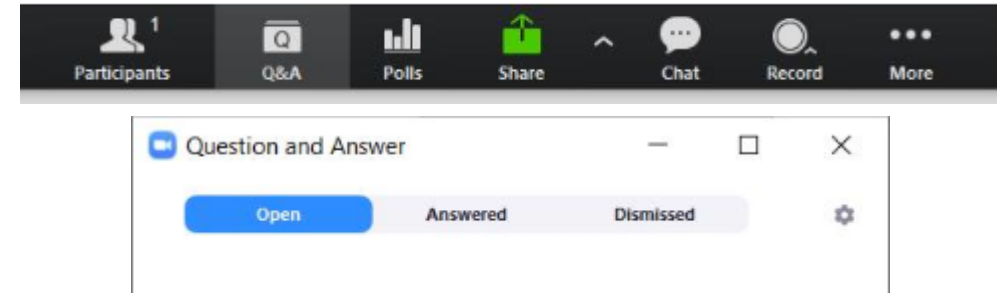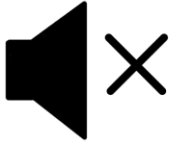# An Introduction to Data Engineering Streaming

# (*AKA* Big Data Streaming)

Ramesh Jha

Informatica Global Customer Support

**Informatica**™

# Housekeeping Tips

➢ Todays Webinar is scheduled to last 1 hour including Q&A

➢ All dial-in participants will be muted to enable the speakers to present without interruption

➢ Questions can be submitted to "All Panelists"  via the Q&A option and we will respond at the end of the presentation

➢ The webinar is being recorded and will be available to view on our INFASupport YouTube channel and Success Portal. The link will be emailed as well.

➢ Please take time to complete the post-webinar survey and provide your feedback and suggestions for upcoming topics.
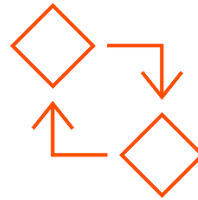
# Success Portal
## https://success.informatica.com
## Learn. Adopt. Succeed.

Bootstrap product trial experience

Enriched Onboarding experience

FREE Product Learning Paths and weekly Expert sessions

Informatica Concierge with Chatbot integrations

Tailored training and content recommendations

Informatica®

# Safe Harbor

The information being provided today is for informational purposes only. The development, release, and timing of any Informatica product or functionality described today remain at the sole discretion of Informatica and should not be relied upon in making a purchasing decision.

Statements made today are based on currently available information, which is subject to change. Such statements should not be relied upon as a representation, warranty or commitment to deliver specific products or functionality in the future.

**Informatica**

# Agenda

- Streaming Overview
- Structured streaming
- Streaming Sources and Targets
- Streaming mapping Configurations
- Window transformation
- Use case & Demo
- Troubleshooting and self-service
- References
- Q&A
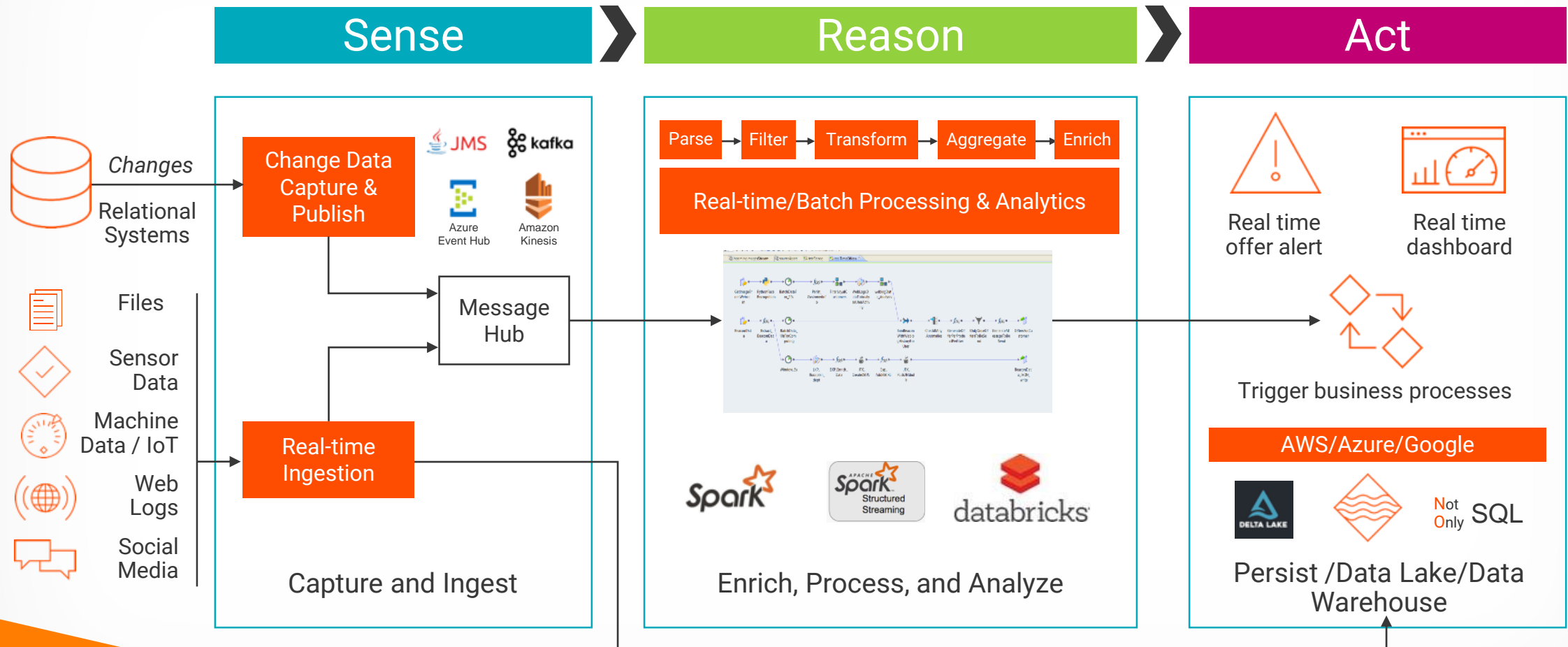
Informatica™

