

# Building Data Pipelines with Spark and StreamSets



StreamSets

Pat Patterson  
Community Champion  
@metadaddy  
pat@streamsets.com

# Agenda

---



Data Drift

StreamSets Data Collector

Running Pipelines on Spark Today

Future Spark Integration

Demo

A solid blue horizontal bar spanning the width of the slide.

# The Evolution of Data-in-Motion



Past



Emerging



Data Sources

Data Stores

Data Consumers

# Data Drift - a Data Engineering Headache

---



The unpredictable, unannounced and unending mutation of data characteristics caused by the operation, maintenance and modernization of the systems that produce the data

Structure  
Drift

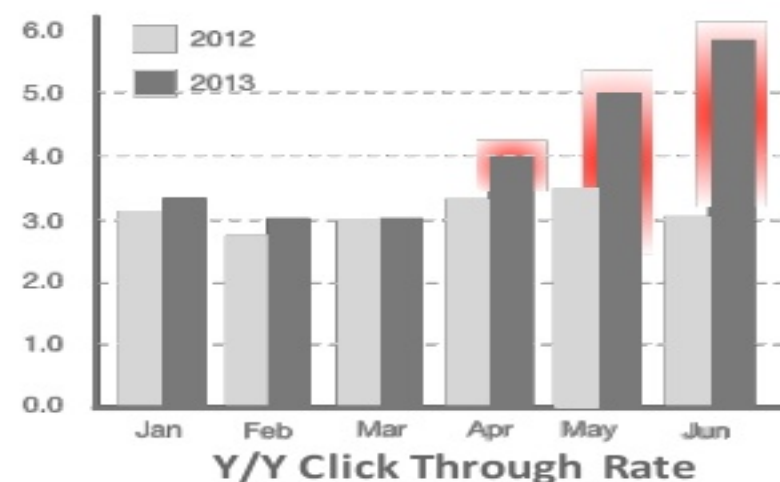
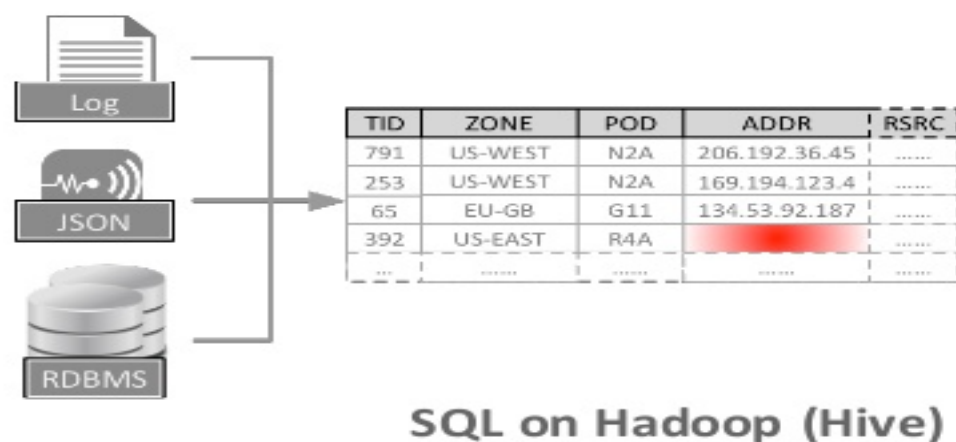
Semantic  
Drift

Infrastructure  
Drift





# Example: Data Loss and Corrosion



80% of analyst time is spent preparing and validating data,  
while the remaining 20% is actual data analysis

