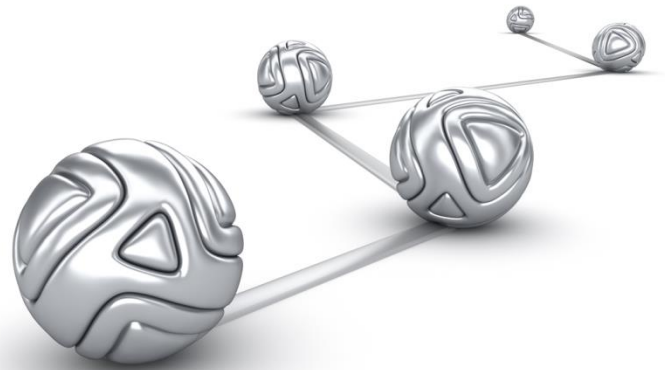


syncsort

DMX-h ETL Use Case Accelerator
Mainframe Variable Processing



© 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

1	Introduction	1
2	Mainframe Variable Processing with DMX-h.....	2
2.1	J_MainframeVariableExtractLocalToHDFS Subjob.....	2
2.1.1	T_MainframeVariableExtractLocalToHDFS Task.....	2
2.2	T_MainframeVariableTotalRoyaltyPerBook Task.....	3

1 Introduction

Mainframe variable length files are EBCDIC-encoded files with varying record lengths. Each record has a 4-byte prefix which provides the length of the record. As is, these files are non-splittable when stored in HDFS, which effectively means you cannot achieve distributed processing when running a Spark/MapReduce job on such a file.

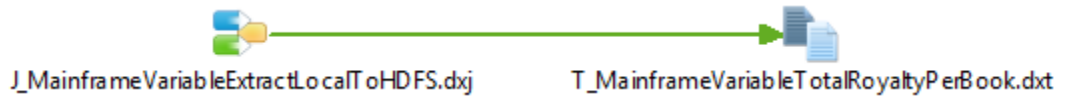
With DMX-h, you can ingest these files into HDFS in a distributable manner by converting them to DMX-h's Mainframe Hadoop Distributable format. You can then run Spark/MapReduce jobs on these EBCDIC-encoded files.

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 the 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](#).

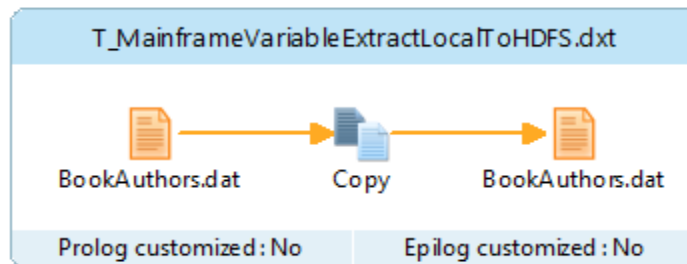
2 Mainframe Variable Processing with DMX-h

The Mainframe variable processing solution consists of a DMX-h job, `J_MainframeVariableProcessingTotalRoyaltyPerBook.dxj`, with one subjob followed by a copy task.



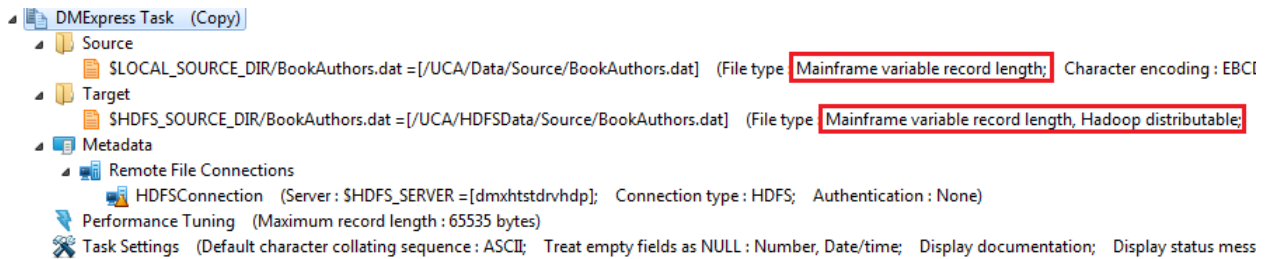
2.1 J_MainframeVariableExtractLocalToHDFS Subjob

This subjob consists of a single copy task that reads the mainframe variable length file and loads it to HDFS. Since it's an ingestion job, DMX-h will run it on the edge node.