# Streaming Overview

# Streaming Overview

Streaming is the processing of live data streams from unbounded data sources like Kafka, Flume, Kinesis, TCP sockets.
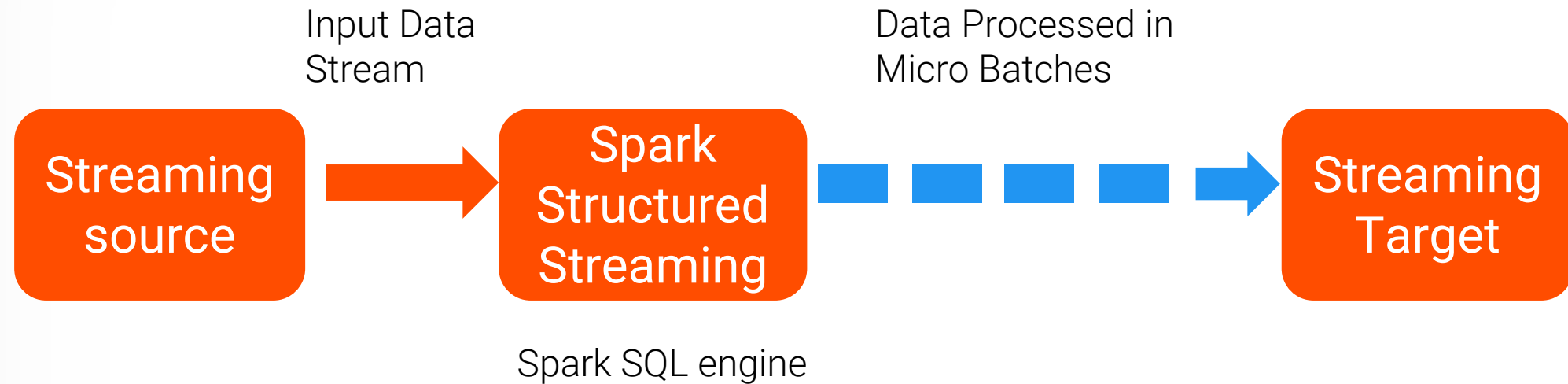


An unbounded data source is one where data is continuously flowing in and there is no definite boundary

# Streaming Overview – Informatica Data Engineering Streaming

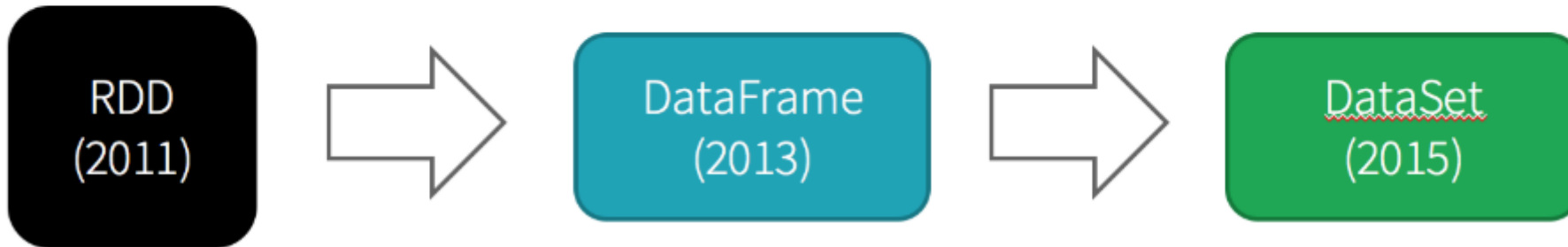| Sense | Reason | Act |
|---|---|---|

**Sense**

Changes

Relational Systems

**Change Data Capture & Publish**

JMS · kafka · Azure Event Hub · Amazon Kinesis

Files

Sensor Data

Machine Data / IoT

Web Logs

Social Media

**Message Hub**

**Real-time Ingestion**

Capture and Ingest

**Reason**

Parse → Filter → Transform → Aggregate → Enrich

**Real-time/Batch Processing & Analytics**

Spark · Spark Structured Streaming · databricks

Enrich, Process, and Analyze

**Act**

Real time offer alert

Real time dashboard

Trigger business processes

**AWS/Azure/Google**

DELTA LAKE · Not Only SQL

Persist /Data Lake/Data Warehouse

Informatica™

# Streaming Process

Input Data
Stream

Data Processed in
Micro Batches

Streaming source → Spark Structured Streaming ⇢ Streaming Target

Spark SQL engine

Spark Structured Streaming receives data from streaming sources such as Kafka and divides the data into micro batches.