# Solving Data Drift



## Data Sources



// DIY Custom Code



## Data Stores



## Data Consumers



Data Drift

Custom code

Fixed-schema

Poor Data Quality

Tools

Applications

Delayed and  
False Insights

# Solving Data Drift



## Data Sources



## Data Stores



## Data Consumers



Data Drift

Intent-Driven

Drift-Handling

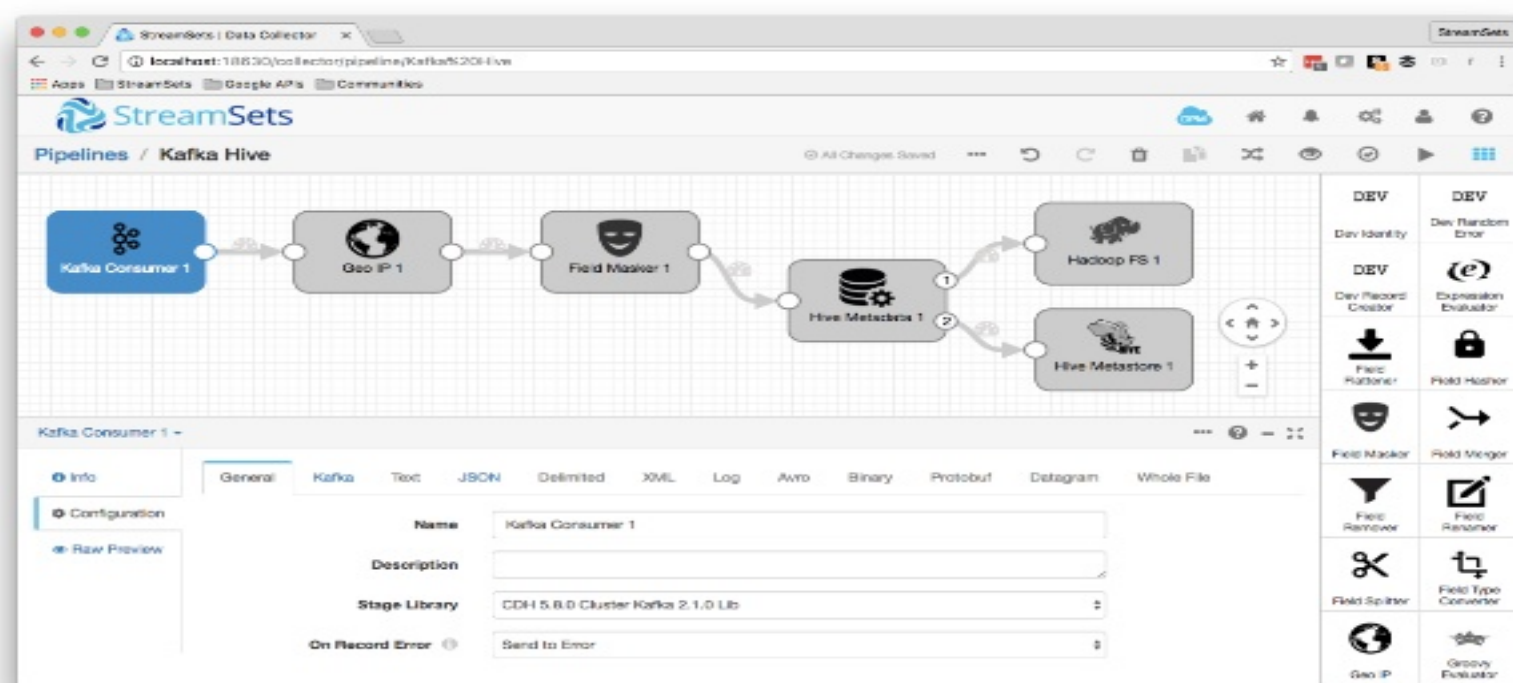
Data KPIs

Tools

Applications

Trusted Insights

# StreamSets Data Collector



Open source software for the rapid development and reliable operation of complex data flows.

