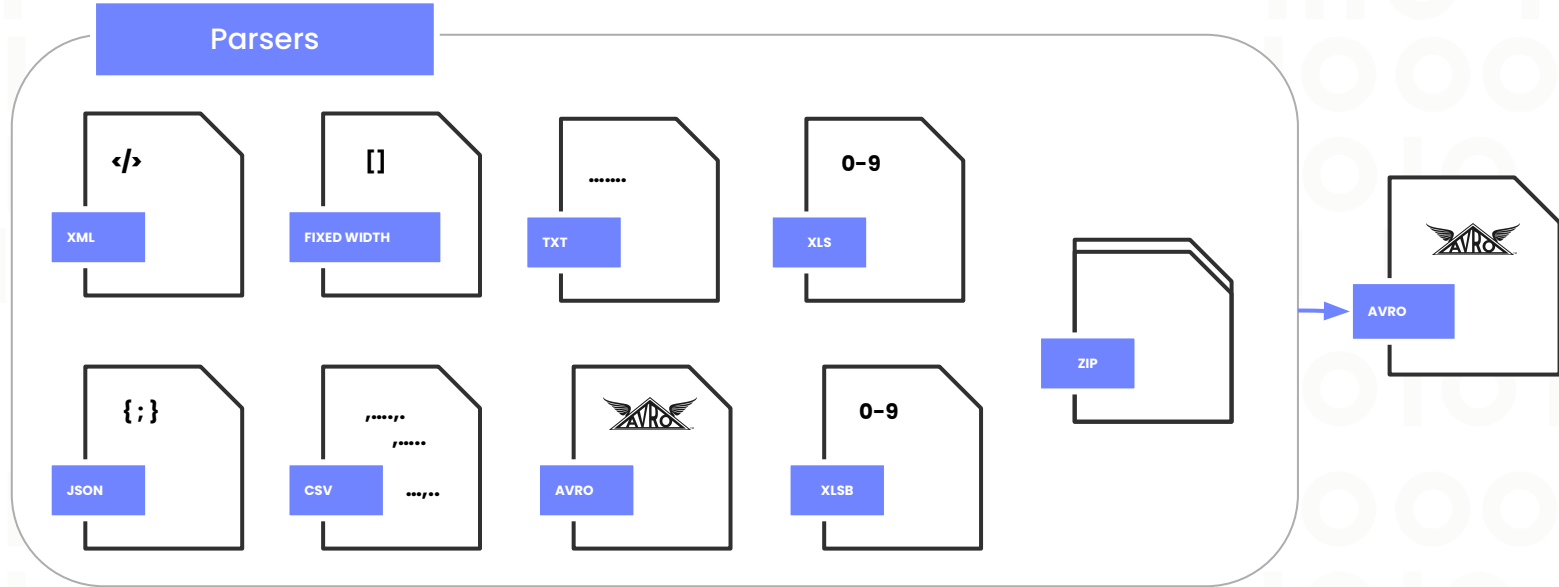
**Unpredictable demand**
Unpredictable surges in transaction volumes

**Business Critical (99.99%)**-
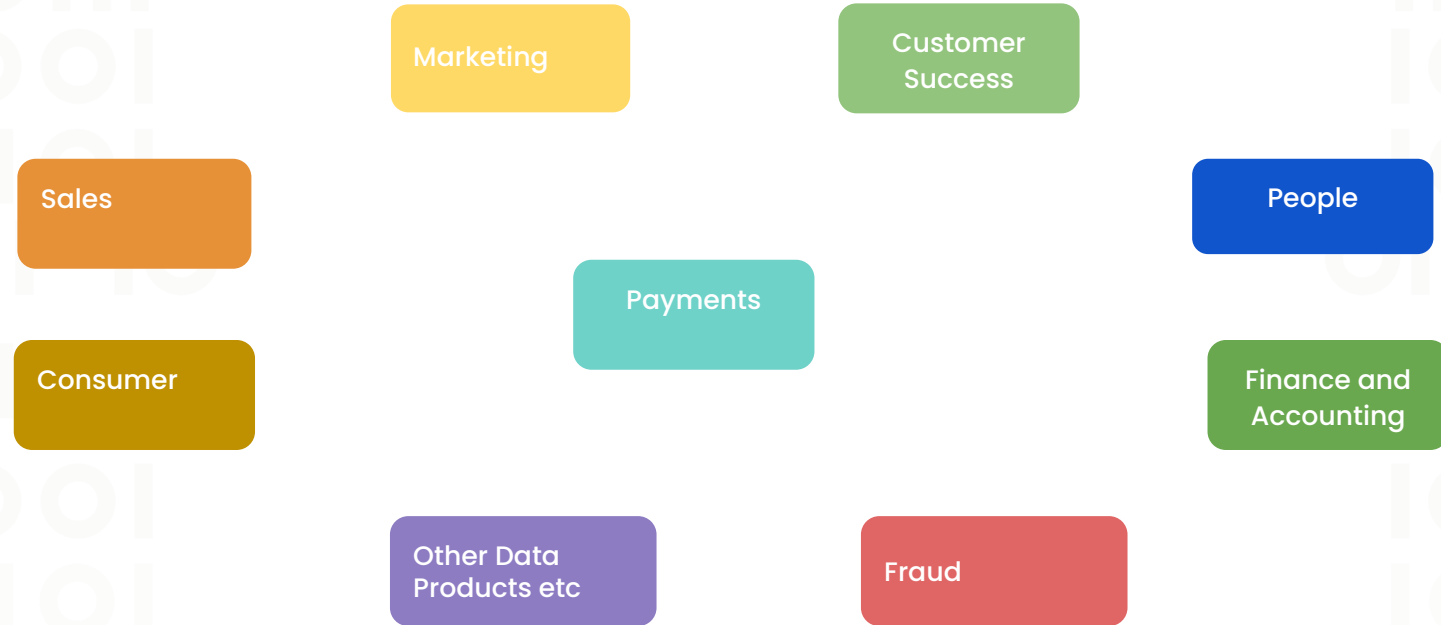Internally we refer to this as building a nuclear power station - it cannot fail.

dojo.

# Abstracting away the complexity

Transformation into a consistent file format - Avro



dojo.

# other Data Domains as well ...

Marketing

Customer Success

Sales

People

Payments

Consumer

Finance and Accounting

Other Data Products etc

Fraud

dojo.

# Data Infrastructure Generations

**Enterprise Data Warehouse**

**Big Data Ecosystem**

**Centralised Data Platform**

**Batch and Real Time Streaming**

**Cloud Based Managed Services**

dojo.

# Challenges: Central Data Platform

**Centralised Team Ownership**

**Data Quality, Accountability and Democratisation**

Scalability

**Siloed and Specialised Data Platform Team**

**Data Quality Issues**

**Data Volume and New Data Sources**

**Stretched Data Platform Team**

**Lack of Accountability when it comes to data issues**

**Cost Efficiency**

**Delays in Data Access and Insights**

dojo.

# What's Next

**Federated Governance policies**

**Domain Ownership**

**Domain  Data Quality and Observability**

**Reduce Complexity**

**Scalability**

**Data Democratisation**

**Innovation and Agility**

**Data Accountability and better Support**

**Data Observability**

**Integrations and Data Contracts**

**Generic Data Infrastructure or Self Serve Data Platforms**

dojo.

# Data Mesh

**Domain Ownership**

**Data as a Product**

**Self-serve Data Platform**

**Federated Computational Governance**

dojo.

