# Scalable Data Ingestion Architecture Using Airflow and Spark

April 17, 2019
Data Council
San Francisco, CA

Johannes Leppä
Data Engineer
johannes.leppa@komodohealth.com

# Agenda

- ❖ Komodo Health

- ❖ Data Ingestion Challenges

- ❖ Data Ingestion System Architecture

- ❖ Lessons Learned and Future Developments

- ❖ Scaling Processes

- ❖ Conclusions

# Komodo Health™ Integrity
# Our Map Links Activities of the Entire Healthcare System

**Payers**
• 500+ payers

**Providers**
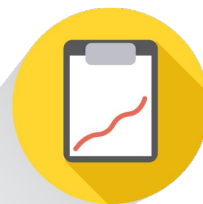• 3.5 M doctors / nurses

**Institutions**
• 450K hospitals / clinics

**Patient-Centric AI powered linkages**

**Biopharma**
• $20B payments

**Clinical Trials**
• 100k+ Clinical Trials

**Scientific Publications**
• 20M publications

komodohealth™

# Agenda

❖ Komodo Health

❖ **Data Ingestion Challenges**

❖ Data Ingestion System Architecture

❖ Lessons Learned and Future Developments

❖ Scaling Processes

❖ Conclusions

komodohealth™

# Variation in data size and cadency
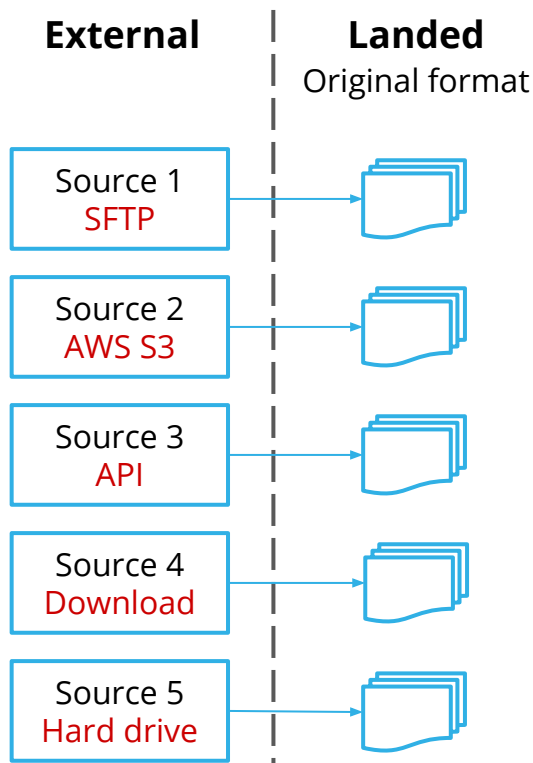
**External**

Source 1

Source 2

Source 3

Source 4

Source 5

- Public and proprietary sources
- Size of data
  - From MBs to TBs
- Refresh cadencies:
  - Daily
  - Weekly
  - Monthly
  - Quarterly
  - Bi-annual
  - One-off
    - Historical drop followed by incremental additions

komodohealth™

# Variation in access to raw data

**External**

**Landed**
Original format

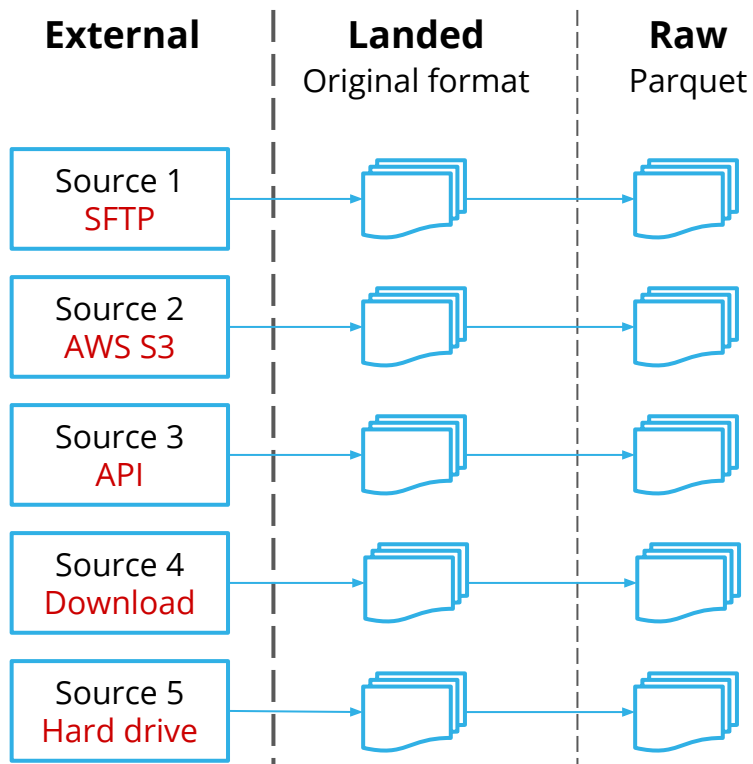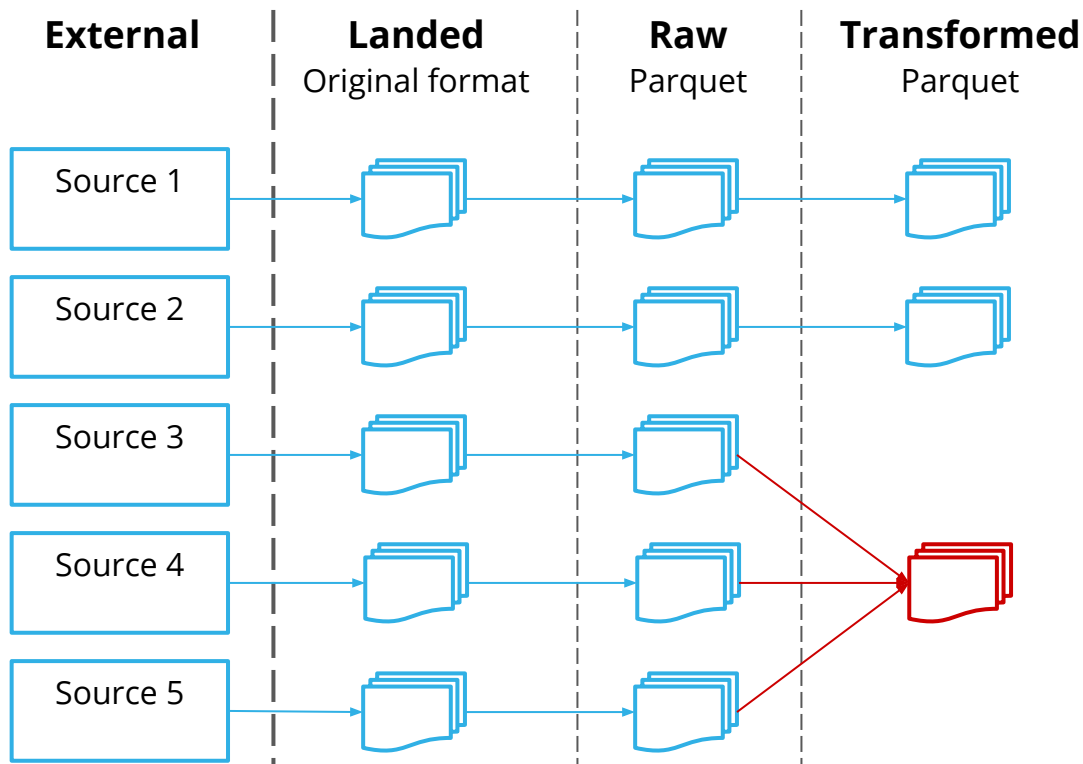| Source 1 SFTP |
| Source 2 AWS S3 |
| Source 3 API |
| Source 4 Download |
| Source 5 Hard drive |