- Intent-driven
- UI Abstraction
- Extensible



# Handling Drift with Hive



- Monitor data structure
- Detect schema change
- Alter Hive Metadata

# Running Pipelines on Spark Today



```
spark-submit  
--num-executors ...  
--archives ...  
--files ...  
--jar ...  
--class ...
```



- Container on Spark
- Leverage Kafka RDD
- Scale out for performance

# SDC on Spark - Connectivity

---



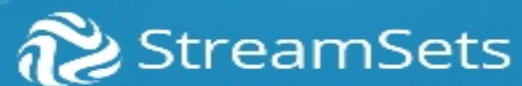
## Sources

- Kafka

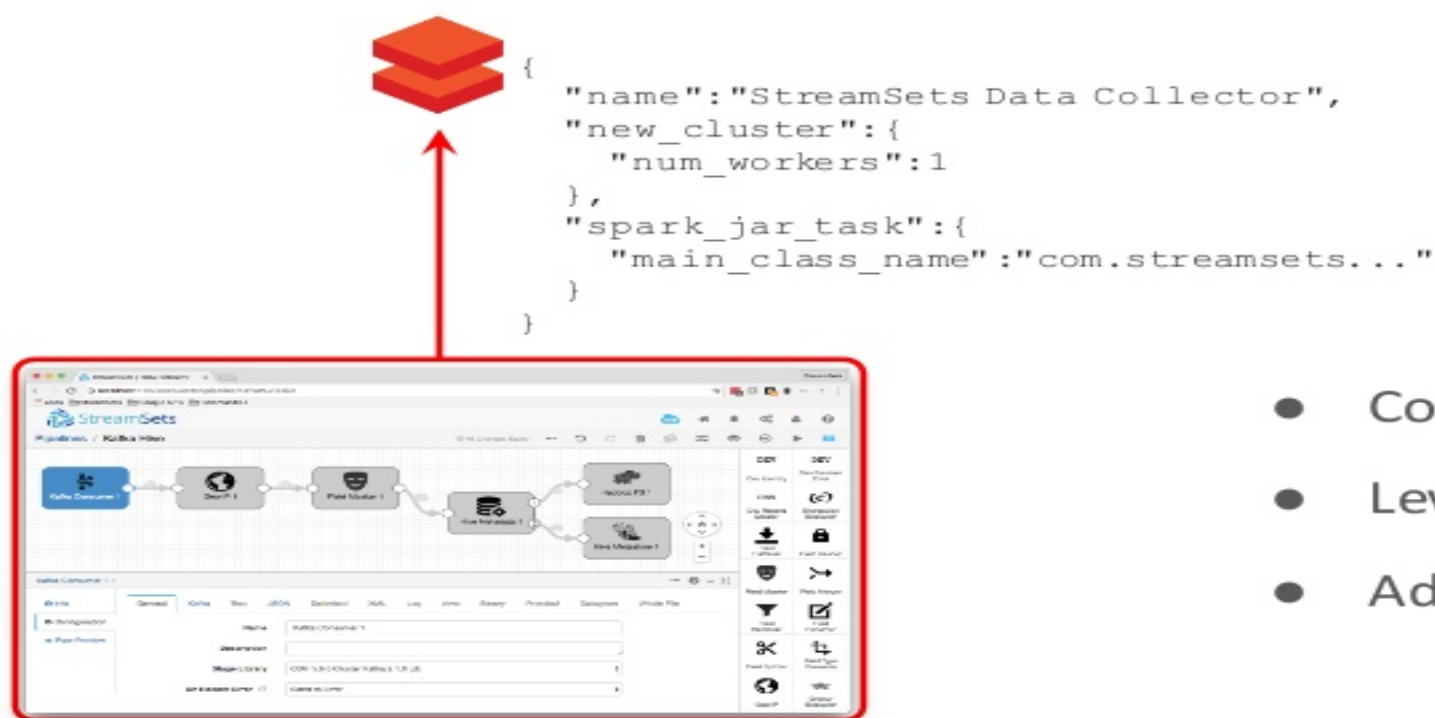
## Destinations

- HDFS
- HBase
- S3
- Kudu
- MapR DB
- Cassandra
- ElasticSearch
- Kafka
- MapR Streams
- Kinesis
- etc, etc, etc!