# Modern Data Stack is Broken



The 2023 MAD (ML/AI/Data) Landscape

# What should we do ?

Build Small and Go Big

Open Source Tools

Early Feedbacks

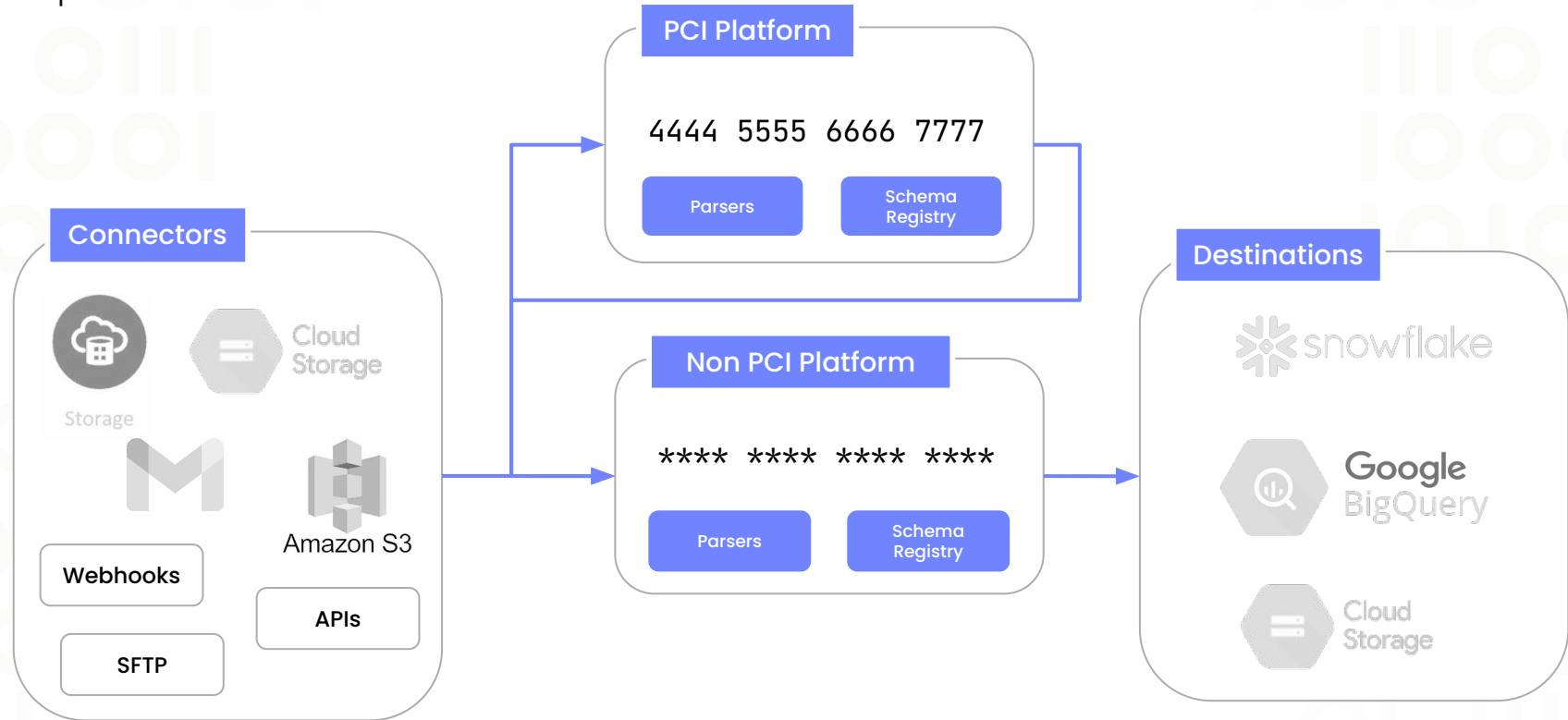Kubernetes is your friend
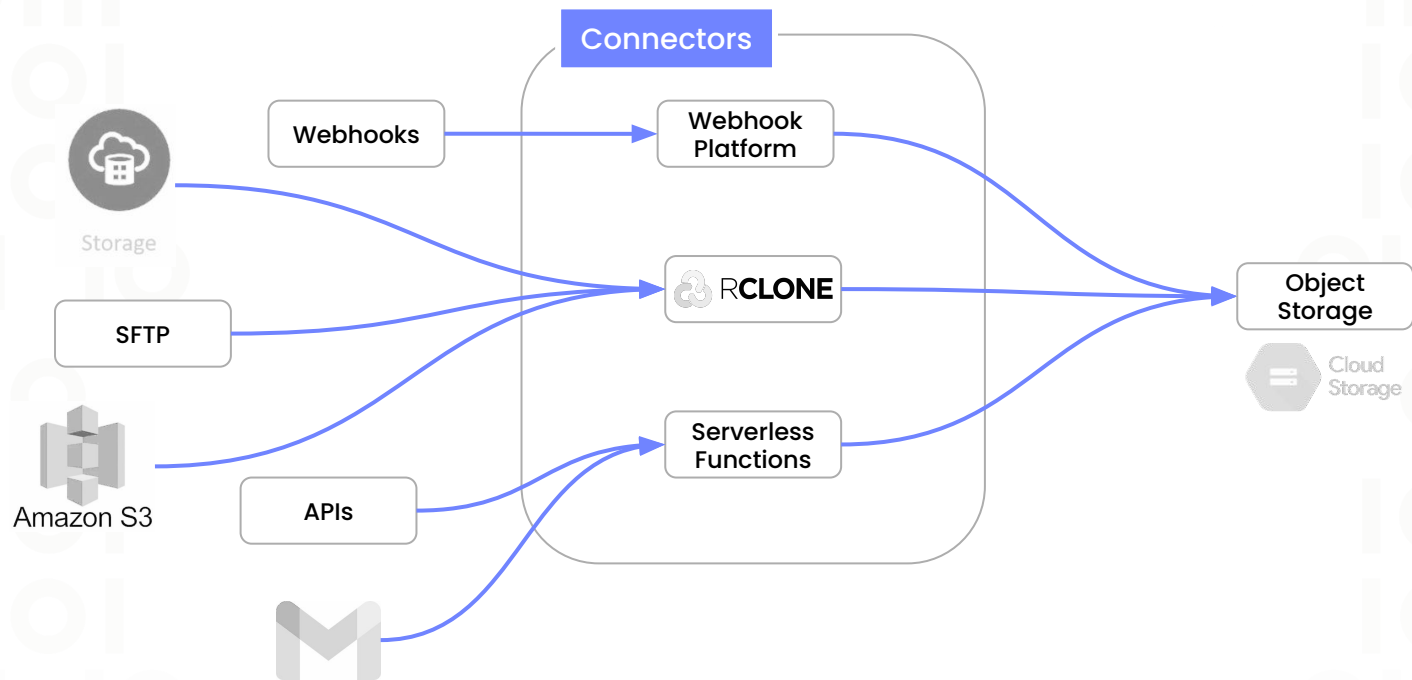
Managed Cloud Services can be handy

Cloud Agnostic

dojo.

# Platform Overview

Component based architecture

**Connectors**

- Storage
- Cloud Storage
- Amazon S3
- Webhooks
- APIs
- SFTP

**PCI Platform**

4444 5555 6666 7777

- Parsers
- Schema Registry

**Non PCI Platform**

**** **** **** ****

- Parsers
- Schema Registry

**Destinations**

- snowflake
- Google BigQuery
- Cloud Storage

dojo.

