



DMX-h ETL Use Case Accelerator
File Lookup Aggregation



© Syncsort® Incorporated, 2015

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	File Lookup Aggregation with DMX-h IX.....	2
2.1	T_FileLookupAggregation Task	2
Appendix A	File Lookup Aggregation with DMX-h User-defined MapReduce	4
A.1	Map Step	4
A.1.1	MT_FileLookupAggregation.dxt	5
A.2	Reduce Step	6
A.2.1	RT_Aggregation.dxt	7

1 Introduction

DMX-h ETL can efficiently lookup a value in a dimension table and subsequently perform an aggregation in Hadoop MapReduce. In this use case accelerator, DMX-h performs a lookup using two TPC Benchmark H (TPC-H) data sets: a line item file and a supplier file. The lookup is followed by an aggregation on the supplier key.

The use case accelerators are developed as standard DMExpress jobs with just minor accommodations that allow them to be run in or outside of Hadoop:

- When specified to run in the Hadoop cluster, DMX-h Intelligent Execution (IX) automatically converts them to be executed in Hadoop using the optimal Hadoop engine, such as MapReduce. In this case, some parts of the job may be run on the Linux edge node if not suitable for Hadoop execution.
- Otherwise, they are run on a Windows workstation or Linux edge node, which is useful for development and testing purposes before running in the cluster.

While the method presented in the next section is the simplest and most efficient way to develop this example, the more complex user-defined MapReduce solution is provided as a reference in Appendix A.

For guidance on setting up and running the examples both outside and within Hadoop, see [Guide to DMExpress Hadoop Use Case Accelerators](#).

For details on the requirements for IX jobs, see “Developing Intelligent Execution Jobs” in the DMExpress Help.

2 File Lookup Aggregation with DMX-h IX

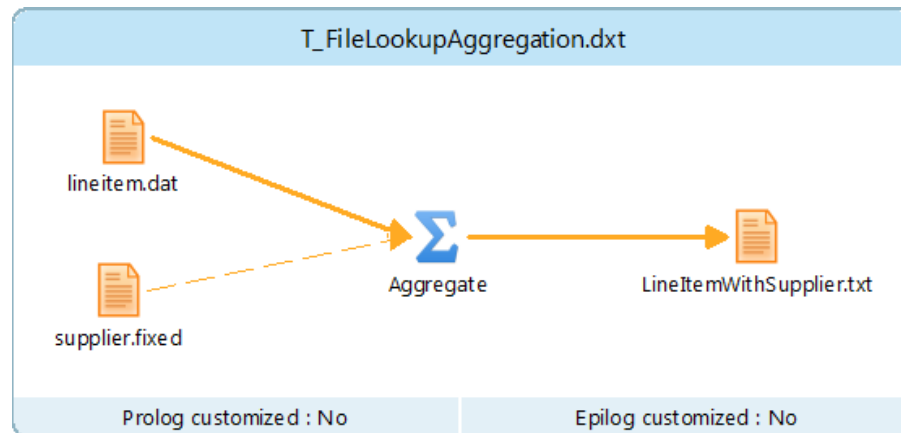
The File Lookup Aggregation solution in DMX-h ETL consists of a job, J_FileLookupAggregation.dxi, with one aggregate task.



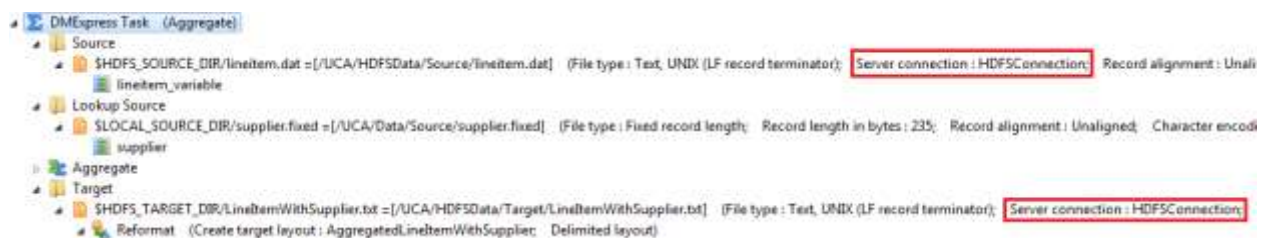
T_FileLookupAggregation.dxt

2.1 T_FileLookupAggregation Task

This task performs the lookup with the fixed length dimension file supplier.fixed and then performs the aggregation.



The HDFS source and target files are defined with a remote file connection to the Hadoop cluster. The Lookup file is located on the edge node.