# Future Directions



# Run Pipelines on Databricks



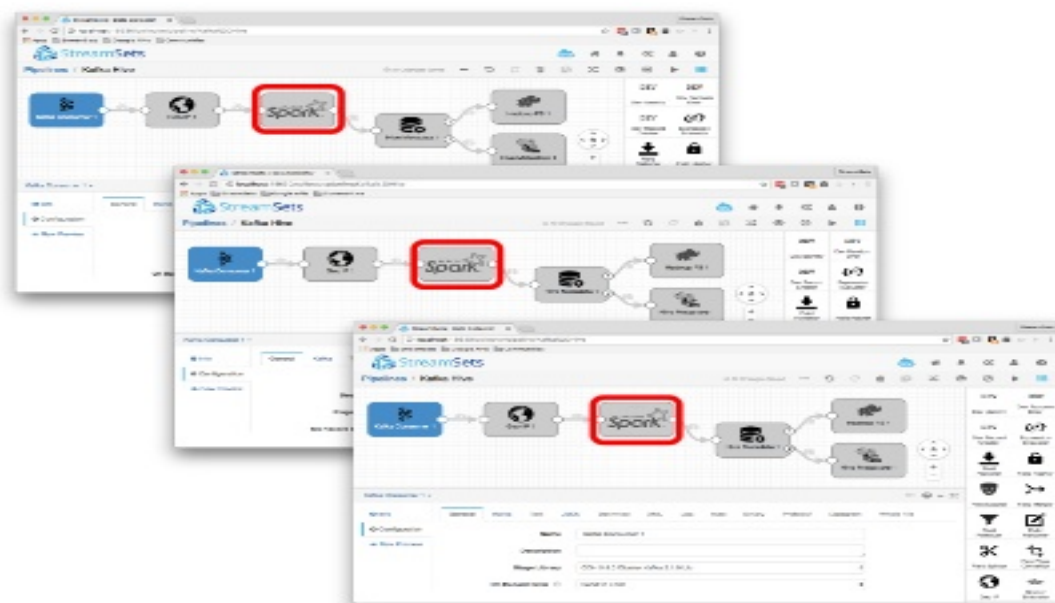
- Container on Databricks
- Leverage REST API
- Add S3 origin



# Break Out Spark Processor



## APACHE **Local Mode**



- Standalone containers, Spark processor
- Leverage Spark code
- Custom RDD
- Start local Spark job for each batch
- Example use cases: *running* image classification, sentiment analysis