Informatica™

# Structured Streaming

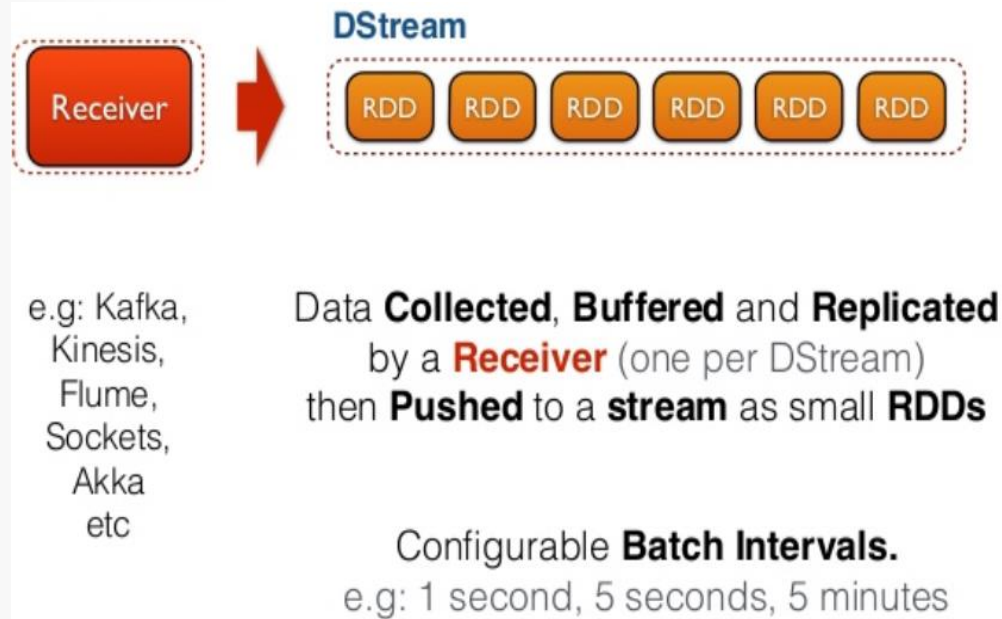# Structured Streaming

## History of Spark APIs



| RDD (2011) | → | DataFrame (2013) | → | DataSet (2015) |
|---|---|---|---|---|

**RDD (2011)**
- Spark Core API
- SparkContext
- Low Level API

**DataFrame (2013)**
- Spark SQL
- RDD + Schema
- SqlContext
- Optimizer support
- High Level(Built on RDD)
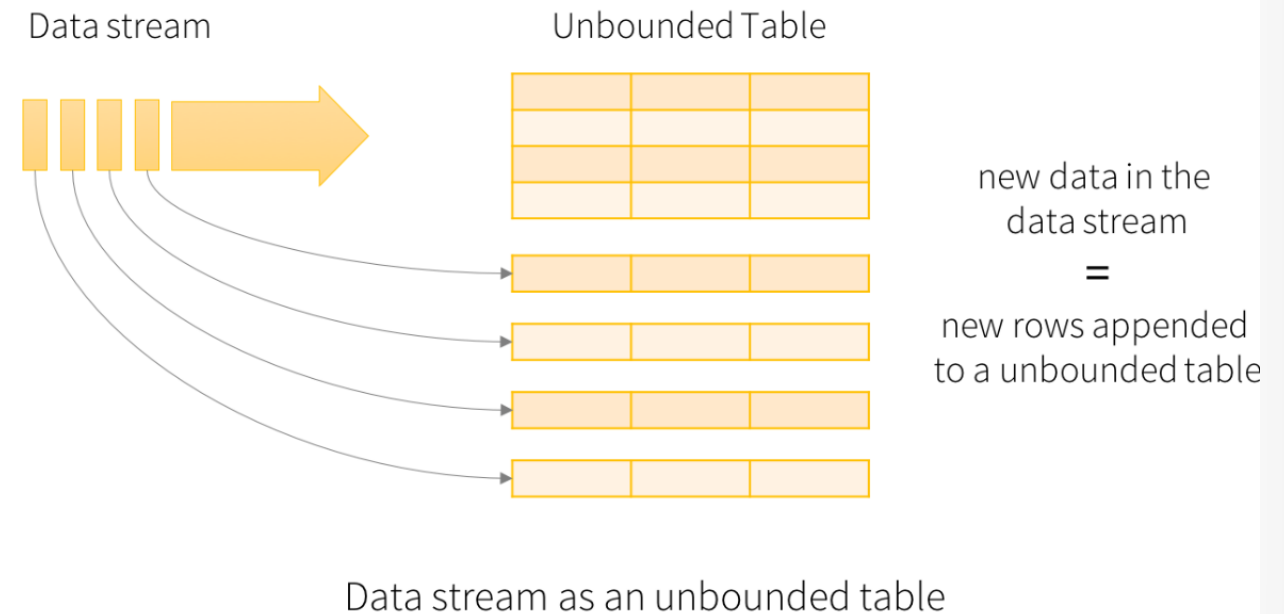
**DataSet (2015)**
- Extension of DataFrame.
- Type Safety

# Structured Streaming

## Spark Streaming (Pre 10.2.2)

**DStream**

Receiver → RDD RDD RDD RDD RDD RDD

e.g: Kafka, Kinesis, Flume, Sockets, Akka etc

Data **Collected**, **Buffered** and **Replicated** by a **Receiver** (one per DStream) then **Pushed** to a **stream** as small **RDDs**

Configurable **Batch Intervals.**
e.g: 1 second, 5 seconds, 5 minutes

**RDD**

## Structured Streaming (10.2.2 & beyond)

Data stream

Unbounded Table

new data in the data stream

=

new rows appended to a unbounded table

Data stream as an unbounded table

**DataFrame**

**Informatica**™

# Structured Streaming – Why ?

**Leverage Spark Optimization**

Dstream cannot leverage the optimizations offered by Spark SQL's Catalyst optimizer and Spark's Tungsten Optimization especially managing Aggregator state management.

DataFrame - Can leverage all the optimizations offered by Spark SQL Catalyst optimizer and Tungsten Optimization.

