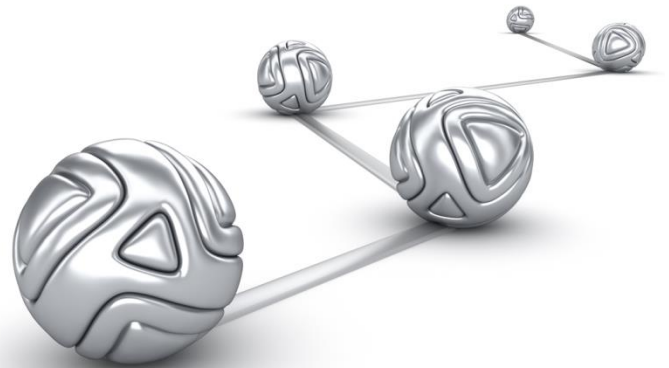




DMX-h ETL Use Case Accelerator  
*Fraud Detection using Apache Kafka*



© Syncsort® Incorporated, 2016

All rights reserved. This document contains proprietary and confidential material, and is only for use by licensees of DMExpress. This publication may not be reproduced in whole or in part, in any form, except with written permission from Syncsort Incorporated. Syncsort is a registered trademark and DMExpress is a trademark of Syncsort, Incorporated. All other company and product names used herein may be the trademarks of their respective owners.

The accompanying DMExpress program and the related media, documentation, and materials ("Software") are protected by copyright law and international treaties. Unauthorized reproduction or distribution of the Software, or any portion of it, may result in severe civil and criminal penalties, and will be prosecuted to the maximum extent possible under the law.

The Software is a proprietary product of Syncsort Incorporated, but incorporates certain third-party components that are each subject to separate licenses and notice requirements. Note, however, that while these separate licenses cover the respective third-party components, they do not modify or form any part of Syncsort's SLA. Refer to the "Third-party license agreements" topic in the online help for copies of respective third-party license agreements referenced herein.

---

## Table of Contents

|          |  |          |
|----------|--|----------|
| <b>1</b> | <b>Introduction .....</b>                                  | <b>1</b> |
| <b>2</b> | <b>Fraud Detection with DMX-h IX.....</b>                  | <b>2</b> |
| 2.1      | J_LoadKafkaMessagesToHDFS.dxj Subjob.....                  | 2        |
| 2.1.1    | T_LoadKafkaMessagesToHDFS.dxt Task .....                   | 2        |
| 2.2      | J_FraudDetection.dxj Subjob .....                          | 3        |
| 2.2.1    | T_FraudDetection.dxt Task.....                             | 3        |
| 2.2.2    | T_LoadPotentialFraudulentTransactionsToKafka.dxt Task..... | 4        |
| 2.2.3    | T_LoadValidTransactionsToHive.dxt Task .....               | 4        |

# 1 Introduction

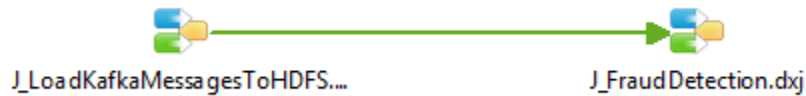
One use case for Apache Kafka message queues is point-of-sale debit card fraud detection, which involves determining whether a debit card is being used to make unauthorized purchases or cash withdrawals based on transaction history and cardholder information. This use case accelerator implements fraud detection in DMX-h ETL by extracting debit card transaction messages from an Apache Kafka topic and joining them with the associated cardholder information. Potentially fraudulent transactions are loaded into another Kafka topic for further inspection, and non-fraudulent transactions are loaded into a Hive table.

DMX-h ETL use case accelerators are developed as standard ETL jobs and can be run on an edge node, single cluster node, or in a Spark/MapReduce cluster without making any changes to the application.

For guidance on setting up and running this and other use case accelerators, see the [Guide to DMX-h ETL Use Case Accelerators](#). Note that the Guide details additional Windows setup that must be done prior to working with this UCA in DMExpress.

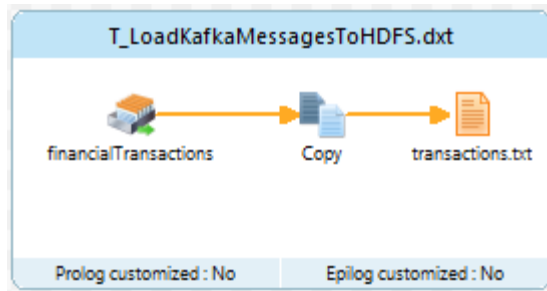
## 2 Fraud Detection with DMX-h IX

The Fraud Detection solution in DMX-h ETL consists of a job J\_FraudDetectionDemo.dxd, containing two subjobs. The first subjob uses a copy task to ingest Kafka messages (debit card transactions) into HDFS, and the second subjob flags the transaction as either valid or potentially fraudulent.



### 2.1 J\_LoadKafkaMessagesToHDFS.dxd Subjob

This subjob consists of one copy task.



#### 2.1.1 T\_LoadKafkaMessagesToHDFS.dxt Task

This copy task reads debit card transactions from the Kafka topic “financialTransactions” and loads them into HDFS.

The batch size is 100 messages, and the group ID is set to allow Kafka to recognize which messages were already read by this consumer. If there are fewer than 100 messages, the job will wait until the queue has at least 100 messages. The financialTransactions topic has exactly 100 messages, so after running the job once, the specified group ID will have consumed all the messages and the queue will be empty. If you want to run the job again, you will need to repopulate the queue by re-running the prep script as described in the Guide to DMX-h ETL Use Case Accelerators.

The Kafka source is defined with a message queue connection, and the HDFS target file is defined with a remote file connection to the Hadoop cluster.





- DMExpress Task (Copy)
  - Source
    - \$(HDFS\_TARGET\_DIR)/validTransactions.txt=\${HDFS\_TARGET\_DIR}/validTransactions.txt (File type : Text, UNIX (LF record terminator); Server connection)
  - Target
    - default.validtransactions (Database Connection : HiveODBC; Insertion Method : Table)
  - Metadata
    - Database Connections
      - HiveODBC (DBMS : Hive; Access method : ODBC; Database : HiveODBC; Authentication : Auto-detect; Connect as : dmxdemo; Password : \*\*\*\*\*)
    - External Metadata (Type : DMExpress; File : T\_FraudDetection.dxt)
    - External Metadata (Type : DMExpress; File : T\_LoadKafkaMessagesToHDFS.dxt)
      - Remote File Connections
        - HDFS\_Connection (Server : \${HDFS\_SERVER}=\${HDFS\_SERVER}); Connection type : Hadoop; File system scheme : hdfs; Authentication : None)



## About Syncsort

Syncsort provides enterprise software that allows organizations to collect, integrate, sort, and distribute more data in less time, with fewer resources and lower costs. Thousands of customers in more than 85 countries, including 87 of the Fortune 100 companies, use our fast and secure software to optimize and offload data processing workloads. Powering over 50% of the world's mainframes, Syncsort software provides specialized solutions spanning "Big Iron to Big Data", including next gen analytical platforms such as Hadoop, cloud, and Splunk. For more than 40 years, customers have turned to Syncsort's software and expertise to dramatically improve performance of their data processing environments, while reducing hardware and labor costs. Experience Syncsort at [www.syncsort.com](http://www.syncsort.com).

**Syncsort Inc.**

**50 Tice Boulevard, Suite 250, Woodcliff Lake, NJ 07677**

**201.930.8200**