# Spark Processor - Connectivity

---



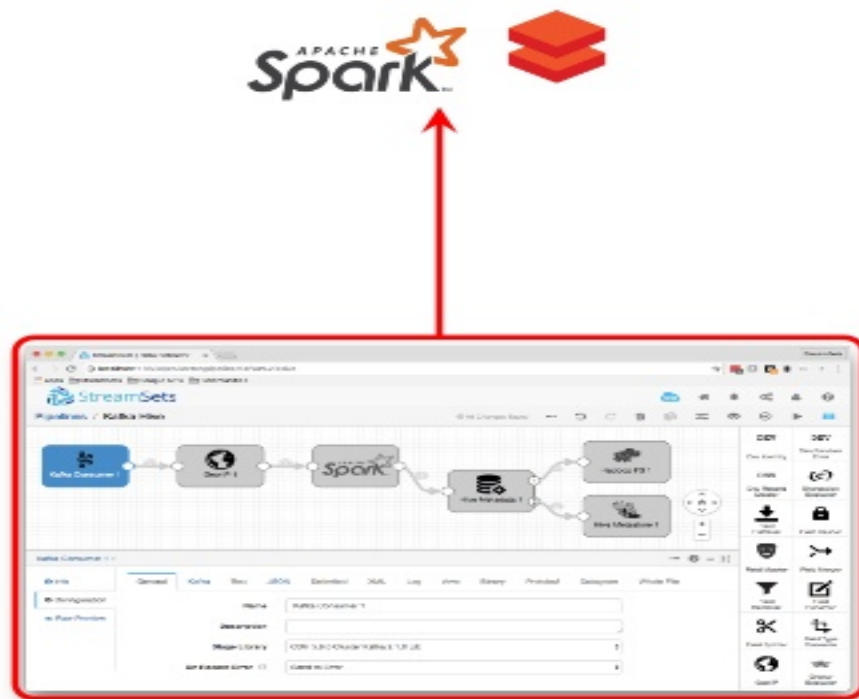
## Sources

- Kafka
- S3
- MapR Streams
- JDBC
- MongoDB
- Local Filesystem
- Redis
- JMS
- HTTP
- UDP
- etc, etc, etc!

## Destinations

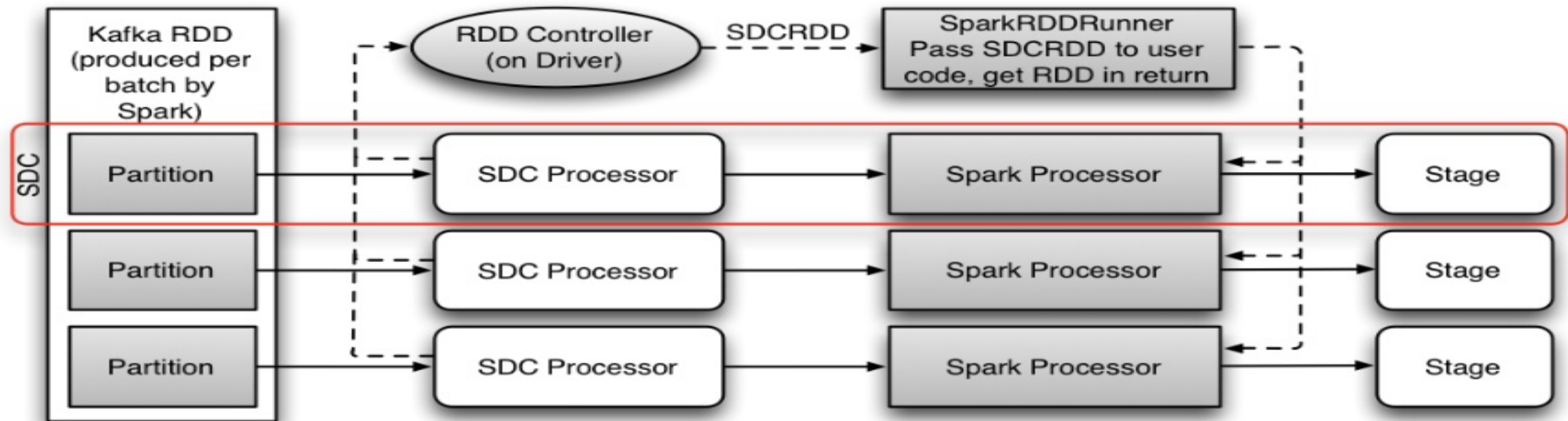
- HDFS
- HBase
- S3
- Kudu
- MapR DB
- Cassandra
- ElasticSearch
- Kafka
- MapR Streams
- JDBC
- etc, etc, etc!

# Deepen Spark Integration



- Container on Spark, Spark processor
- Leverage Spark code
- Custom RDD
- Start Spark job 'on cluster' for each pipeline
- Example use cases: *training* image classification, sentiment analysis

# Spark Integration Architecture



# SDC on Spark - Connectivity Tomorrow

---



## Sources


- Kafka
- S3
- *MapR Streams*
- *JDBC*
- *MongoDB*
- *Redis*
- *JMS*
- *HTTP*
- *UDP*
- ...any partitionable data source...

## Destinations

- HDFS
- HBase
- S3
- Kudu
- MapR DB
- Cassandra
- ElasticSearch
- Kafka
- MapR Streams
- JDBC
- etc, etc, etc!