- Public and proprietary sources
- Size of data
  - From MBs to TBs
- Refresh cadencies:
  - Daily
  - Weekly
  - Monthly
  - Quarterly
  - Bi-annual
  - One-off
    - Historical drop followed by incremental additions
- **Several interfaces for data extraction**

komodohealth™

# Variation in file formats

| External | Landed<br>Original format | Raw<br>Parquet |
|---|---|---|
| Source 1<br>SFTP | | |
| Source 2<br>AWS S3 | | |
| Source 3<br>API | | |
| Source 4<br>Download | | |
| Source 5<br>Hard drive | | |

- Original file formats
  - CSV
  - XML
  - SAS
  - Fixed-width
  - Parquet
- Various compression formats
- Encrypted data

komodohealth™

# Cover several aspects of healthcare system

| **External** | **Landed**<br>Original format | **Raw**<br>Parquet | **Transformed**<br>Parquet |
|---|---|---|---|



- Several datasets covering a single aspect of healthcare
  - Different schemas
  - Different conventions
- Need to transform to **common schema**

komodohealth™

# Security and privacy



**External**

**Landed**
Original format

**Raw**
Parquet

**Transformed**
Parquet

Source 1

Source 2

Source 3

Source 4

Source 5

- Security and privacy
  - Access control
  - Data encryption
  - Compliances

komodohealth™

# Prior to centralized data ingestion system

- Eternal question: What is the priority?
  - Scalability, maintainability, robustness, reliability
  - Rapid development

# Prior to centralized data ingestion system

- Eternal question: What is the priority?
  - Scalability, maintainability, robustness, reliability
  - **Rapid development ← startup choice**
    - Provide value to customers and show progress to investors
    - React to changing requirements

# Prior to centralized data ingestion system

- Eternal question: What is the priority?
  - Scalability, maintainability, robustness, reliability
  - **Rapid development ← startup choice**
    - Provide value to customers and show progress to investors
    - React to changing requirements
- Consequences:
  - Specialized pipelines
  - Manual operations
  - Variation in technologies and how to use them
  - Less reusable code

# Why did we build a centralized ingestion system?

- Previous approach hard to maintain
  - Overhead in onboarding engineers to processes
  - Accumulation of manual tasks
- Project to integrate a few new data sources
  - Daily increments
  - Similar data sources
  - **Opportunity**: build system for these sources and migrate other sources later
- Pros of in-house implementation
  - Flexibility
  - Integrate with our tech stack
    - Leverage previous experience