The Lookup() function is used to look up the account balance in the supplier dimension file that matches each supplier key in the line item source file. A total of these account balances is computed for each supplier by aggregating on the supplier key and adding the supplier account balance. The supplier key and account balance total are included in the reformat for the HDFS target.

- ▲ 📁 Lookup Source
 - ▲ 📄 SLOCAL_SOURCE_DIR/supplier.fixed =[/UCA/Data/Source/supplier.fixed] (File type : Fixed record length; Record length in bytes : 235;
 - 📄 supplier
 - ▲ 📊 Aggregate
 - ▲ 📊 Group by
 - 🔗 lineitem_variable.L_SUPPKEY
 - ▲ 📊 Summarize
 - 📊 total(SupplierAccountBalance)
 - ▲ 📁 Target
 - ▲ 📄 SHDFS_TARGET_DIR/LineItemWithSupplier.txt =[/UCA/HDFSData/Target/LineItemWithSupplier.txt] (File type : Text, UNIX (LF record term
 - 🔗 Reformat (Create target layout : AggregatedLineItemWithSupplier; Delimited layout)
 - 📄 lineitem_variable.L_SUPPKEY
 - 📄 total(SupplierAccountBalance)
 - ▲ 📁 Metadata
 - ▶ 📄 Source Record Layouts
 - ▲ 📄 Target Record Layouts
 - ▶ 📄 AggregatedLineItemWithSupplier
 - ▲ 📄 Values
 - 🔗 SupplierAccountBalance (Lookup('Any',supplier.S_ACCTBAL,'supplier.fixed',supplier.S_SUPPKEY = lineitem_variable.L_SUPPKEY))

Appendix A File Lookup Aggregation with DMX-h User-defined MapReduce

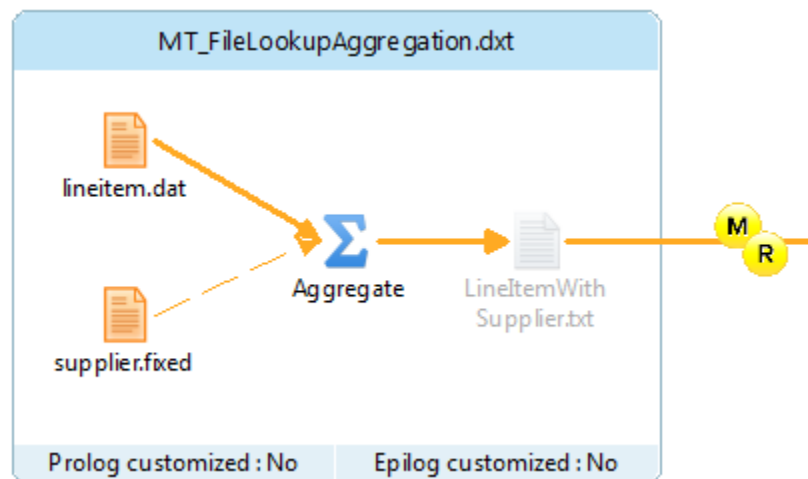
The File Lookup Aggregation job in DMX-h contains a map step and a reduce step, separated by a MapReduce data flow connector, as follows:



The map step performs a lookup and a partial aggregation, while the reduce step performs the final aggregation.

A.1 Map Step

The File Lookup Aggregation map step in DMX-h consists of one task that performs a lookup and a partial aggregation.



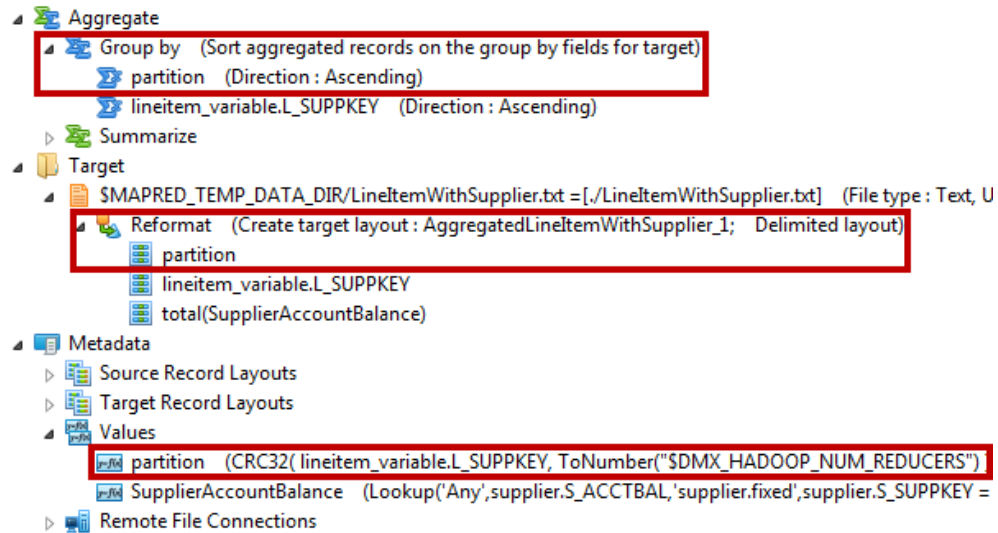
To optimize performance, the first phase of aggregation occurs in the map step as the aggregation reduces the amount of data going over the network to the reducer.