# Connectors Overview

A unified approach to ingesting data from multiple sources

# PCI DSS Level 1

Process, store, or transmit credit card or cardholder data maintain a secure environment
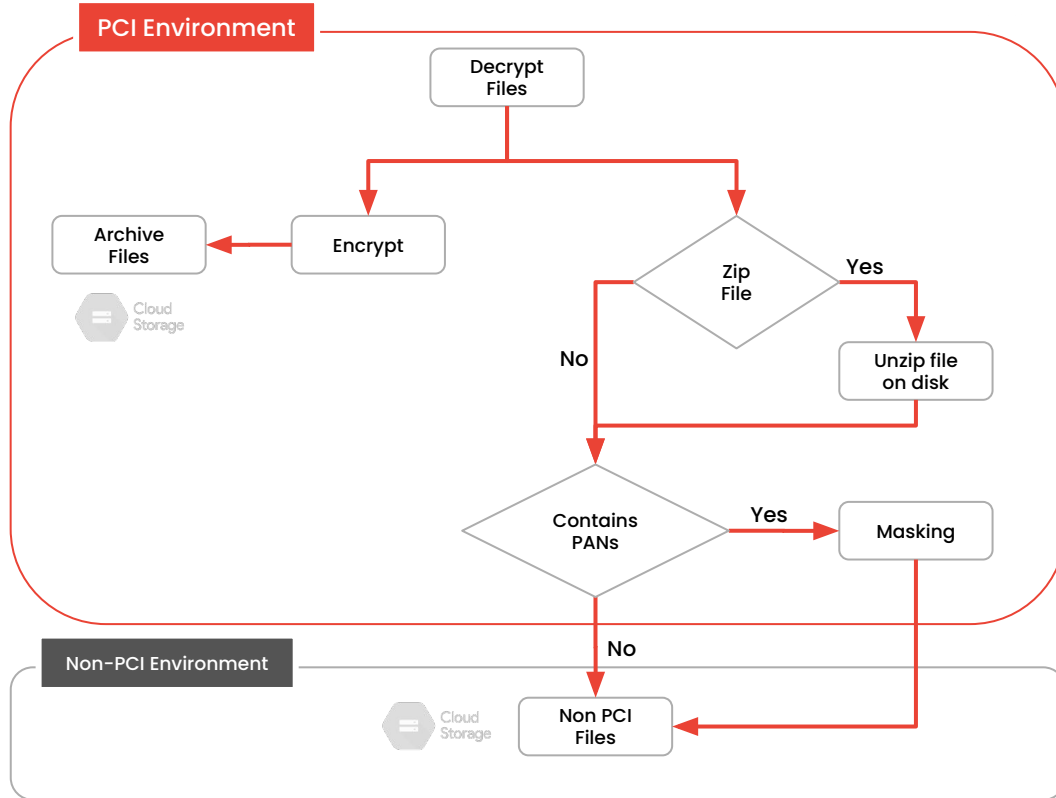
Card data transmitted securely into the Data Platform

Strong Encryption, Monitoring and Security testing of Data Platform

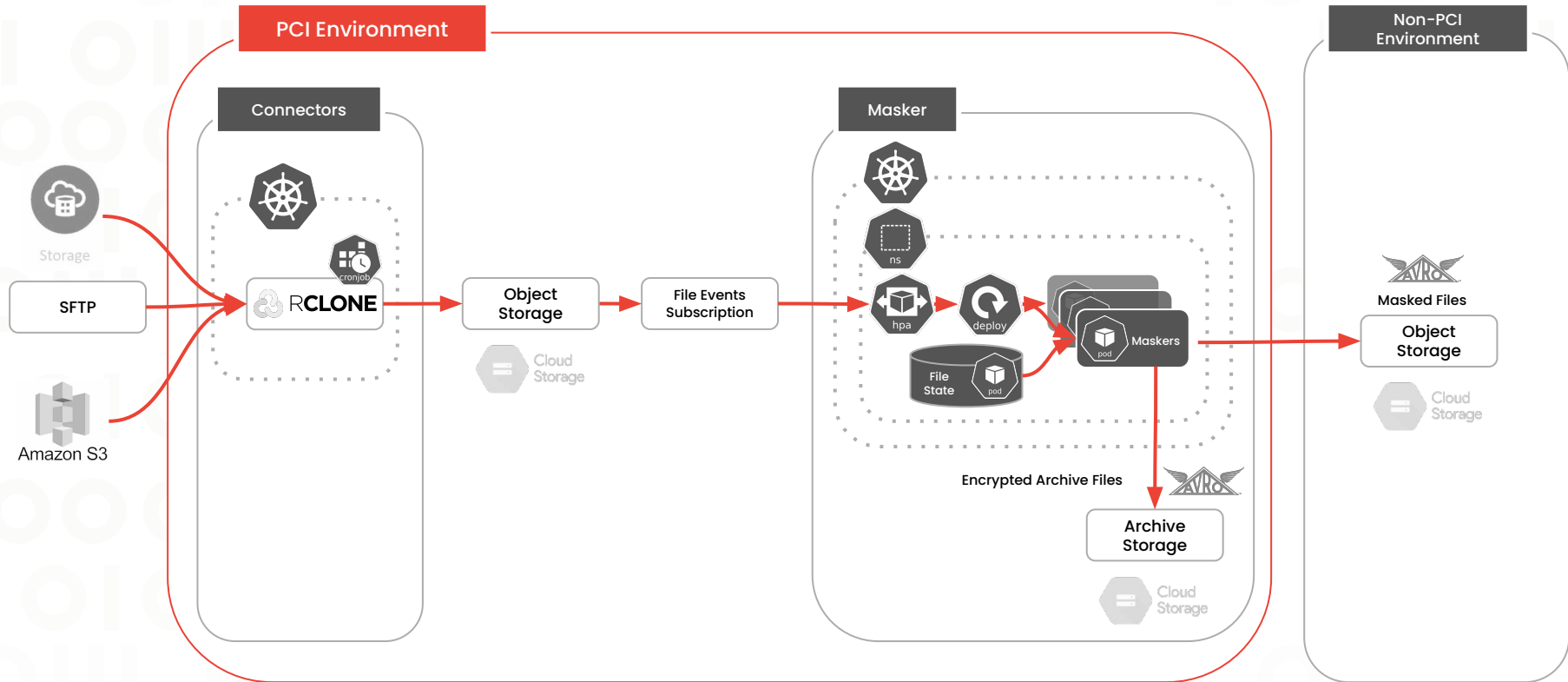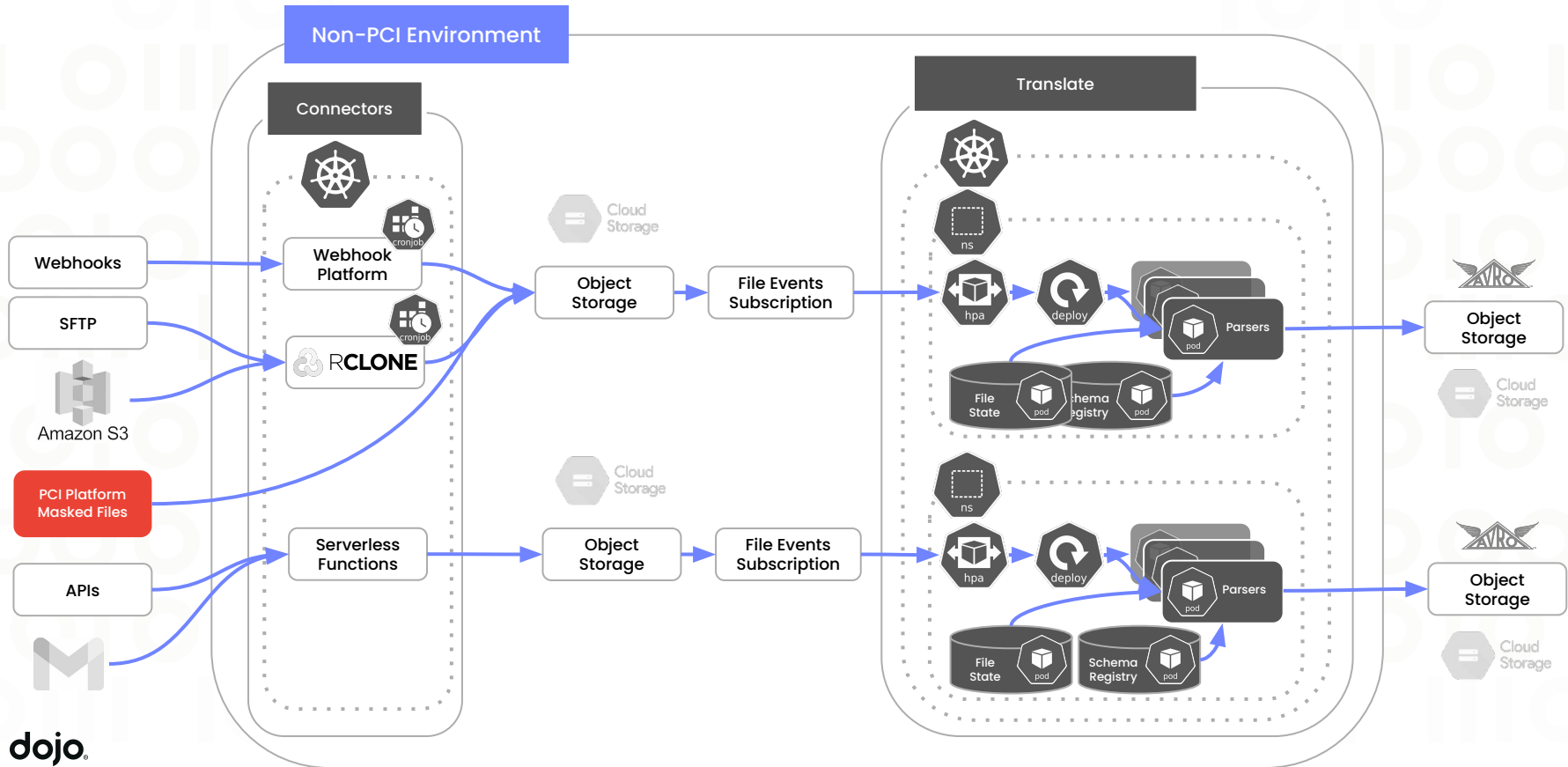Yearly Audits to verify the security of the Data Platform