**Handling Late Data**

This is exclusively a Structured Streaming , we can control how late the window can wait before it can evicted from Result Table and written to target through Watermark property

**Output mode**

There is no output mode in Dstream. It is append by default.
Determine how and when data needs to be evicted from Result Table to the target. Supported output modes are Append , Update and Complete

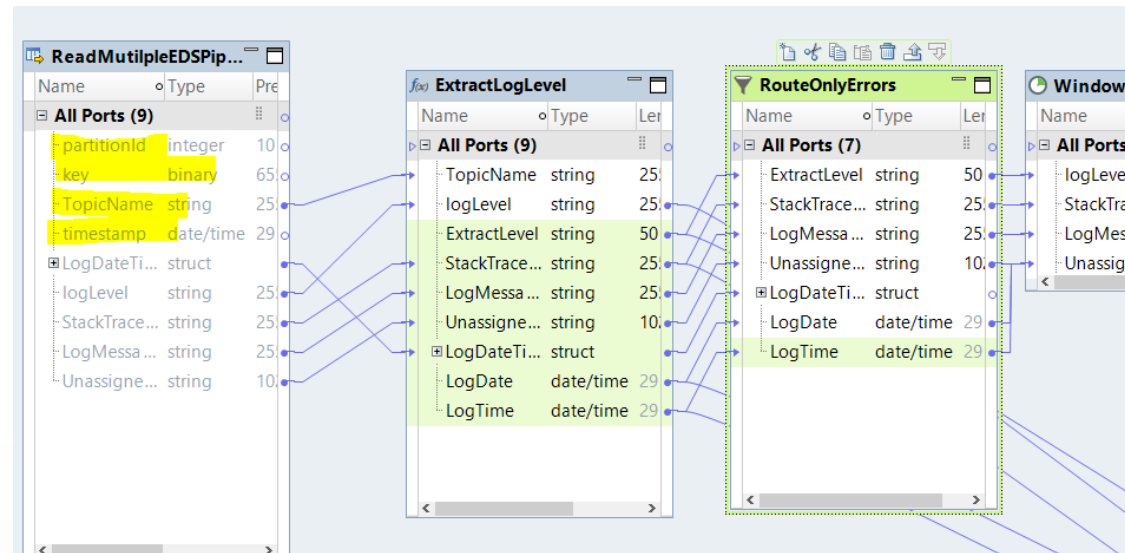**Message delivery**

Guarantees exactly once message delivery

**Informatica**™

# Structured Streaming – Why ?  Contd..

**Message Header Support**
- Enable developers to use message headers from streaming sources
- Transformations can be applied on message header data

**How does It help ?**
- Customers can now use message metadata for better analytics on the data.
- No need to parse the whole message.

# Streaming Sources and Targets

# Streaming Sources and Targets

Spark

**Sources**

- Kafka
- JMS
- Amazon Kinesis
- Azure Event Hubs
- Confluent Kafka
- MapR Streams

**Targets**

- Kafka
- JMS
- Amazon Kinesis
- Azure Event Hubs
- Confluent Kafka
- HBase
- MapR Streams
- Amazon S3
- Complex file Data Object
- ADLS Gen1,Gen2
- Hive
- JDBC compliant Relational Database.
- Snowflake

Informatica™

# Streaming Sources and Targets

## DataBricks (Azure)

### Sources

- Azure Event Hubs

### Targets

- Azure Event Hubs
- ADLS Gen2
- Databricks Delta Lake

# Streaming Sources and Targets : File Formats

The following table shows the different file formats supported in Data Engineering Streaming :

| Format | Schema Type | Amazon Kinesis Firehose | Amazon S3 | Azure Data Lake Store | Azure Event Hub | Complex File | JMS | Kafka | MapR Streams |
|--------|-------------|-------------------------|-----------|-----------------------|-----------------|--------------|-----|-------|--------------|
| Avro | Flat | Not supported | Supported | Supported | Supported | Supported | Not supported | Supported | Supported |
| Avro | Hierarchical | Not supported | Supported | Supported | Supported | Supported | Not supported | Supported | Supported |
| Binary | Binary | Supported | Not Supported | Supported | Supported | Supported | Supported | Supported | Supported |
| Flat | Flat | Not Supported | Supported | Not Supported | Supported | Not supported | Supported | Supported | Not Supported |
| JSON | Flat | Supported | Supported | Supported | Supported | Supported | Supported | Supported | Supported |
| JSON | Hierarchical | Supported | Supported | Supported | Supported | Supported | Supported | Supported | Supported |
| XML | Flat | Not supported | Not Supported | Supported | Supported | Supported | Supported | Supported | Supported |
| XML | Hierarchical | Not supported | Not Supported | Supported | Supported | Supported | Supported | Supported | Supported |

Informatica™

# Streaming mapping Configurations

# Streaming mapping Configurations

- It must have a streaming source.
- For File based Targets, DES provides rollover mechanism of the output file, for downstream application to consume the data seamlessly.

  - Complex file Data object
  - S3
  - ADLS gen1,gen2



© Informatica. Proprietary and Confidential.

# Streaming mapping Configurations

## Streaming properties

- Batch interval
- Cache refresh interval
- State Store Connection
- Checkpoint Directory



© Informatica. Proprietary and Confidential.