# Agenda

❖ Komodo Health

❖ Data Ingestion Challenges

❖ **Data Ingestion System Architecture**

❖ Lessons Learned and Future Developments

❖ Scaling Processes

❖ Conclusions

komodohealth™

# Overview of the system infrastructure

- Airflow
  - Organize workflows
  - Automation
  - Alerting
- Spark
  - Distributed processing
- Kubernetes
  - Container management
- AWS
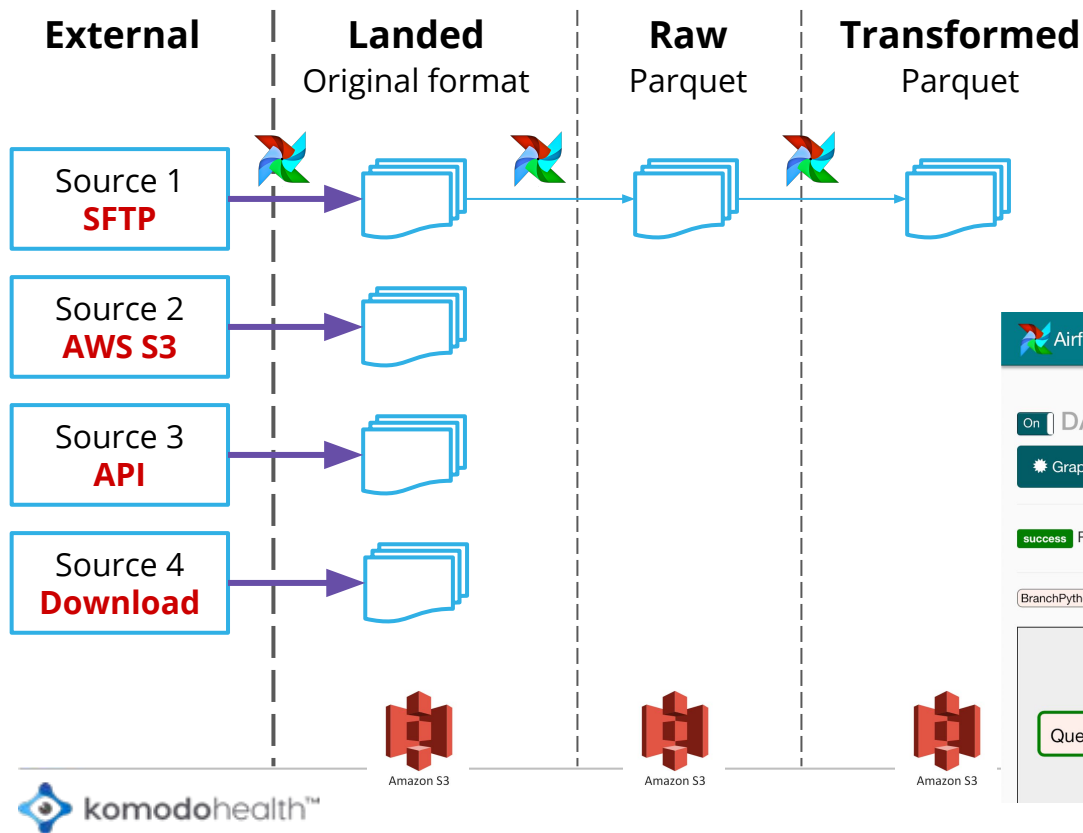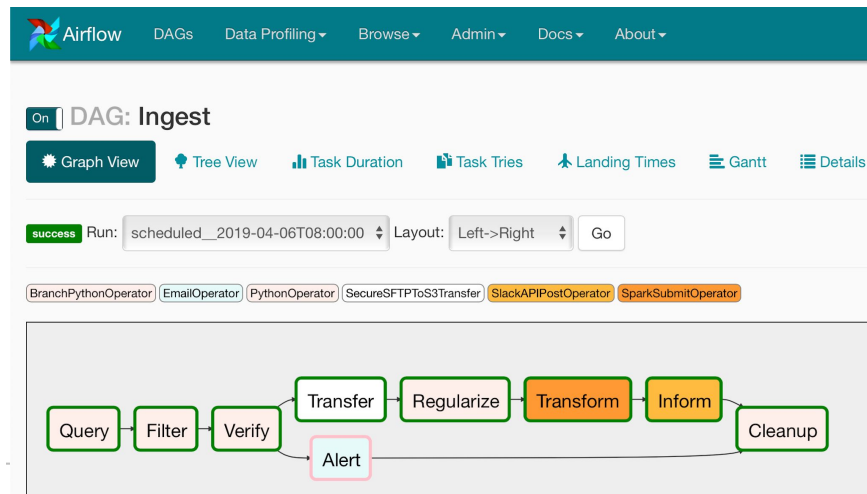  - EC2 - servers
  - S3 - store data

Apache **Airflow**

APACHE **Spark**™

**kubernetes**

Amazon EC2    Amazon S3    **amazon** web services

komodohealth™

# Airflow: Schedule workflows



| External | Landed<br>Original format | Raw<br>Parquet | Transformed<br>Parquet |
|---|---|---|---|
| Source 1<br>**SFTP** | | | |
| Source 2<br>**AWS S3** | | | |
| Source 3<br>**API** | | | |
| Source 4<br>**Download** | | | |

Amazon S3    Amazon S3    Amazon S3

**Pros:**
- DAGs written in Python
- **Hooks** to integrate with sources
- **Operators** for common tasks
- Alert on success/failure
- Monitoring
- Parallelize DAGs and tasks

Airflow   DAGs   Data Profiling ▾   Browse ▾   Admin ▾   Docs ▾   About ▾

On | DAG: Ingest

✻ Graph View    🍀 Tree View    Task Duration    Task Tries    Landing Times    Gantt    Details

success Run: scheduled__2019-04-06T08:00:00 ⇕   Layout: Left->Right ▾   Go

BranchPythonOperator   EmailOperator   PythonOperator   SecureSFTPToS3Transfer   SlackAPIPostOperator   SparkSubmitOperator

Query → Filter → Verify → Transfer → Regularize → Transform → Inform → Cleanup
                              ↓
                           Alert

komodohealth™

# Airflow: Schedule workflows

| External | Landed Original format | Raw Parquet | Transformed Parquet |
|---|---|---|---|