PCi Security Standards Council®

dojo.

# PCI Management Process

Secure and prevent the egress of PCI sensitive data

# PCI Platform Overview

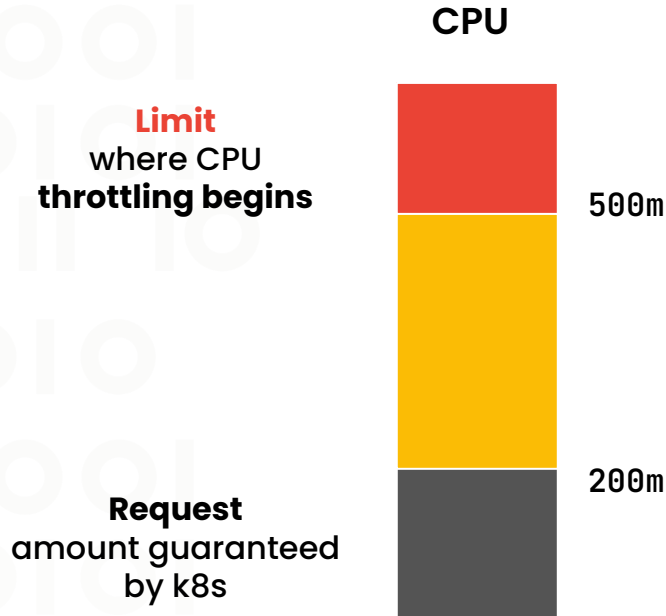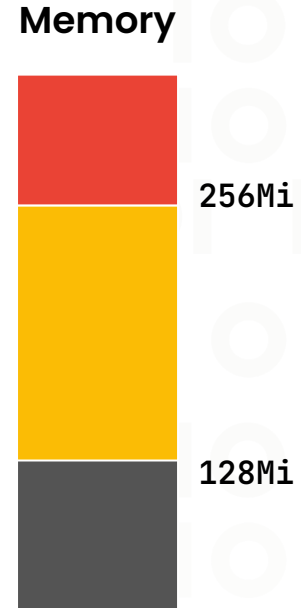# Transformation Platform Overview

# Autoscaling Challenges

Capacity planning is essential to scaling efficiently

## CPU

**Limit**
where CPU
**throttling begins**

500m

200m

**Request**
amount guaranteed
by k8s

## Memory

**Limit**
where **processes are killed**
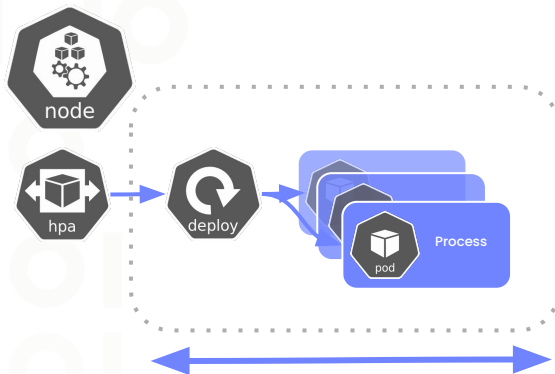
256Mi

128Mi

**Request**
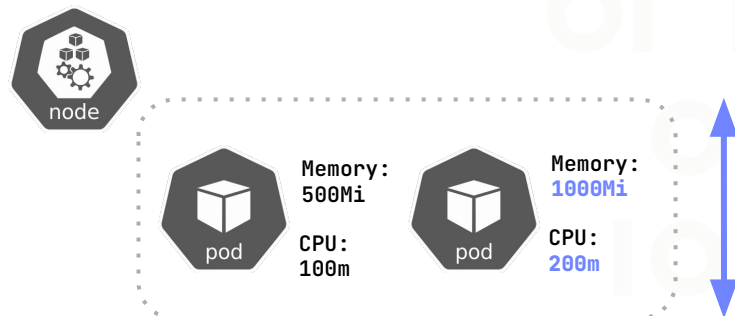amount guaranteed
by k8s

dojo.

# Autoscaling

Two approaches to autoscaling in Kubernetes - Horizontal (HPA) or Vertical (VPA)

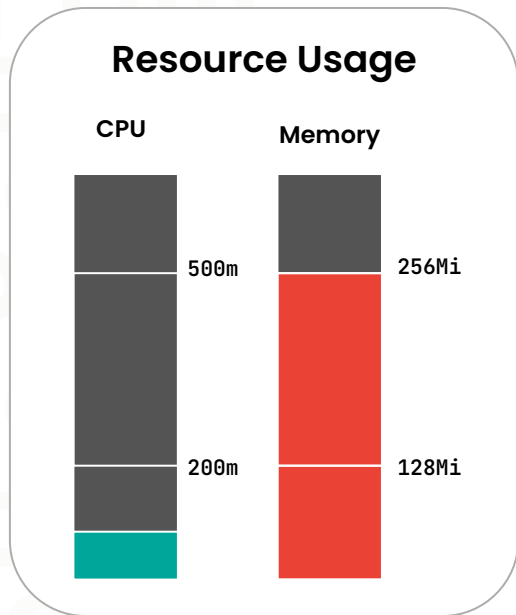**Scale Out:** Increase / decrease number of pods based on metrics

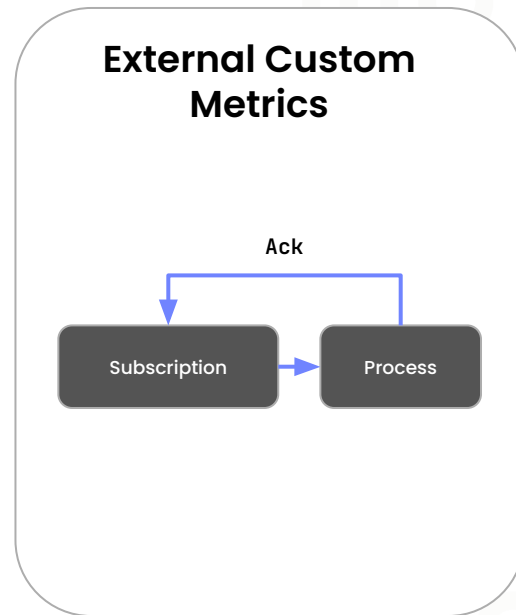**Scale Up:** Increase resources assigned to workload



dojo.