- Window transformation

# Window Transformation

In a streaming mapping, depending on your use case, you might want to apply some aggregation over data collected by time (say, every 5 minutes or every hour), e.g

- Average speed of vehicles every 5 min
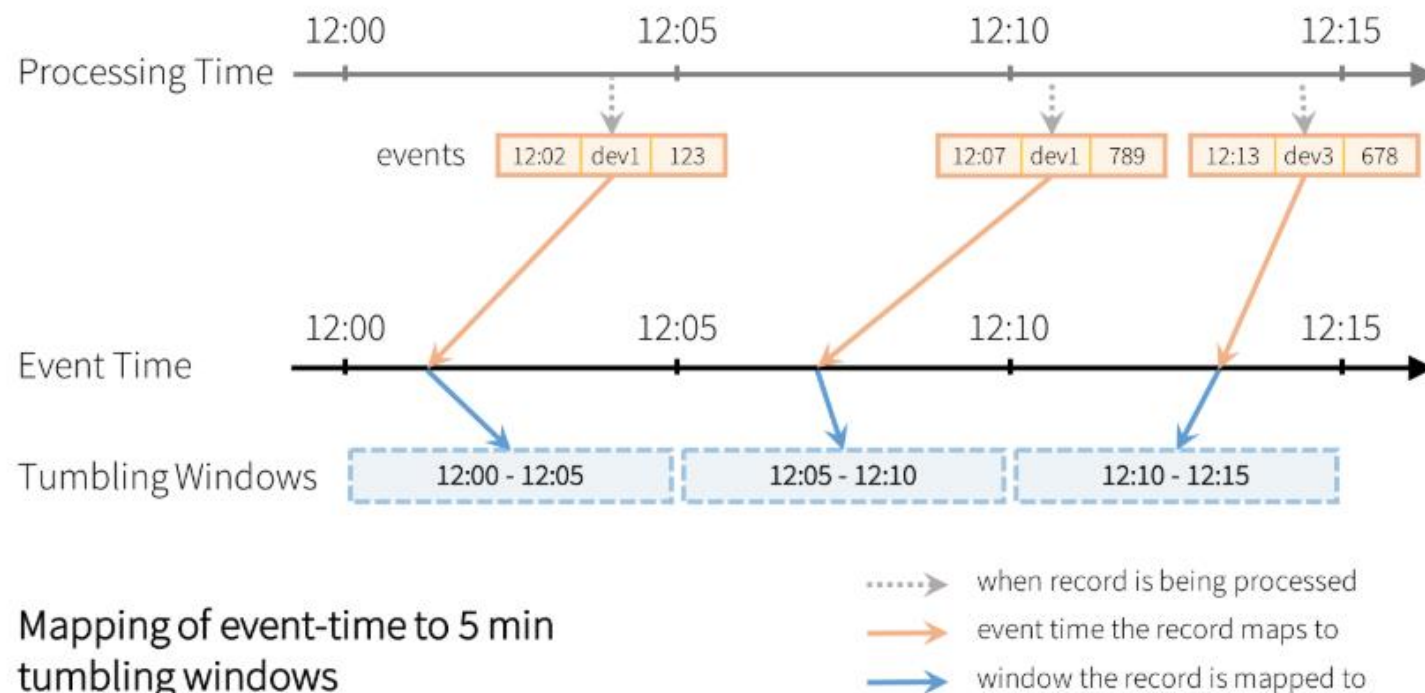- Calculate Maximum value of a stock every min

So, To introduce bounded intervals to unbounded data, use a **Window transformation.**

**Window Types:**

- Tumbling : Max value of a stock price **every five minutes** for stock prices **collected over a five-minute** time interval
- Sliding :  Max value of a stock price **every minute** for stock prices **collected over a five-minute** time interval