# Demo

 StreamSets

## Conclusion

---



StreamSets Data Collector brings a UI abstraction to Spark

Standalone container + local Spark Processor bring wide connectivity to Spark code

Spark Container + Spark Processor allow iterative Spark code in pipelines



# Resources

---



**Download StreamSets Data Collector**

<https://streamsets.com/opensource>

**Contribute Code**

<https://github.com/streamsets/datacollector>

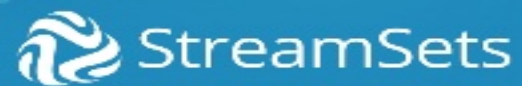
**Get Involved**

<https://streamsets.com/community>

A solid blue horizontal bar spanning the width of the page.



Thank You!



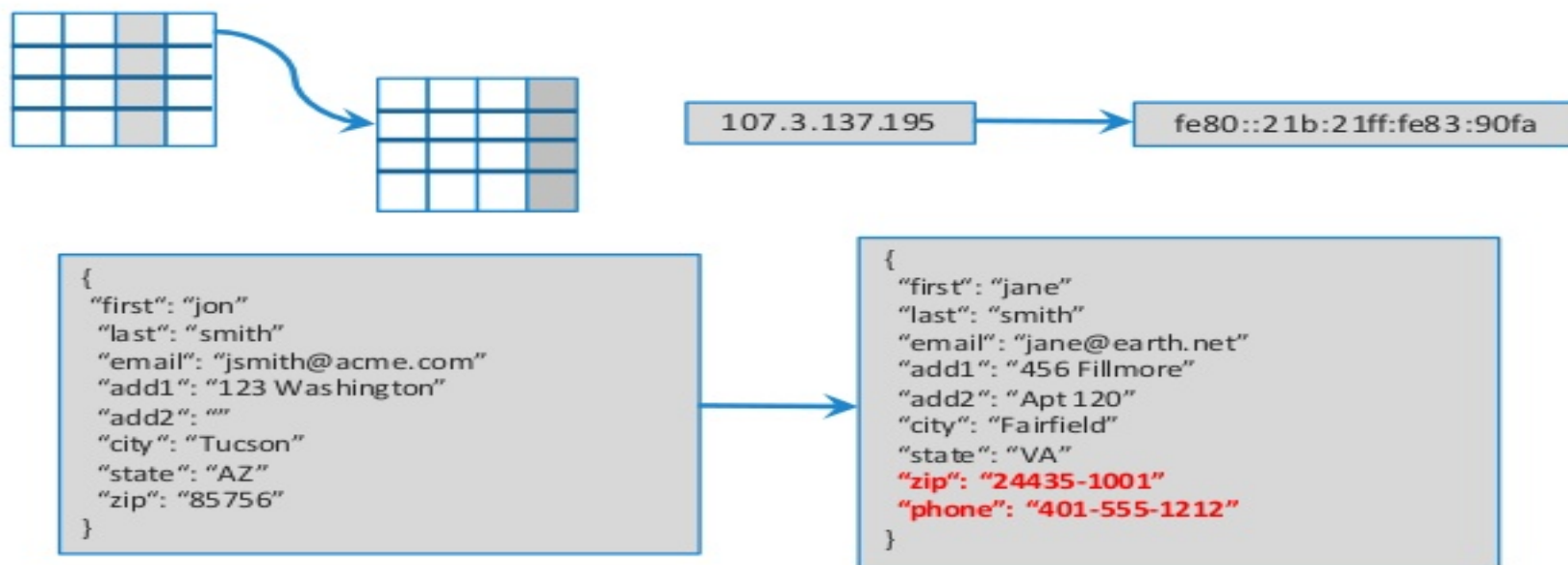


# Backup Slides





# Structure Drift



Data Structure Evolution

## Structure Drift

Data structures and formats evolve and change unexpectedly

**Implication:**  
Data Loss  
Data Squandering

# Semantic Drift



24122-52172 → 00-24122-52172