# Autoscaling Triggers

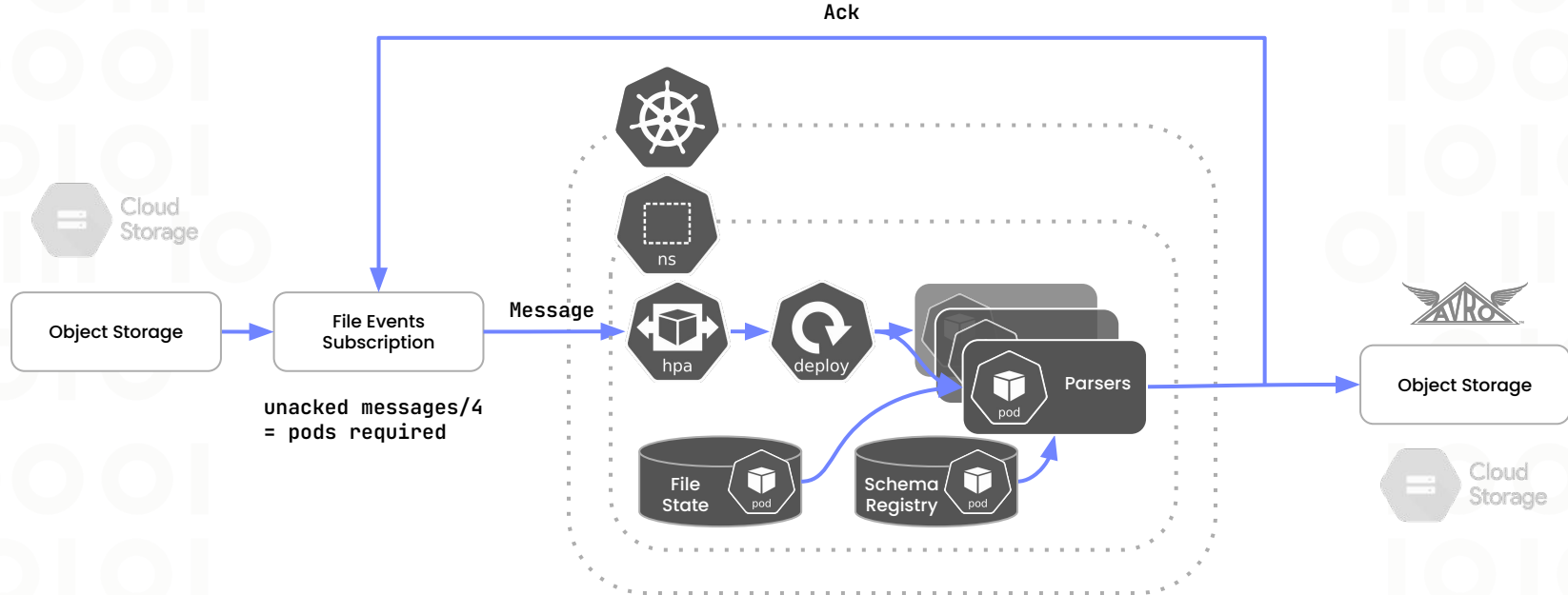Kubernetes provides three solutions depending on the use case



## Resource Usage

CPU        Memory

500m        256Mi

200m        128Mi

metrics.k8s.io

## Custom Metrics

ing

hits-per-second

custom.metrics.k8s.io

## External Custom Metrics

Ack

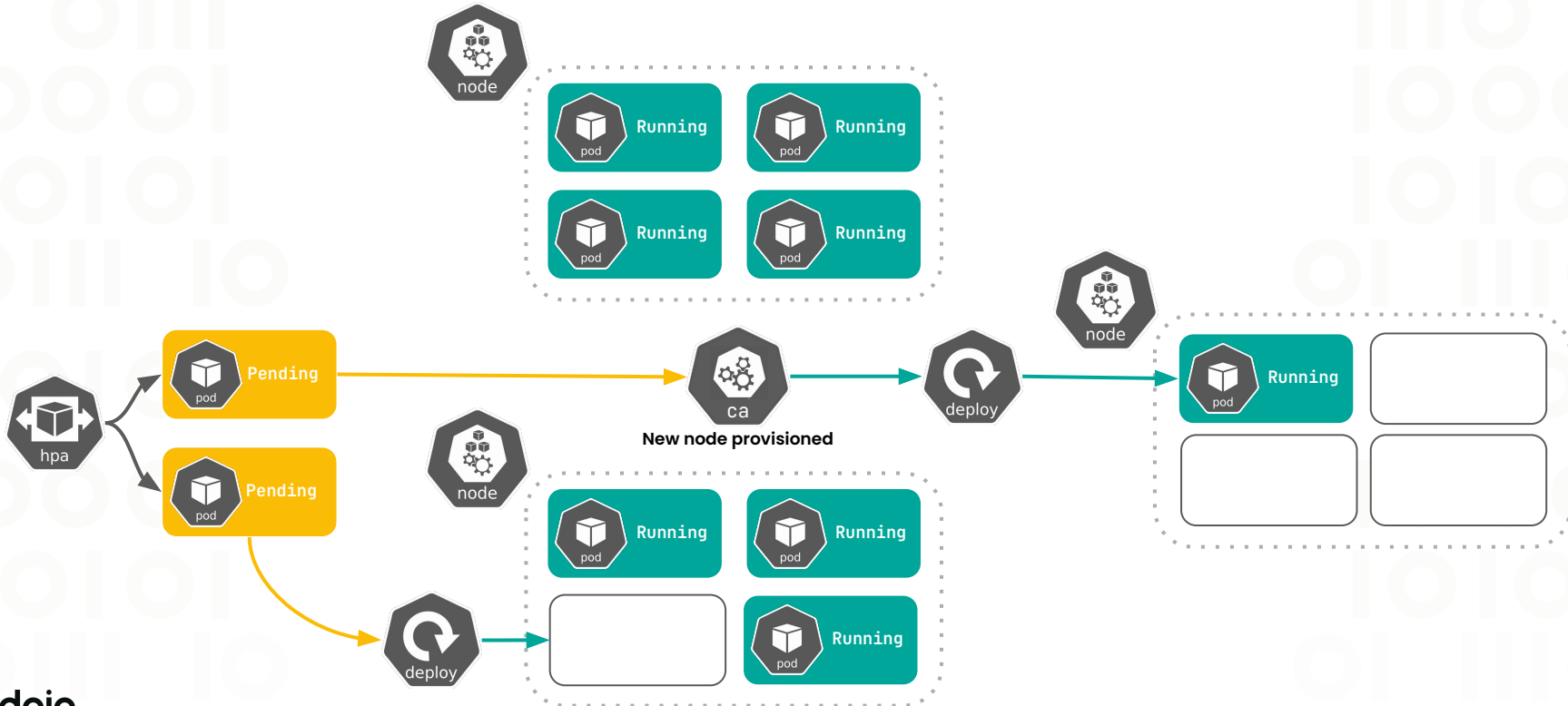Subscription → Process

external.metrics.k8s.io

dojo.