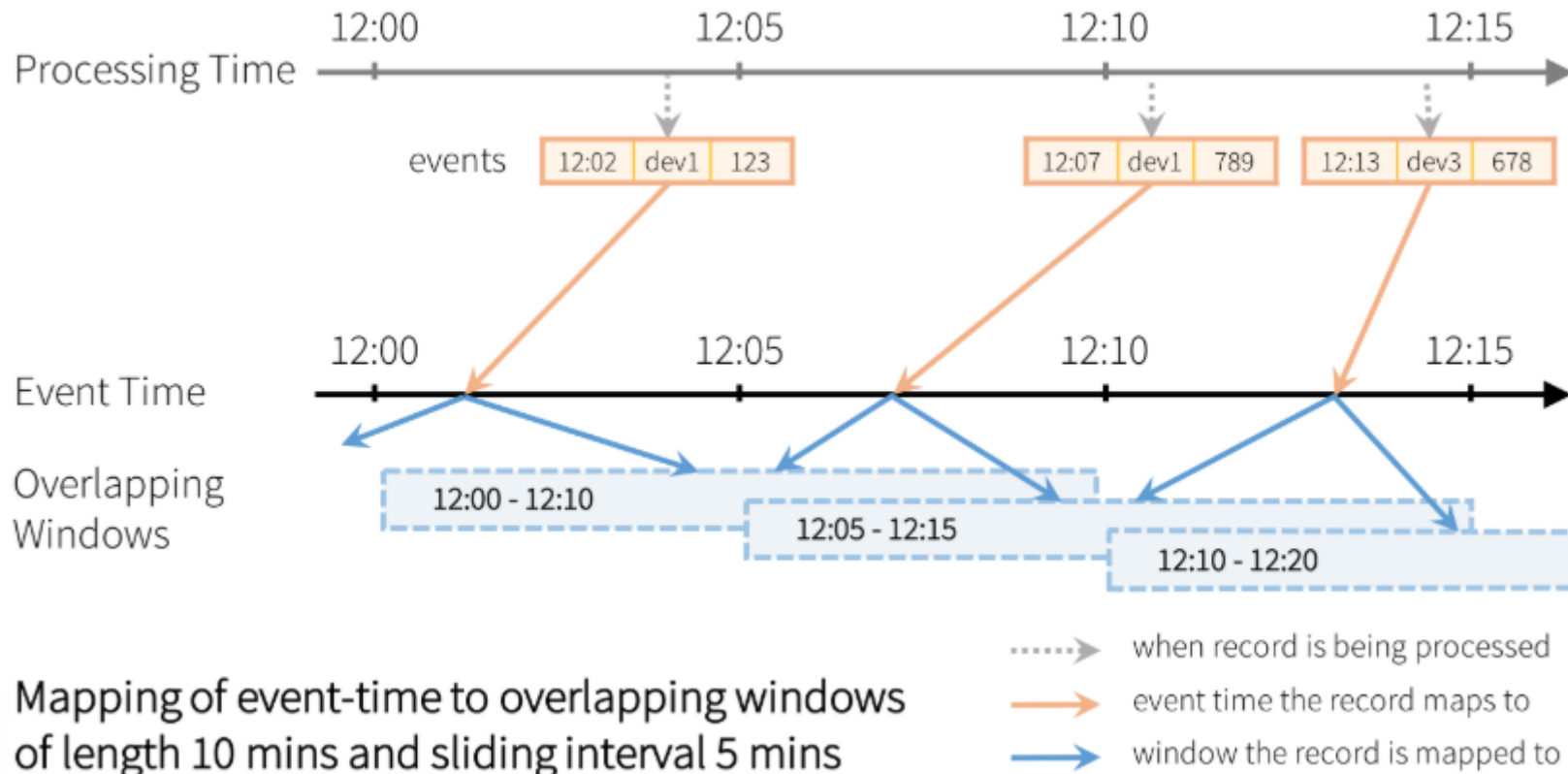Informatica

# Window Transformation - Tumbling

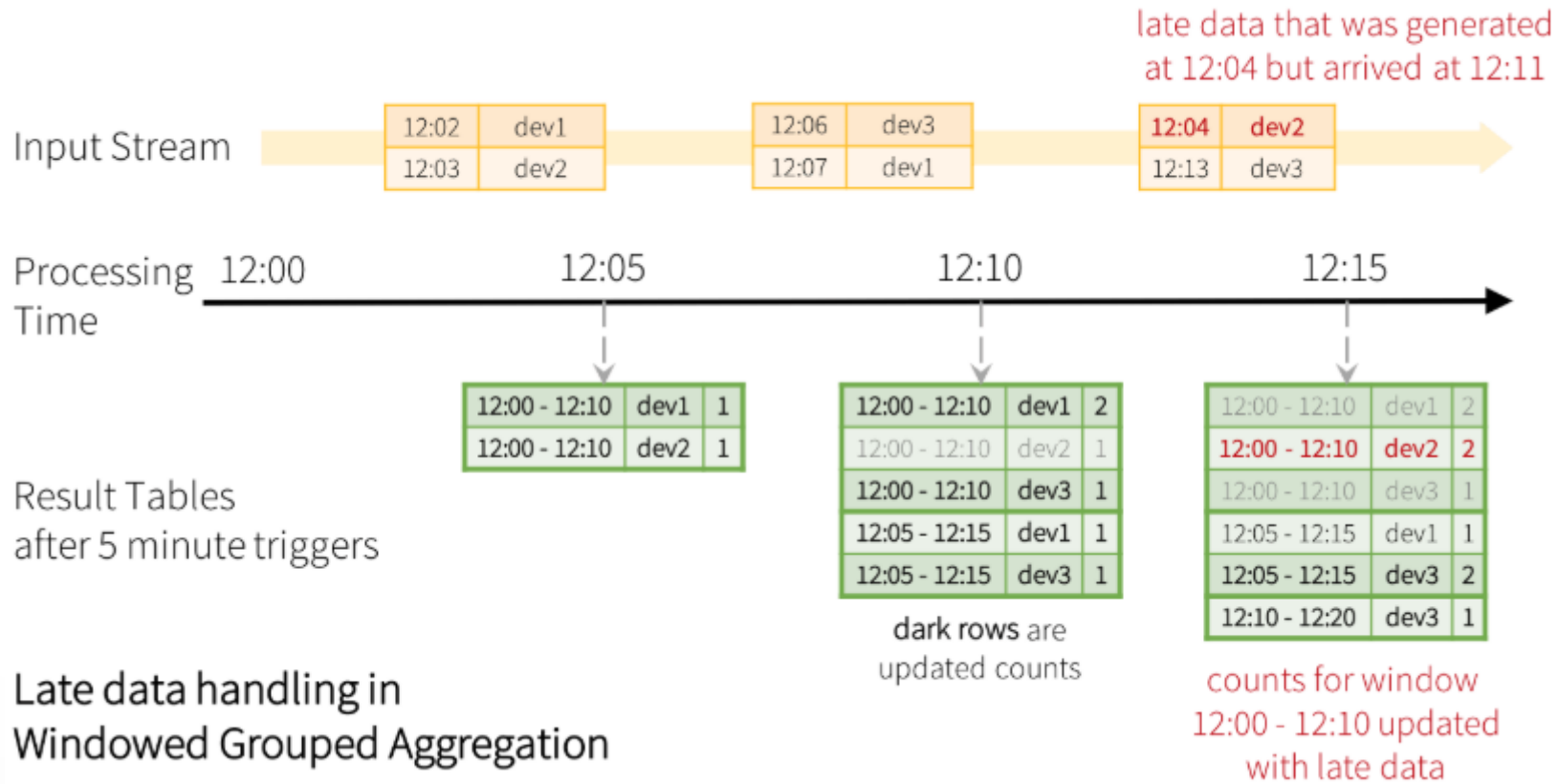Every record is going to be assigned to a 5 minute tumbling window as illustrated below



© Informatica. Proprietary and Confidential.

# Window Transformation - Sliding

Every record will be assigned to multiple overlapping windows as illustrated below.