Source 1 **SFTP**

Source 2 **AWS S3**

Source 3 **API**

Source 4 **Download**
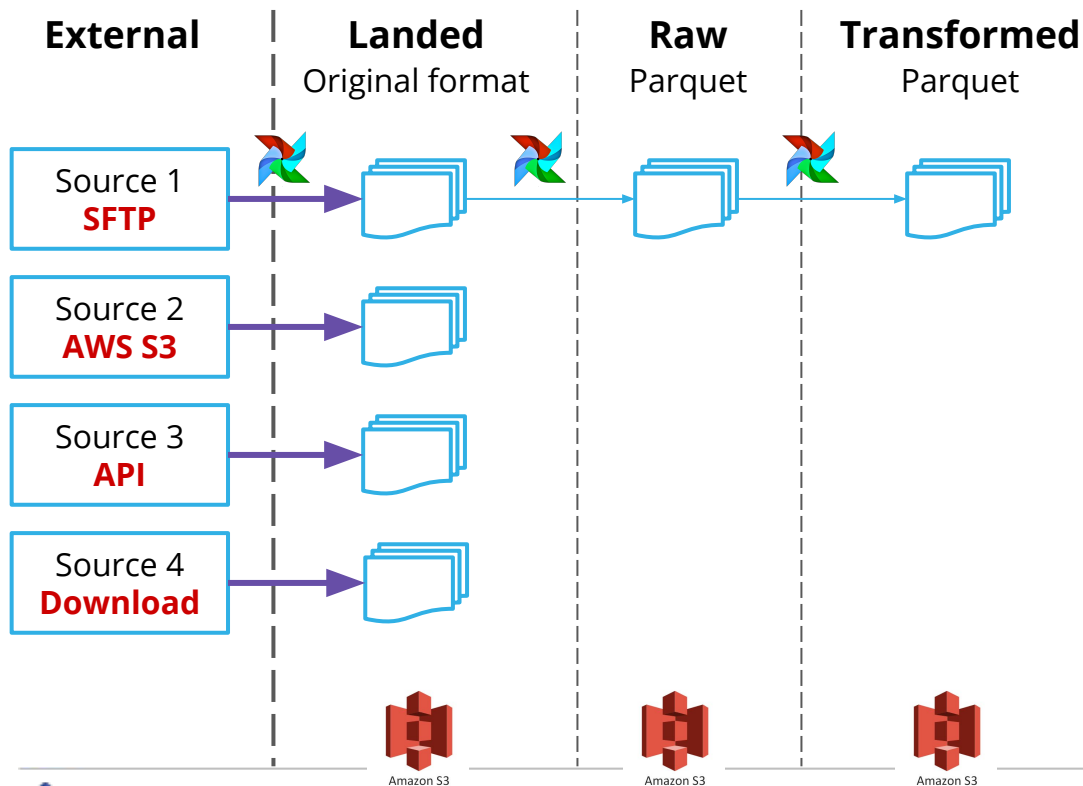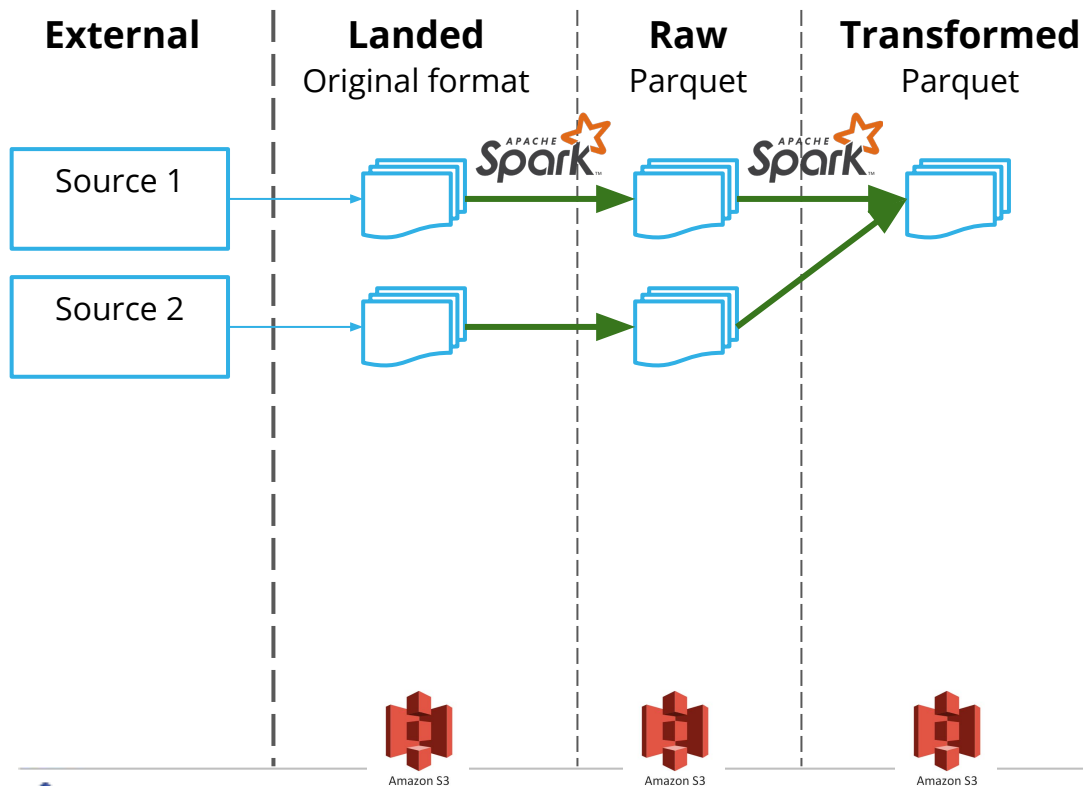
Amazon S3

Amazon S3

Amazon S3

**Pros:**
- DAGs written in Python
- **Hooks** to integrate with sources
- **Operators** for common tasks
- Alert on success/failure
- Monitoring
- Parallelize DAGs and tasks

**Cons:**
- Had to customize hooks and operators
  - Handling credentials
  - Needing additional S3 metadata

komodohealth™

# Spark: Distributed processing

| External | Landed | Raw | Transformed |
|----------|--------|-----|-------------|
| | Original format | Parquet | Parquet |



**Pros:**
- Reliable
- Python and Scala APIs

**Cons:**
- Performance tuning can be tricky

komodohealth™

# Kubernetes: Container management

**Pros:**
- Environments isolated to namespaces
- Node selectors for resource allocation
  - Nodes labeled based on the Auto Scaling Groups instances are tied to
- **Self-healing of pods!**

**Cons:**
- Occasional stability issues
  - Networking issues
- Difficult to troubleshoot

**Spark Master**

Pod

**Airflow Scheduler**

Pod

Node

Amazon EC2

**Airflow WebUI**

Pod

Node

Amazon EC2

kubernetes

komodohealth™

# So far so good

- ☑ Scheduled execution
- ☑ Parallelized tasks
- ☑ Scalable resources
- ☑ Alerting
- ☑ Monitoring
- ☑ Resilient infrastructure
- ☑ Isolated environments