# Using HPA to elastically scale

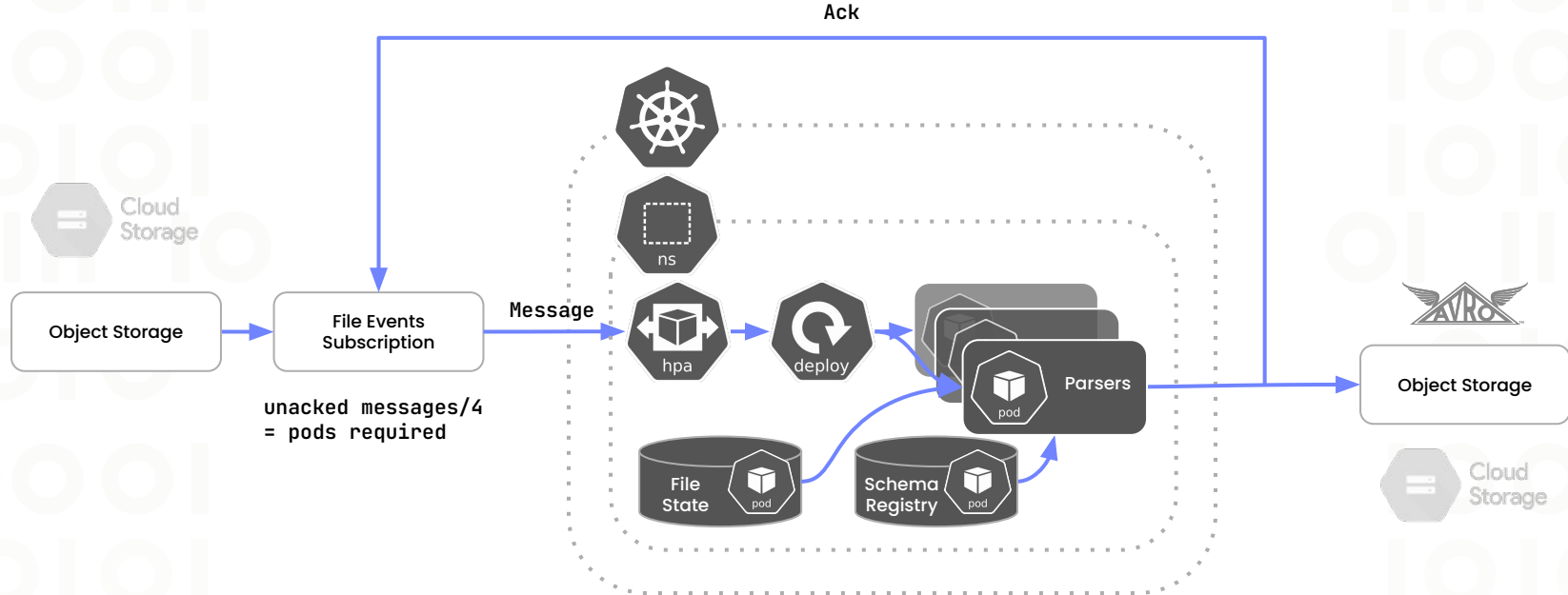Pods required determined by number of unacknowledged messages in queue

# Cluster Autoscaler

New nodes provisioned based on pods in **Pending** state
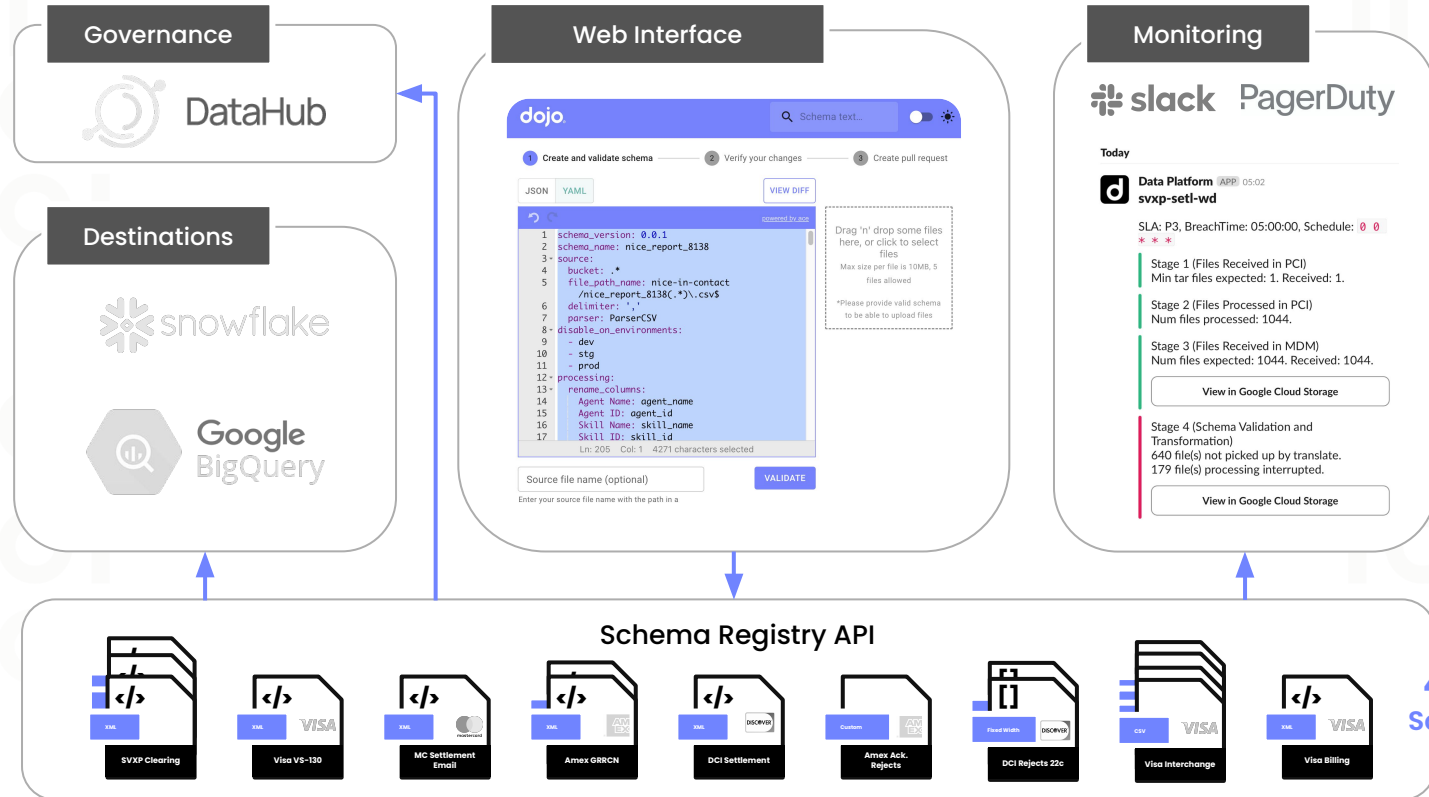