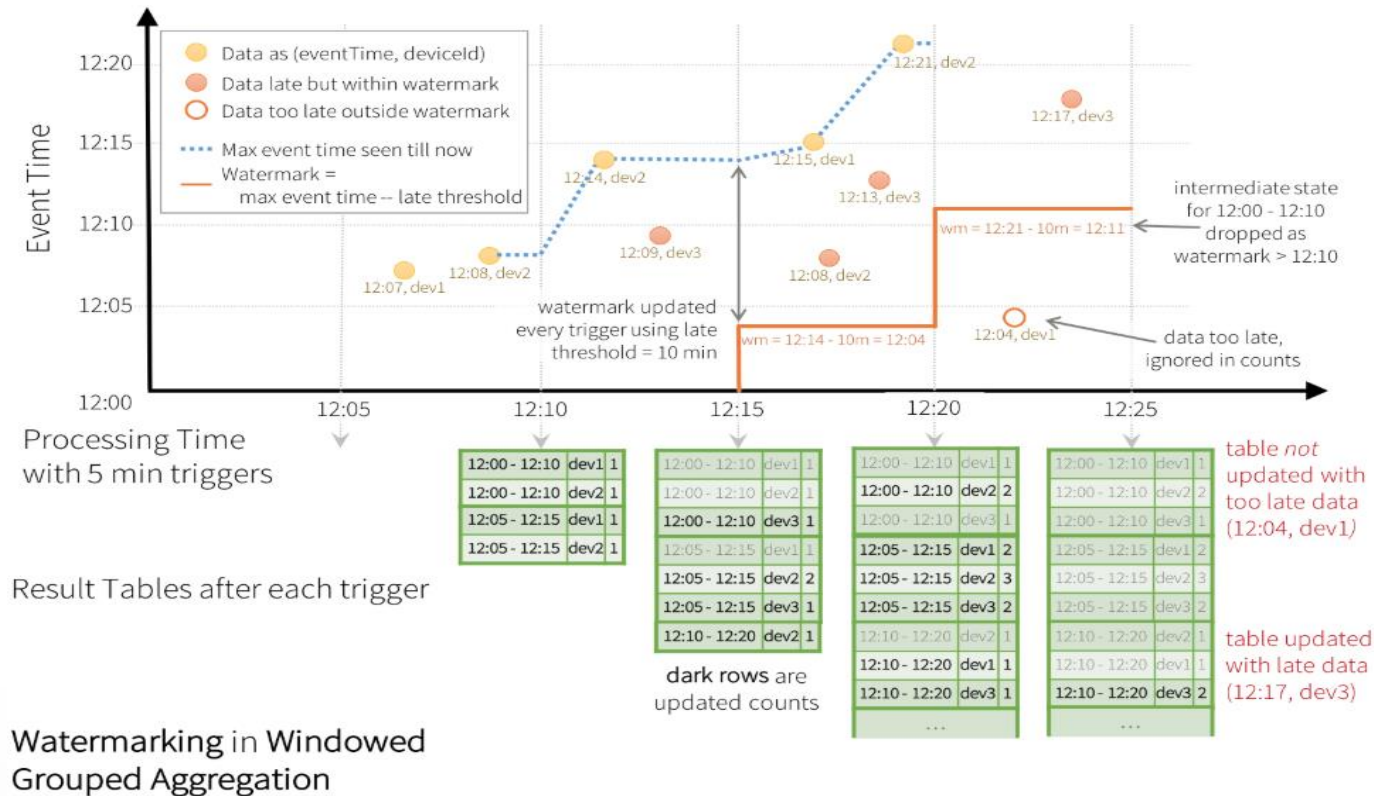© Informatica. Proprietary and Confidential.

# Window Transformation - Sliding

Automatically handles late and out-of-order data



© Informatica. Proprietary and Confidential.

# Window Transformation - WaterMark

The watermark delay defines threshold time for a delayed event to be accumulated into a data group. *"**Watermark delay**"* gets computed at the beginning of every batch based on the latest data arrived in the previous batch.



Watermarking in Windowed Grouped Aggregation

# Window Transformation – sum up

- **Window Type**
- **Window Size**
- **Sliding Interval**
- **Watermark Delay**



© Informatica. Proprietary and Confidential.

# Use case & Demo

# Use Case & Demo

Imagine you started a ride hauling company and need to check if the vehicles are over-speeding. We will create a simple near real-time streaming application to calculate the maximum speed of vehicles every few seconds, while talking about **the concept of window transformation**

© Informatica. Proprietary and Confidential.