# Agenda

❖ Komodo Health

❖ Data Ingestion Challenges

❖ Data Ingestion System Architecture

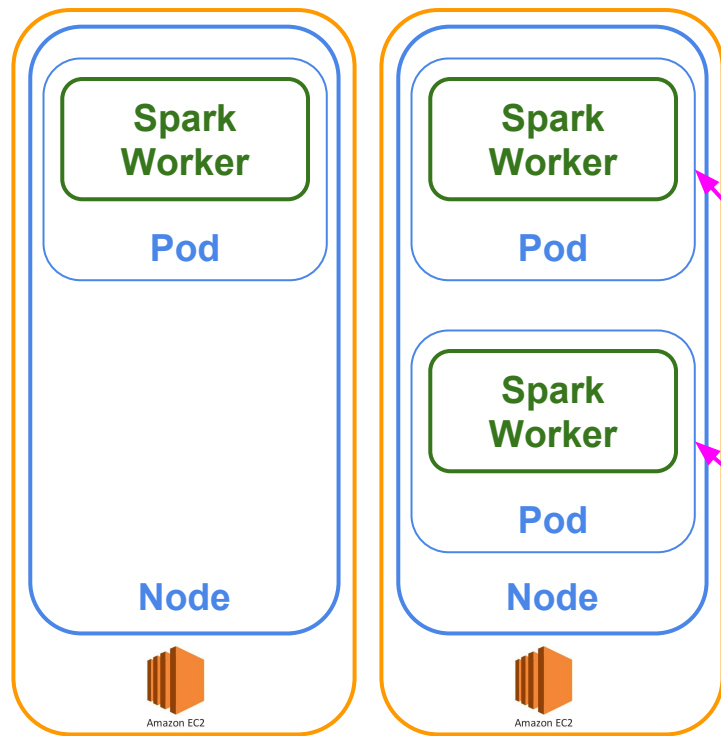❖ **Lessons Learned and Future Developments**

❖ Scaling Processes

❖ Conclusions

komodohealth™

# Infra limitation: Spark scaled manually

**Spark Worker**

Pod

Node

Amazon EC2
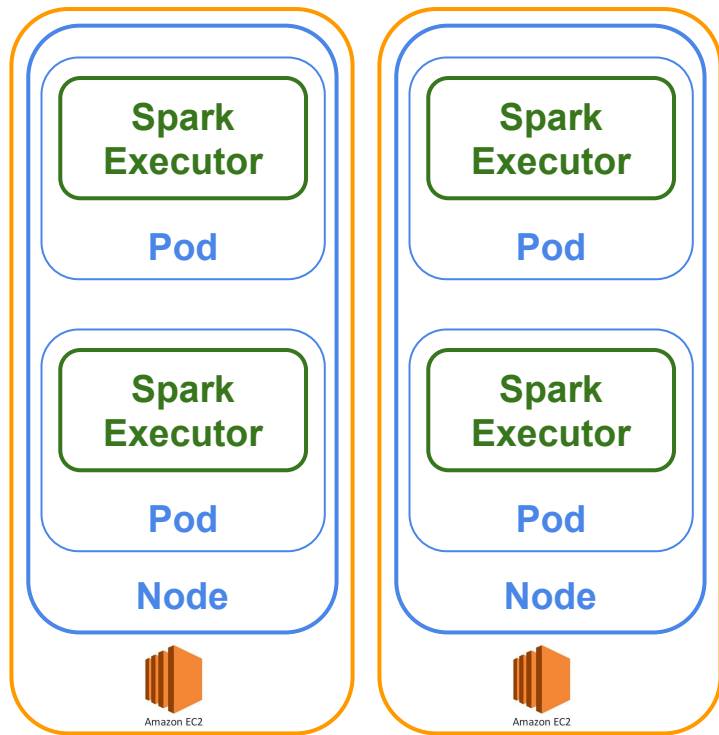
**Spark Worker**

Pod

Node

Amazon EC2

## Big spikes in resource usage

- Wasteful to keep scaled up
- Scaling down is tricky
- Currently run big workloads on separate cluster
  - Manual operation :(

komodohealth™

# Infra limitation: Spark scaled manually

**Spark Worker**

Pod

**Spark Worker**

Pod

**Spark Worker**

Pod

Node

Node

Amazon EC2

Amazon EC2

## Big spikes in resource usage

- Wasteful to keep scaled up
- Scaling down is tricky
- Currently run big workloads on separate cluster
  - Manual operation :(

Two Spark workers on the same node resulted in double counting Spark resources

komodohealth™

# Automatic scaling under development

**Spark Executor**

Pod

**Spark Executor**

Pod

Node

Amazon EC2

**Spark Executor**
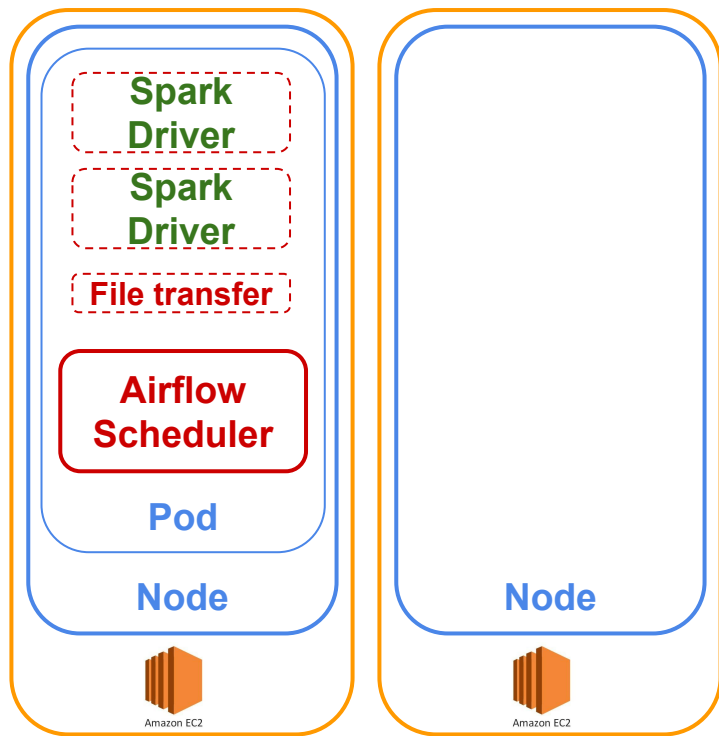
Pod

**Spark Executor**

Pod

Node

Amazon EC2

## Big spikes in resource usage

- Wasteful to keep scaled up
- Scaling down is tricky
- Currently run big workloads on separate cluster
    - Manual operation :(