New node provisioned

# Using HPA to elastically scale

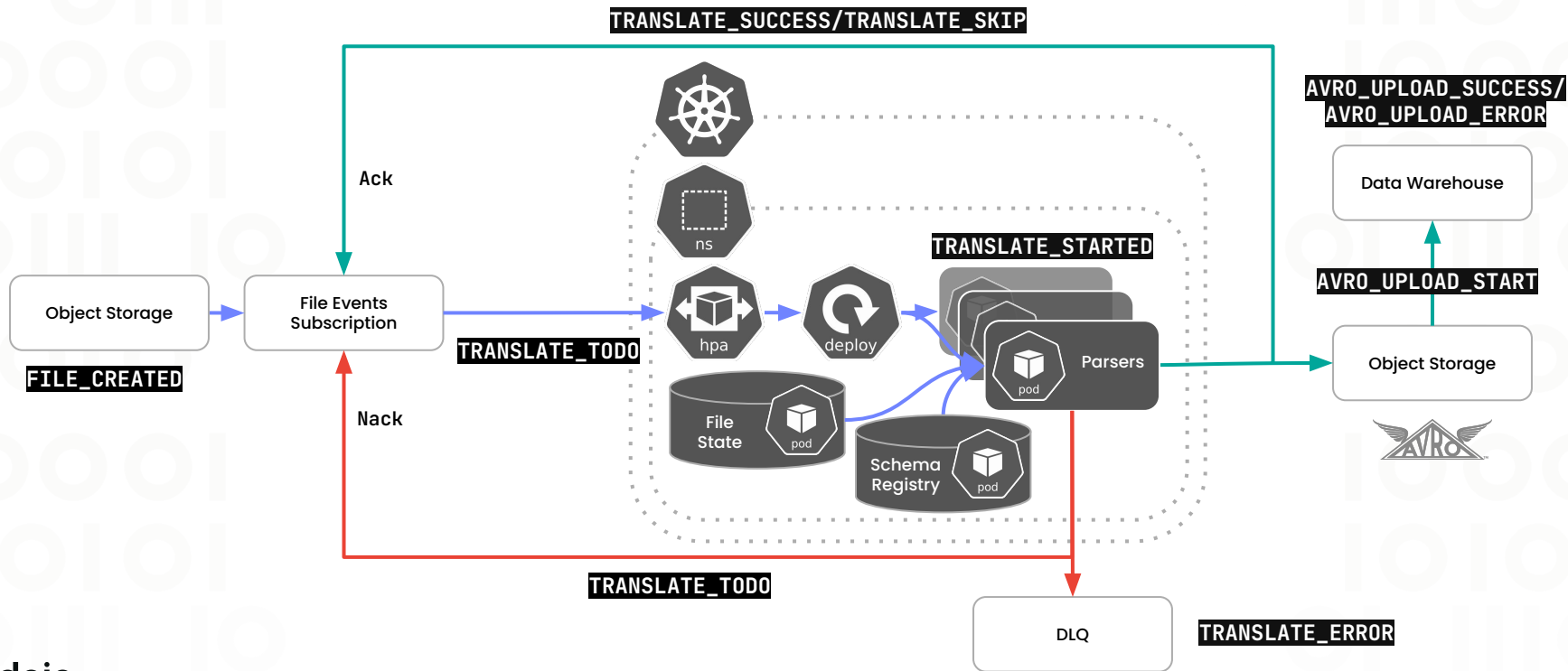Pods required determined by number of unacknowledged messages in queue

# Schema Registry

A central source of metadata for file lifecycle management

# State Management

State of every file recorded throughout the transformation process



TRANSLATE_SUCCESS/TRANSLATE_SKIP

AVRO_UPLOAD_SUCCESS/
AVRO_UPLOAD_ERROR

Data Warehouse

TRANSLATE_STARTED

AVRO_UPLOAD_START

Object Storage

File Events
Subscription

Ack

Nack

TRANSLATE_TODO

ns

hpa

deploy

Parsers

pod

File
State

pod

Schema
Registry

pod

FILE_CREATED

TRANSLATE_TODO

DLQ

TRANSLATE_ERROR

dojo.

# Real time File Monitoring

State of every file recorded throughout the transformation process

# E2E File Monitoring

**PagerDuty**   ❖ **slack**

P0        P1              P2              P3

❖ Looker

---

Today

⬤ **Data Platform** APP 05:02
**svxp-setl-wd**

SLA: P3, BreachTime: 05:00:00, Schedule: 0 0
* * *

Stage 1 (Files Received in PCI)
Min tar files expected: 1. Received: 1.

Stage 2 (Files Processed in PCI)
Num files processed: 1044.

Stage 3 (Files Received in MDM)
Num files expected: 1044. Received: 1044.

| View in Google Cloud Storage |

Stage 4 (Schema Validation and
Transformation)
640 file(s) not picked up by translate.
179 file(s) processing interrupted.

| View in Google Cloud Storage |

Stage 5 (Loaded into BigQuery)
Files (887) are uploaded to BQ successfully.

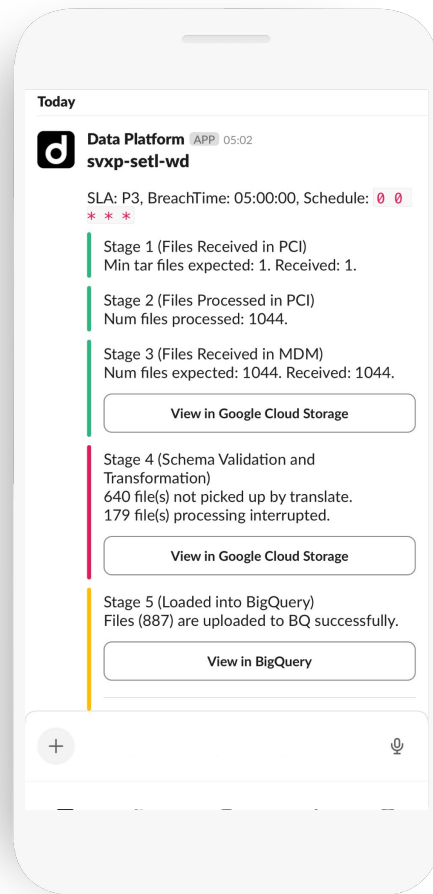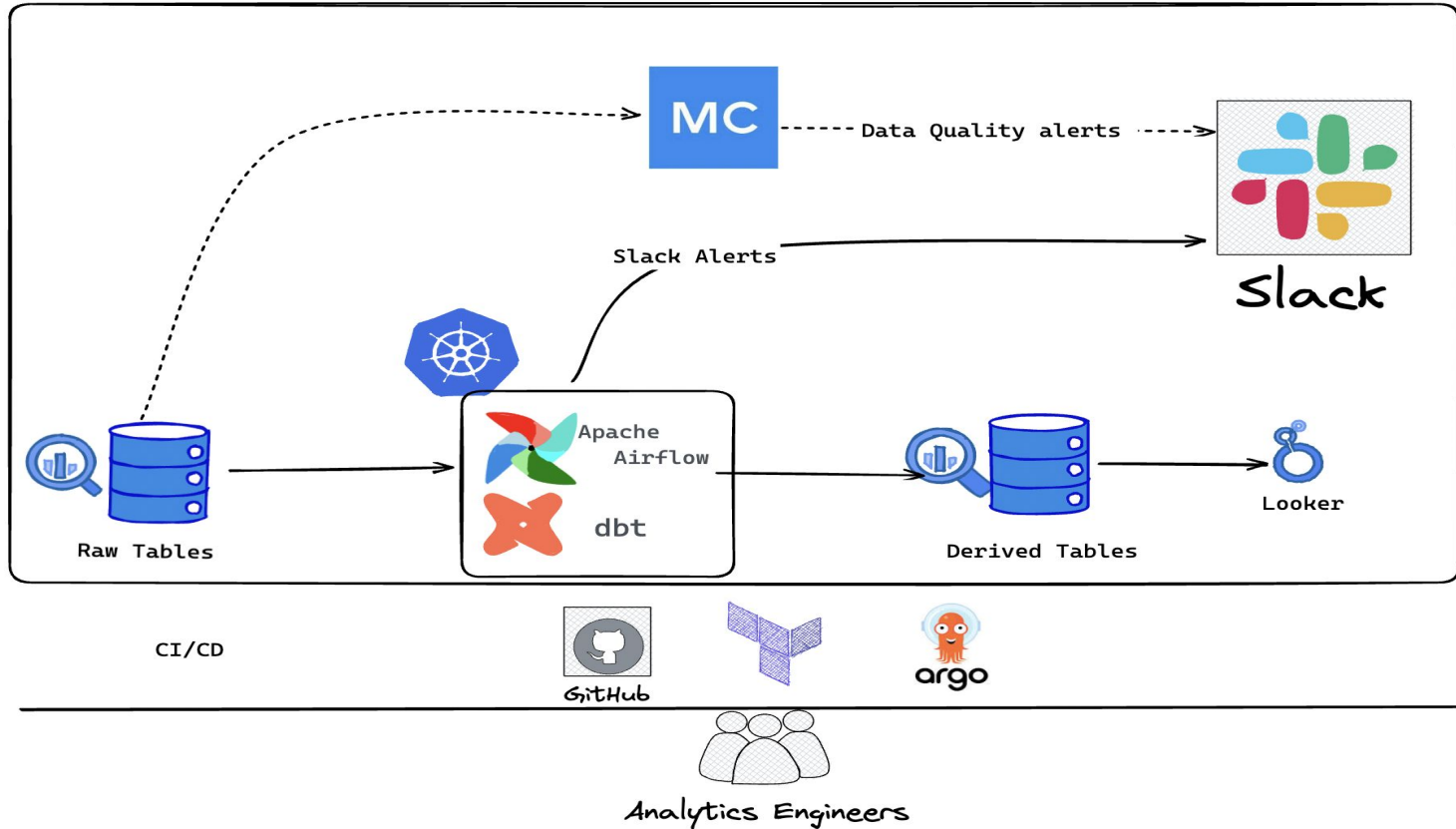| View in BigQuery |