A.1.1 MT_FileLookupAggregation.dxt

This task performs the lookup with the fixed length dimension file and performs the first phase of aggregation:

The Lookup function is used to look up the account balance in the supplier dimension file that matches each supplier key in the line item source file. A partial total of these account balances is computed for each supplier in the subset of data known to each mapper:

“Partition” is the first field in the target as well as the first field of the group by fields. The partition ID determines the reducer to which the data goes. All records that share the same key must go to the same reducer. In this case, the records with the same supplier key value must have the same partition ID; thus, we create a CRC32() hash value based on this key:

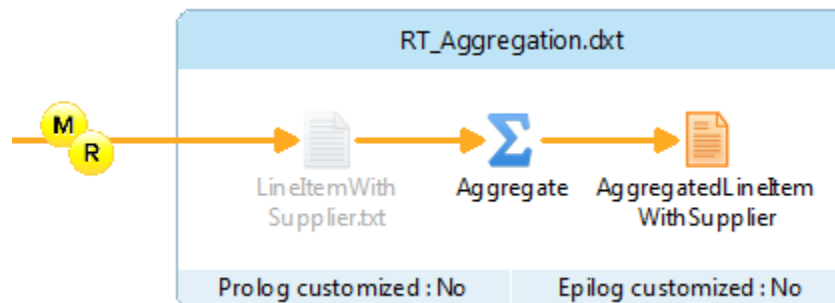


In addition to the value to be hashed, the CRC32() function allows you to provide a number to determine a range of hash values. For partition, the range is defined as the environment variable DMX_HADOOP_NUM_REDUCERS, which provides the number of reducers invoked for your MapReduce job. As a result, the CRC32() function returns a number from 0 to the value of DMX_HADOOP_NUM_REDUCERS.

The data must be ordered based on this partition ID so that the framework can properly distribute the data to the correct reducer. We include the partition ID as the first group by field and sort records based on the group by fields to ensure the data is in the correct sorted order.

A.2 Reduce Step

The File Lookup Aggregation reduce step in DMX-h ETL consists of one aggregate task that computes the final total account balance for each supplier.



A.2.1 RT_Aggregation.dxt

This task reads the line item records containing supplier data from the map step.

- DMExpress Task (Aggregate)
 - Source
 - SMAPRED_TEMP_DATA_DIR/LineItemWithSupplier.txt = [./LineItemWithSupplier.txt] (File type : Text, UNIX (LF record terminator); Record alignment : Unaligned)
 - AggregatedLineItemWithSupplier_1
 - Aggregate
 - Group by (Sort aggregated records on the group by fields for target)
 - AggregatedLineItemWithSupplier_1.L_SUPPKEY (Direction : Ascending)
 - Summarize
 - total(AggregatedLineItemWithSupplier_1.total_SupplierAccountBalance)
 - Target
 - SHDFS TARGET DIR/AggregatedLineItemWithSupplier = [DMXhTestDrive;/UCA/HDFSData/Target/AggregatedLineItemWithSupplier] (File type : Text, UNIX)
 - Reformat (Create target layout : LineItemsWithSupplierAccountBalanceTotals_1; Delimited layout)
 - AggregatedLineItemWithSupplier_1.L_SUPPKEY
 - total(AggregatedLineItemWithSupplier_1.total_SupplierAccountBalance) (Target field name : total_SupplierAccountBalance)
 - Metadata
 - Target Record Layouts
 - External Metadata (Type : DMExpress; File : MT_FileLookupAggregation.dxt)
 - Record Layouts
 - AggregatedLineItemWithSupplier_1
 - partition (Data type : Number; Format : Decimal, edited numeric; Nullability : Empty field; Calculated field number : 1)
 - L_SUPPKEY (Data type : Number; Format : Decimal, edited numeric; Nullability : Empty field; Calculated field number : 2)
 - total_SupplierAccountBalance (Data type : Number; Format : Decimal, edited numeric; Nullability : Empty field; Calculated field number : 3)
 - Remote File Connections

The task aggregates on the supplier key and totals the supplier account balance based on all partial totals computed in the mappers.

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