## Future solution:

- Run Spark directly on Kubernetes
    - Introduced in Spark 2.4.0 for client mode
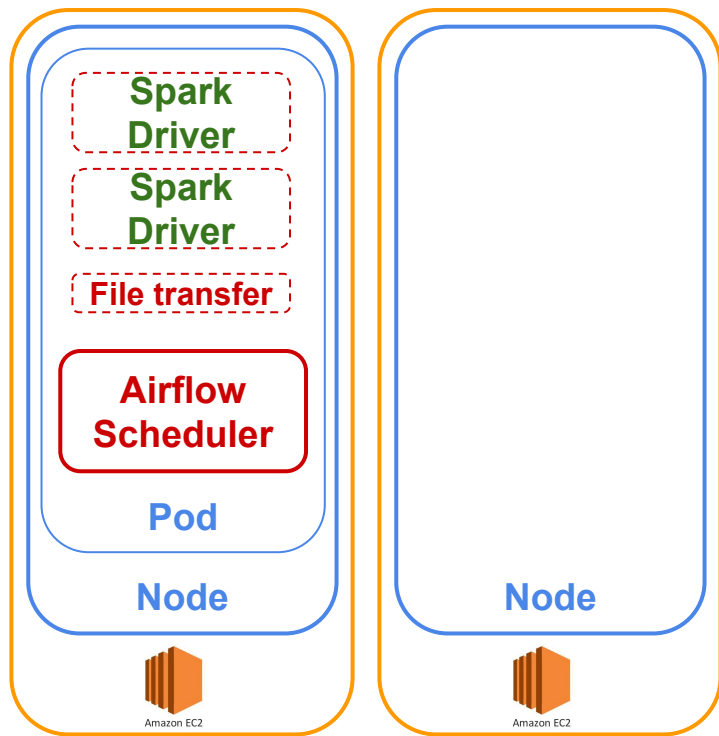- K8s autoscaler to scale nodes

komodohealth™

# Infra limitation: Scheduler a single point of failure

**Spark Driver**

**Spark Driver**

**File transfer**

**Airflow Scheduler**

**Pod**

**Node**

Amazon EC2

**Node**

Amazon EC2

## Using local executor

- Tasks executed as subprocesses of scheduler
- Scale resources vertically
- Self-healing on failures? It depends…

komodo*health*™

# Infra limitation: Scheduler a single point of failure

**Spark Driver**

**Spark Driver**

**File transfer**

**Airflow Scheduler**

**Pod**
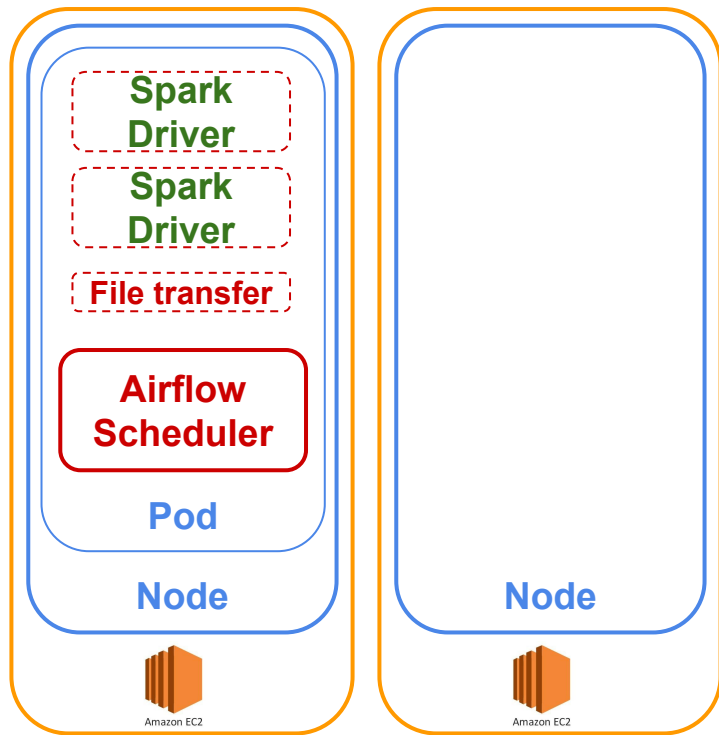
**Node**

Amazon EC2

**Node**

Amazon EC2

## Using local executor

- Tasks executed as subprocesses of scheduler
- Scale resources vertically
- Self-healing on failures? It depends...

## Issues in self-healing:

- Inconsistency in Airflow database
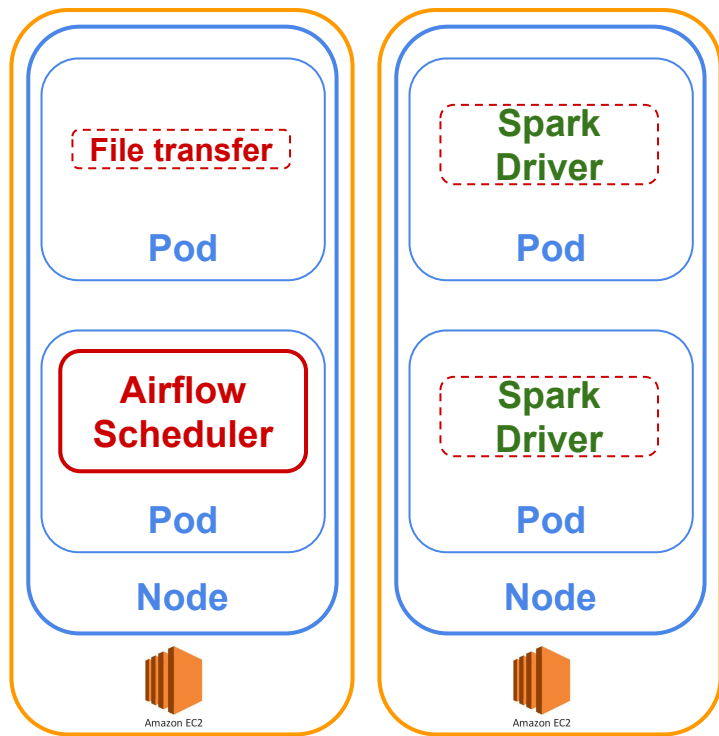- Dependency on lost local file
- Pod evicted due to disk pressure