Account Number Expansion

M134: user {jsmith} read access granted {ac:24122-52172}

M134: user {jsmith} read access granted {ca.ac:24122-52172}

Namespace Qualification

```
.....  
...,3588310669797950,$91.41,jcb,K1088-W#9,...  
...,6759006011936944,$155.04,switch,A6504-Y#9,...  
...,677111111151415,$37.78,laser,Q9936-T#9,...  
...,3585905063294299,$164.48,jcb,S4643-H#9,...  
...,5363527828638736,$117.52,mastercard,X3286-P#9,...  
...,4903080150282806,$168.03,switch,I9133-W#3,...  
.....
```

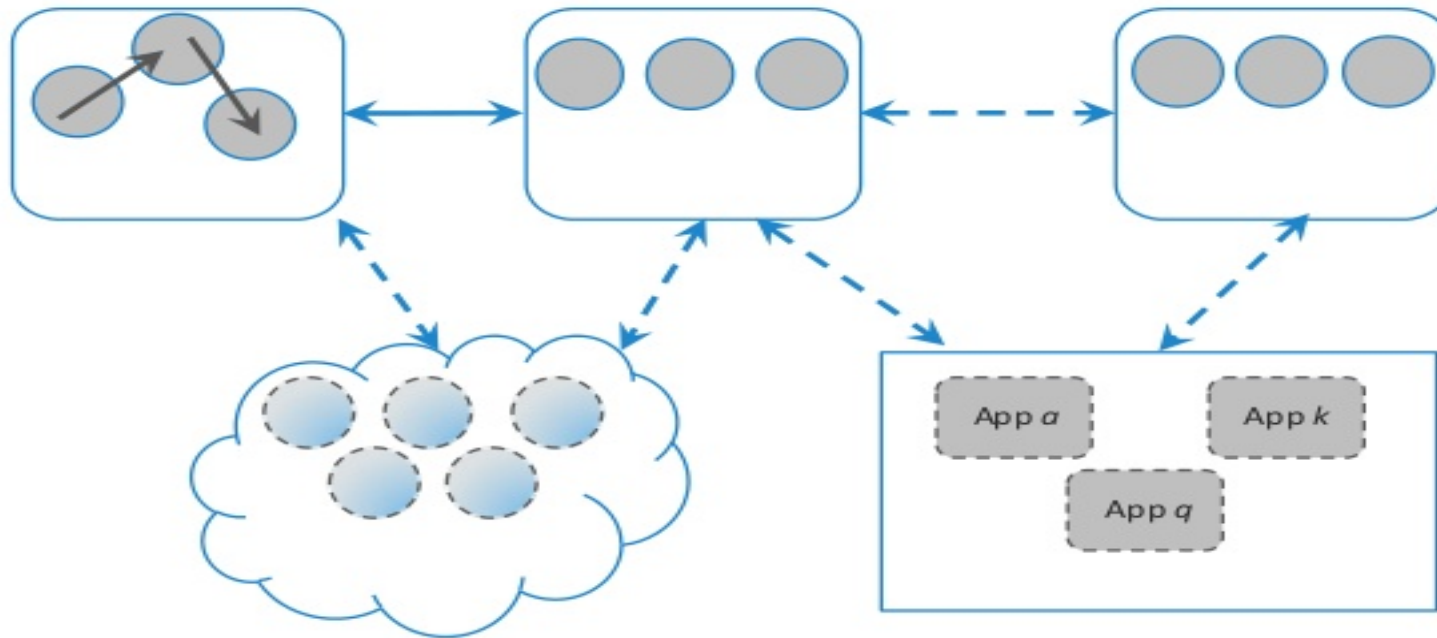
Outlier / Anomaly Detection

## Semantic Drift

Data semantics change with evolving applications

**Implication:**  
Data Corrosion  
Data Loss

# Infrastructure Drift



## Infrastructure Drift

Physical and Logical Infrastructure changes rapidly

**Implication:**  
Poor Agility  
Operational Downtime