2.1.1 T_MainframeVariableExtractLocalToHDFS Task

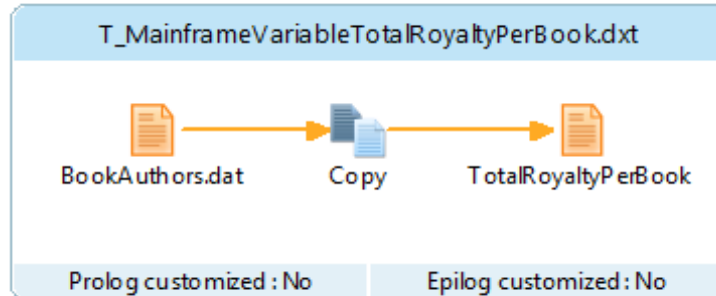
This copy task reads a mainframe variable length file from the local file system and ingests it into HDFS in “Mainframe variable record length, Hadoop distributable” format. You can specify a mainframe server connection in the source file if you are reading directly from the mainframe instead of local disk.



Since this task is ingesting a copy of the mainframe data into Hadoop without interpreting the records, there's no need to provide a copybook to this task.

2.2 T_MainframeVariableTotalRoyaltyPerBook Task

This task, run in the cluster, reads the output of the subjob, which is a Mainframe Hadoop Distributable file. It interprets the EBCDIC data using a copybook, and calculates the total royalty percentage on a book.



Royalty percentage is an element in an OCCURS DEPENDING ON array, the length of which depends on NUM-AUTHORS. We add the royalty percentage for each of the authors in the array to find the total royalty percentage per book. The target of this task is also a mainframe EBCDIC-encoded file, which can later be put back on the mainframe.

- DMExpress Task (Copy)
 - Source
 - SHDFS_SOURCE_DIR/BookAuthors.dat = [/UCA/HDFSData/Source/BookAuthors.dat] (File type : Mainframe variable record leng)
 - {BOOK-SALES-AUTHORS-REC}
 - Target
 - SHDFS_TARGET_DIR/TotalRoyaltyPerBook = [/UCA/HDFSData/Target/TotalRoyaltyPerBook] (File type : Mainframe variable recc)
 - Reformat (Create target layout : TotalRoyaltyOnBook_Layout; Fixed position layout)
 - {BOOK-SALES-AUTHORS-REC}.ISBN
 - {BOOK-SALES-AUTHORS-REC}.TITLE
 - TOTAL_ROYALTY
 - Metadata
 - Target Record Layouts
 - Values
 - TOTAL_ROYALTY ((BOOK-SALES-AUTHORS-REC).AUTHORS[1].{ROYALY-PCT} + {BOOK-SALES-AUTHORS-REC}.AUTHORS
 - Remote File Connections
 - HDFSConnection (Server : SHDFS_SERVER = [dmxhtstdrvhdp]; Connection type : HDFS; Authentication : None)
 - External Metadata (Type : COBOL, VS COBOL II Release 4 data format; File : BookSalesAuthors.cpy)
 - Record Layouts
 - BOOK-SALES-AUTHORS-REC
 - ISBN (Data type : Number; Format : Decimal, unsigned; Length : 13 bytes)
 - TITLE (Data type : Text; Length : 64 bytes)
 - COPYRIGHT-YEAR (Data type : Number; Format : Decimal, unsigned; Length : 4 bytes)
 - COPIES-PRINTED (Data type : Number; Format : Decimal, packed; Length : 5 bytes)
 - COPIES-SOLD (Data type : Number; Format : Decimal, packed; Length : 5 bytes)
 - COVER-PRICE (Data type : Number; Format : Decimal, packed; Length : 4 bytes; Scaling factor order : -2)
 - NUM-AUTHORS (Data type : Number; Format : Decimal, packed; Length : 2 bytes)
 - AUTHORS (Array of up to 4 elements determined by {NUM-AUTHORS})
 - FNAME (Data type : Text; Length : 24 bytes)
 - MNAME (Data type : Text; Length : 24 bytes)
 - LNAME (Data type : Text; Length : 32 bytes)
 - ROYALY-PCT (Data type : Number; Format : Decimal, unsigned; Length : 3 bytes; Scaling factor order : -1)

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.

Syncsort Inc.

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

201.930.8200