komodohealth™

# Why are you using local executor?

**Spark Driver**

**Spark Driver**

**File transfer**

**Airflow Scheduler**

**Pod**

**Node**

Amazon EC2

**Node**

Amazon EC2

## It has served us well, so far

- It was enough when we started
- Did not want to add complexity

komodohealth™

# Automatic scaling under development, again



**It has served us well, so far**

- It was enough when we started
- Did not want to add complexity

**Future solution:**

- Kubernetes executor
  - Introduced in Airflow 1.10.0
- K8s autoscaler to scale nodes

File transfer

Pod

Airflow Scheduler

Pod

Node

Amazon EC2

Spark Driver

Pod
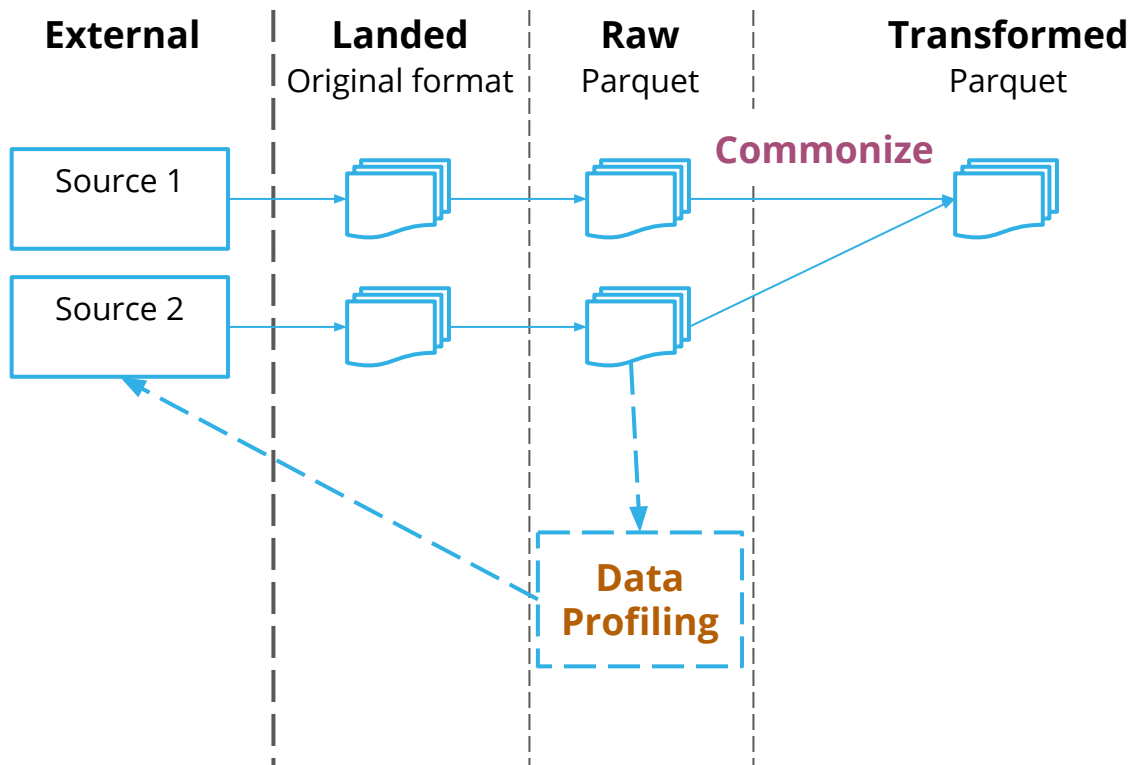
Spark Driver

Pod

Node

Amazon EC2

komodohealth™

# Agenda

- ❖ Komodo Health

- ❖ Data Ingestion Challenges

- ❖ Data Ingestion System Architecture

- ❖ Lessons Learned and Future Developments

- ❖ **Scaling Processes**

- ❖ Conclusions

komodohealth™

# Beyond infra - Scaling the ingestion processes

- Our data ingestion priorities:
  - Speed of data delivery
  - Data quality
  - Security and privacy