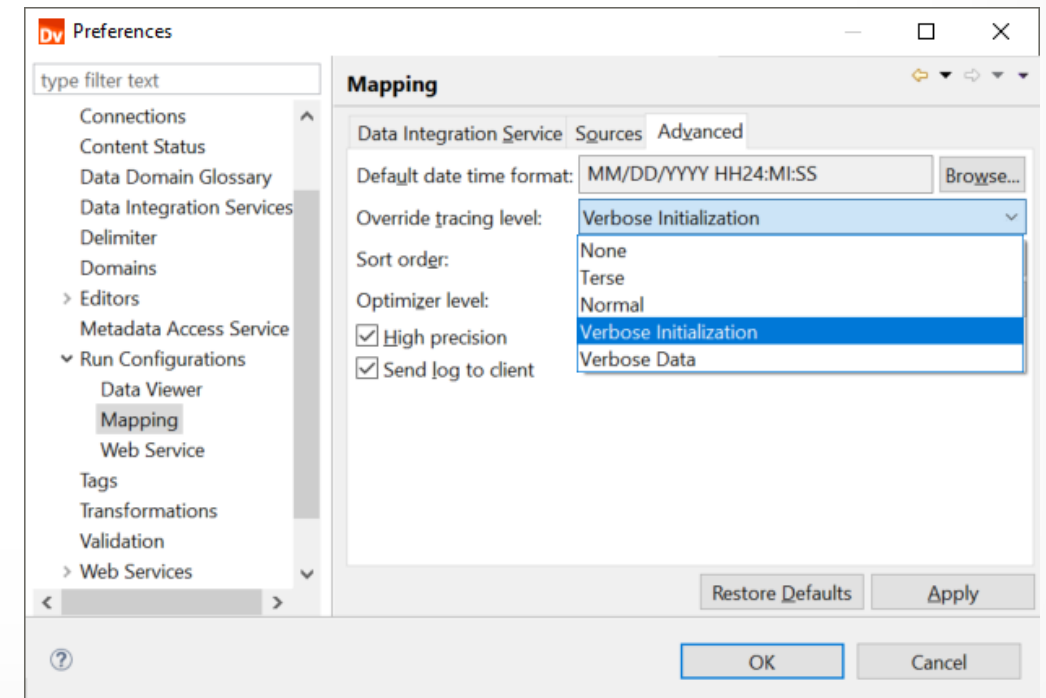# Troubleshooting / Self-Service

# Troubleshooting

## Logs

- Mapping log

- Spark application log

## Override Tracing – Log level

- Normal - INFO

- Verbose Init - DEBUG [Recommended for debugging]

- Verbose Data - DEBUG

# Troubleshooting

## spark.driver.extraJavaOptions | spark.executor.extraJavaOptions

**Advanced Properties**

| Name | Value |
|------|-------|
| infaspark.pythontx.exec | |
| infaspark.pythontx.executorEnv.LD_PRELOAD | |
| infaspark.pythontx.executorEnv.PYTHONHOME | $INFA_HOME/services/shared/spark/python/ |
| infaspark.pythontx.submit.lib.JEP_HOME | |
| spark.driver.extraJavaOptions | -Djava.security.egd=file:/dev/./urandom  -XX:MaxMetaspaceSize=256M  -XX:+UseG1GC  -XX:MaxGCPauseMillis=500  -verbose:class |
| spark.driver.maxResultSize | 4G |
| spark.driver.memory | 4G |
| spark.dynamicAllocation.maxExecutors | 1000 |
| spark.eventLog.enabled | true |
| spark.executor.cores | 2 |
| spark.executor.extraJavaOptions | -Djava.security.egd=file:/dev/./urandom  -XX:MaxMetaspaceSize=256M  -XX:+UseG1GC  -XX:MaxGCPauseMillis=500  -verbose:class |
| spark.executor.memory | 6G |

Hadoop connection

Spark application log

Log Type: stdout
Log Upload Time: Tue Jul 09 11:13:00 -0500 2019
Log Length: 4036717
Showing 4096 bytes of 4036717 total. Click here for the full log.
jar]
[Loaded org.apache.spark.deploy.yarn.ApplicationMaster$$anonfun$finish$2 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-yarn_2.11-2.3.1.jar]
[Loaded org.apache.spark.deploy.yarn.ApplicationMaster$$anon$4$$anonfun$run$9 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-yarn_2.11-2.3.1.jar]
[Loaded java.util.IdentityHashMap$IdentityHashMapIterator from /data/yarn/nm/filecache/83/infa_rpm.tar/jre/lib/rt.jar]
[Loaded java.util.IdentityHashMap$KeyIterator from /data/yarn/nm/filecache/83/infa_rpm.tar/jre/lib/rt.jar]
[Loaded org.apache.hadoop.util.ShutdownHookManager$2 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/hadoop/CDH_5.13/lib/hadoop-common-2.6.0-cdh5.13.0.jar]
[Loaded org.apache.spark.util.SparkShutdownHookManager$$anonfun$runAll$1 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-core_2.11-2.3.1.jar]
[Loaded org.apache.spark.util.SparkShutdownHookManager$$anonfun$runAll$1$$anonfun$apply$mcV$sp$1 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-core_2.11-2.3.1.jar]
[Loaded org.apache.spark.deploy.yarn.YarnRMClient$$anonfun$1 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-yarn_2.11-2.3.1.jar]
[Loaded org.apache.spark.deploy.yarn.ApplicationMaster$$anonfun$unregister$1 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-yarn_2.11-2.3.1.jar]
[Loaded org.apache.spark.deploy.yarn.ApplicationMaster$$anonfun$unregister$1$$anonfun$apply$5 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-yarn_2.11-2.3.1.jar]
[Loaded org.apache.spark.deploy.yarn.ApplicationMaster$$anonfun$unregister$1$$anonfun$apply$6 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-yarn_2.11-2.3.1.jar]
[Loaded org.apache.spark.deploy.yarn.ApplicationMaster$$anonfun$unregister$2 from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/spark-yarn_2.11-2.3.1.jar]
[Loaded org.apache.hadoop.yarn.api.protocolrecords.impl.pb.FinishApplicationMasterRequestPBImpl from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/hadoop/CDH_5.13/lib/hadoop-yarn-common-2.6.0-cdh5.13.0.jar]
[Loaded org.apache.hadoop.yarn.proto.YarnServiceProtos$FinishApplicationMasterRequestProto$Builder from file:/data/yarn/nm/filecache/83/infa_rpm.tar/services/shared/spark/lib_spark_2.3.1_hadoop_2.7.0/hadoop-yarn-api-2.7.3.jar]

*Informatica*

# Troubleshooting

- If you are upgrading from 10.2.1 -> 10.2.2 & later release, recreate the Data Objects for message header support.

- [Common Issues](#)

© Informatica. Proprietary and Confidential.

# References

- [Data Engineering Streaming User Guide](#)

- [Structured Streaming](#)

- [Product Availability Matrix](#)

- [Data Engineering community forum](#)

- [Release notes](#)

Q&A

Informatica

Thank You

Informatica™