dojo.

# Analytics Platform

# Infrastructure Observability

Monitoring and alerting on Data Platform resources

# Challenges: Embracing Data Mesh

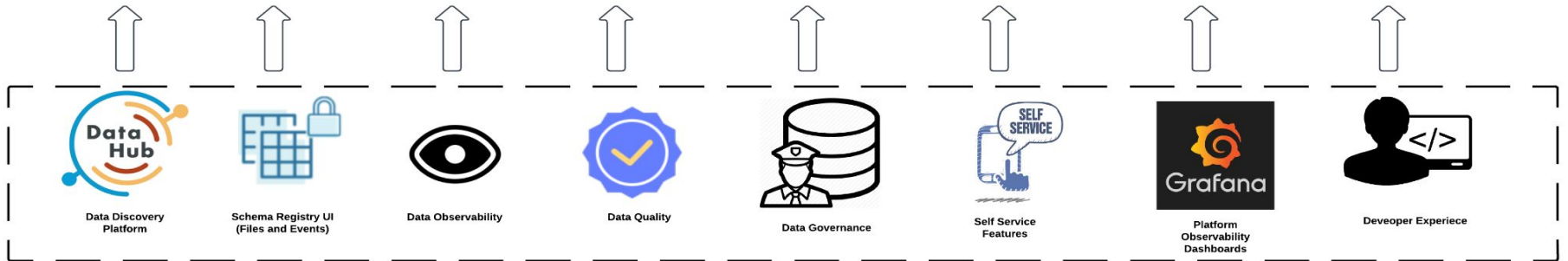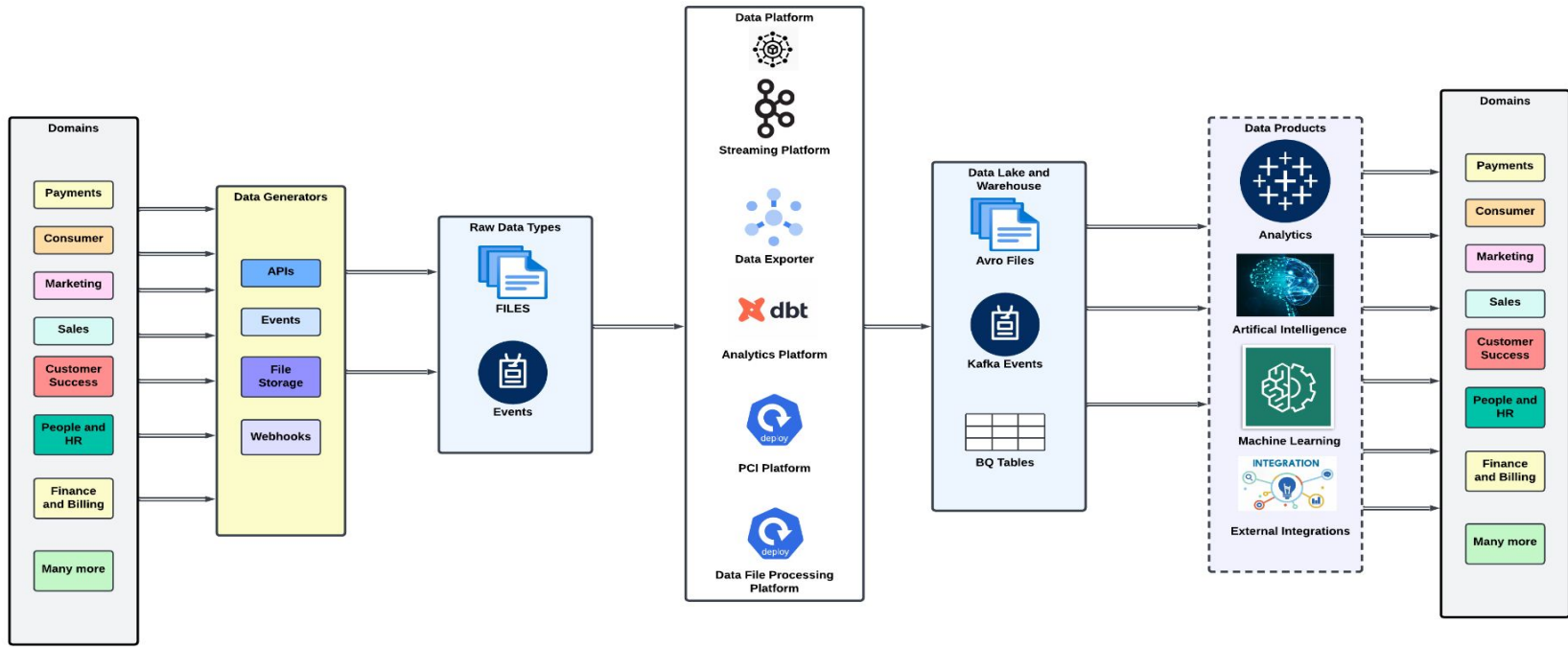| | | |
|---|---|---|
| Cultural Shift | Talent and Skill Gaps | Data Governance Complexity |
| Data Mesh Tooling | Data Privacy and Security | Data Integrations |
| ROI Measurement | Data Cataloging and Data Discovery | Legacy system Integrations |

dojo

# Thank you for your time!

WE'RE HIRING!

https://www.dojo.careers

dojo®