- Bottleneck is engineering time spent on integrating new data sources
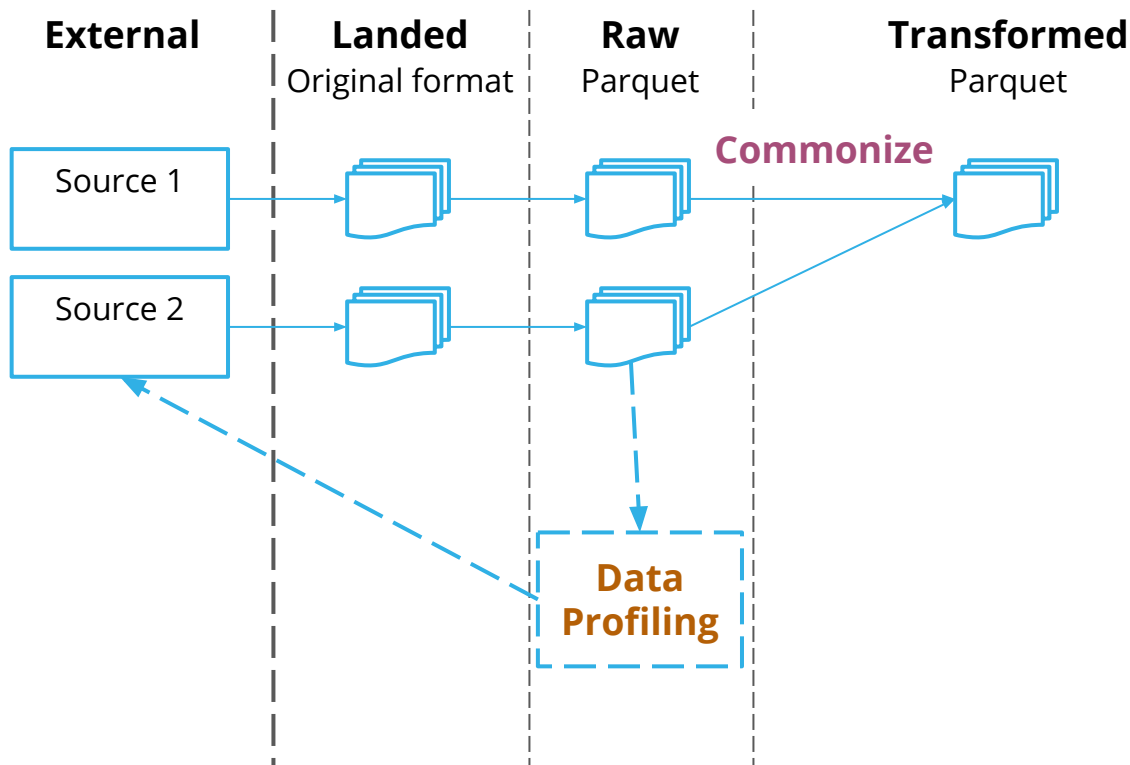  - Tools to simplify processes

# Early and fast iterations

| External | Landed<br>Original format | Raw<br>Parquet | Transformed<br>Parquet |
|----------|---------------------------|----------------|------------------------|



Commonize

Source 1

Source 2

Data Profiling

## Data profiling tool:

- Recognize columns
  - Simplifies commonization
- Validate raw data
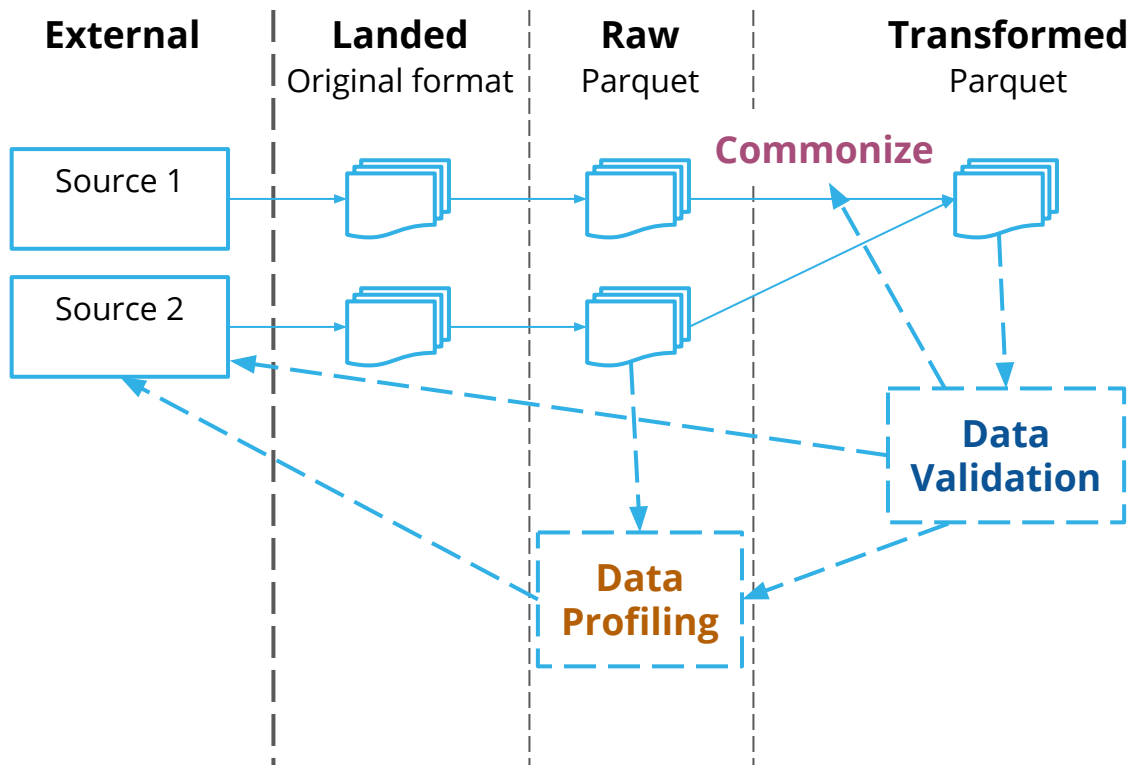  - Communicate issues with source
  - Compliance risks

komodohealth™

# Avoid repeated work

**External** | **Landed** | **Raw** | **Transformed**
Original format | Parquet | Parquet

Source 1

Source 2

**Commonize**

Data
Profiling

**Commonization tool:**

- Similar data to common schema
- Based on configuration file
  - Very little code needed

komodohealth™

# Emphasis on data quality



**External** | **Landed** Original format | **Raw** Parquet | **Transformed** Parquet

Source 1

Source 2

Commonize

Data Validation

Data Profiling

## Data validation tool:

- Validate against data standard
  - Catch bugs in commonization
  - Improve data profiling
  - Communicate issues with source

komodohealth™

# Conclusions

- ❖ Architecture with Airflow, Spark and Kubernetes very flexible for complex data ingestion

- ❖ Lots of nuances with these technologies and their interactions

- ❖ These technologies are constantly improving

- ❖ Not just infra that needs to scale, but also the processes

- ❖ Make sure you know your specific priorities

komodohealth™

# Thank you for your attention!

❖ Architecture with Airflow, Spark and Kubernetes very flexible for complex data ingestion

❖ Lots of nuances with these technologies and their interactions

❖ These technologies are constantly improving

❖ Not just infra that needs to scale, but also the processes

❖ Make sure you know your specific priorities