

# Offload Data Warehousing to Hadoop by using DataStage

See also a guided tour „Offload Data Warehousing to Hadoop by using DataStage” Use IBM® InfoSphere® DataStage® to load Hadoop and use YARN to manage DataStage workloads in a Hadoop cluster (a registered IBM Cloud Id is needed!):

<https://www.ibm.com/cloud/garage/dte/producttour/offload-data-warehousing-hadoop-using-datastage>

See also a YouTube video which describes the details of this guided tour:

<https://www.youtube.com/watch?v=QyWdzCeD6cU>

## Environment

This environment is a front-end Windows workstation that has the DataStage Designer version 11.5 installed. It will connect to a back-end Linux Server that is running the InformationServer Server components, and a configuration known as BigIntegrate. The DataStage Engine will run inside a remote Hadoop Cluster, and participate in the resource management of the Yarn service within Hadoop.

- Pre-configured, Autostarted software and services
- Allows you to see the configuration necessary to utilize the Hadoop cluster for running ETL processes

## Tutorial

In this demo, you use DataStage to complete extract, transform, and load (ETL) data processing in a traditional enterprise data warehouse. You then offload the data and ETL processing into scalable, high-value Hadoop clusters and data lakes.

In this product tour, you will walk through the following tasks:

In this product tour, you get experience with the following features:

- Learn how to run DataStage traditional ETL jobs
- Configure DataStage to run inside Hadoop Clusters
- Examine execution logs to ensure configuration worked correctly

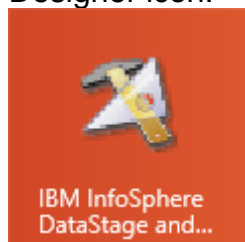
Follow the instructions in this pane to walk through the demo in the left pane.

Run a traditional Data Warehouse ETL job

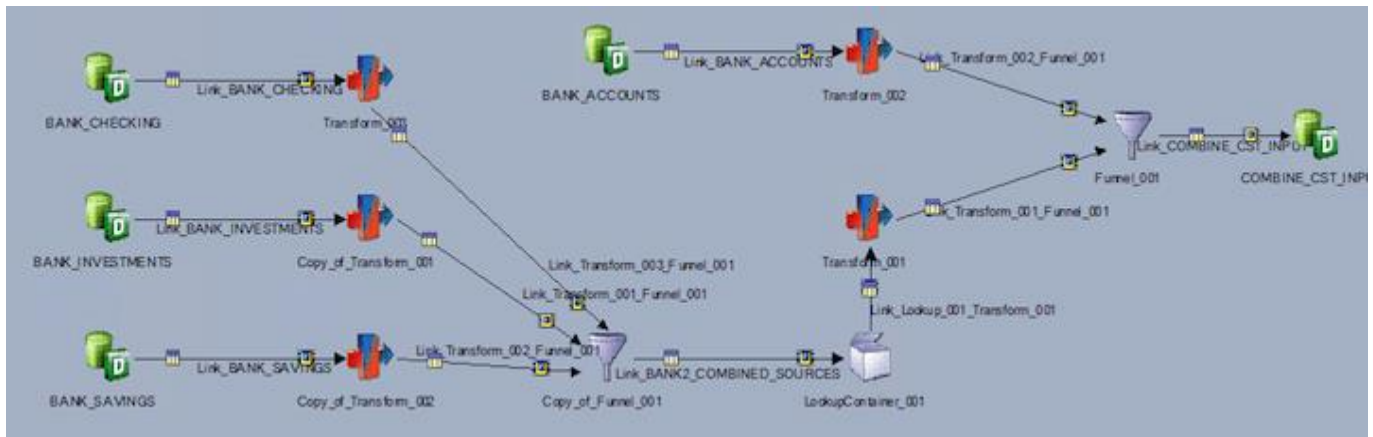
First, you review a DataStage job that combines data from two lines of business into a table. You then run the job to create a combined data repository.

1. Click **Start demo now**

2. From the Start menu, click the IBM InfoSphere DataStage and QualityStage Designer icon.



The DS07\JK\BANK1\And\JK\BANK2\To\COMBINE\CST\INPUT job is displayed.

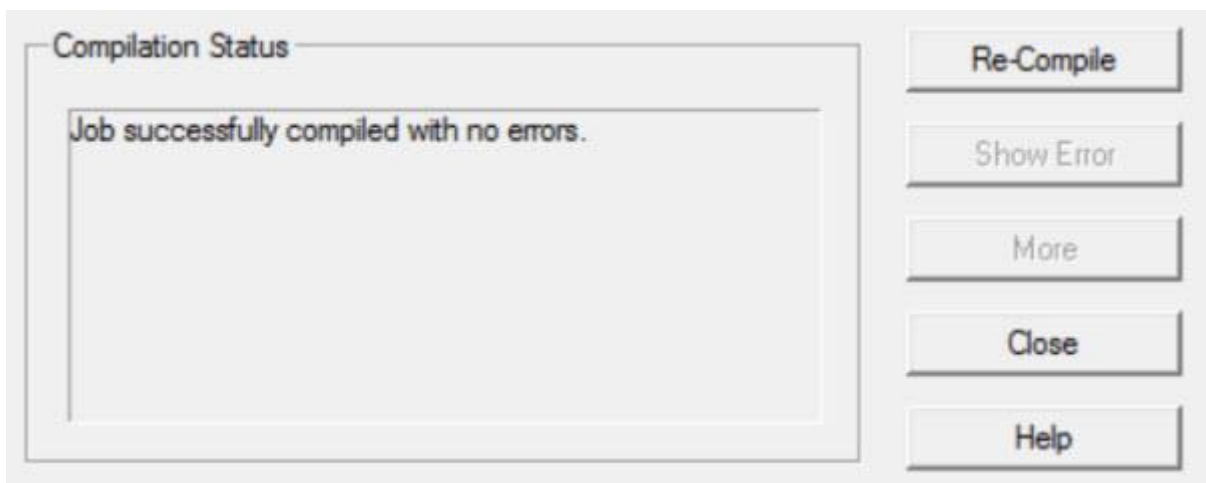


The job has these data sources:

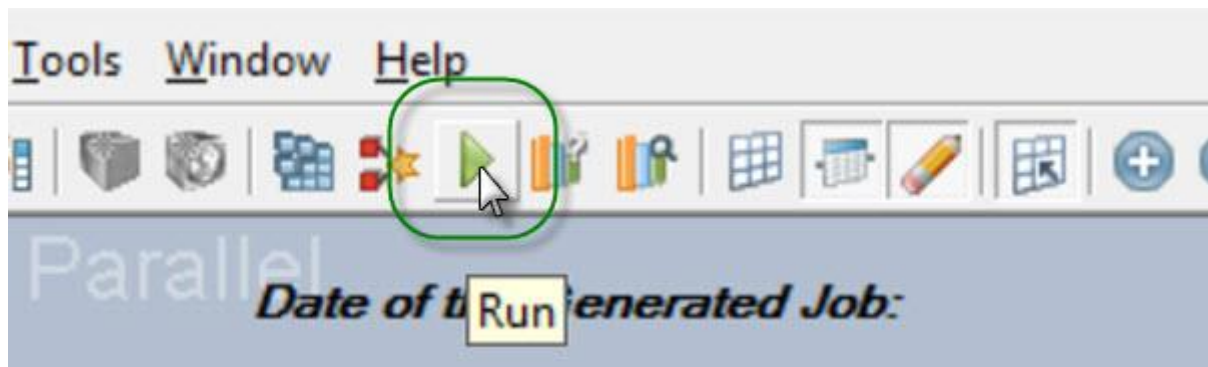
- \* BANK1 ACCOUNTS data, which is labeled as `BANK\_ACCOUNTS`.
- \* BANK2 CHECKING customer data, which is labeled as `BANK\_CHECKING`.
- \* BANK2 INVESTMENT customer data, which is labeled as `BANK\_INVESTMENTS`.
- \* BANK2 SAVINGS customer data, which is labeled as `BANK\_SAVINGS`.

The job's data target is BANK1 and BANK2 combined customer data.

3. Compile the job by clicking the **Compile icon** on the toolbar.



4. Run the job by clicking the **Run** icon on the toolbar.



5. In the Job Run Options window, you can provide values for the runtime job parameters. Use the default values. For the `JKLW\_DBS\_PWD` parameter, type `inf0server`. Click **Run** to run the job.

Name	Value
JKLW_DBS_PWD	
APT_YARN_CONF	/opt/IBM/InformationServer/Server/PXEngine/etc/yam_conf/yamconfig.cf
Configuration file	/opt/IBM/InformationServer/Server/Configurations/default.ap
parameters paramet	(As pre-defined)
hostname	bigdata.ibm.com
Use IPv4 only	True

Workload management

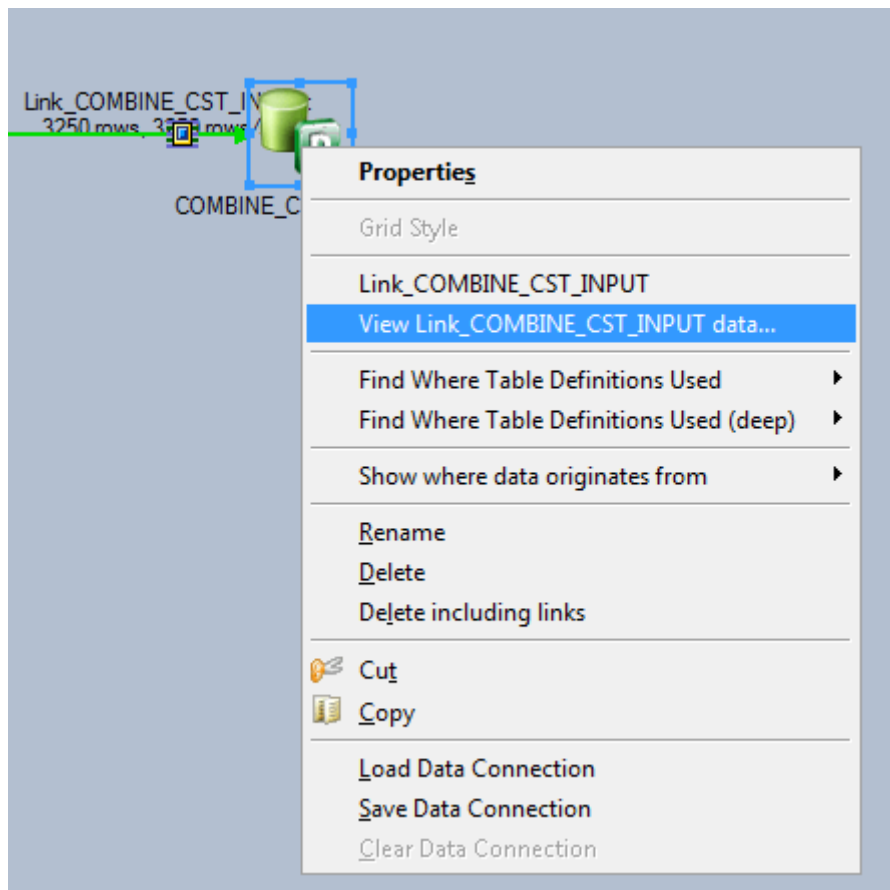
Queue

Project default (HighPriorityJobs)

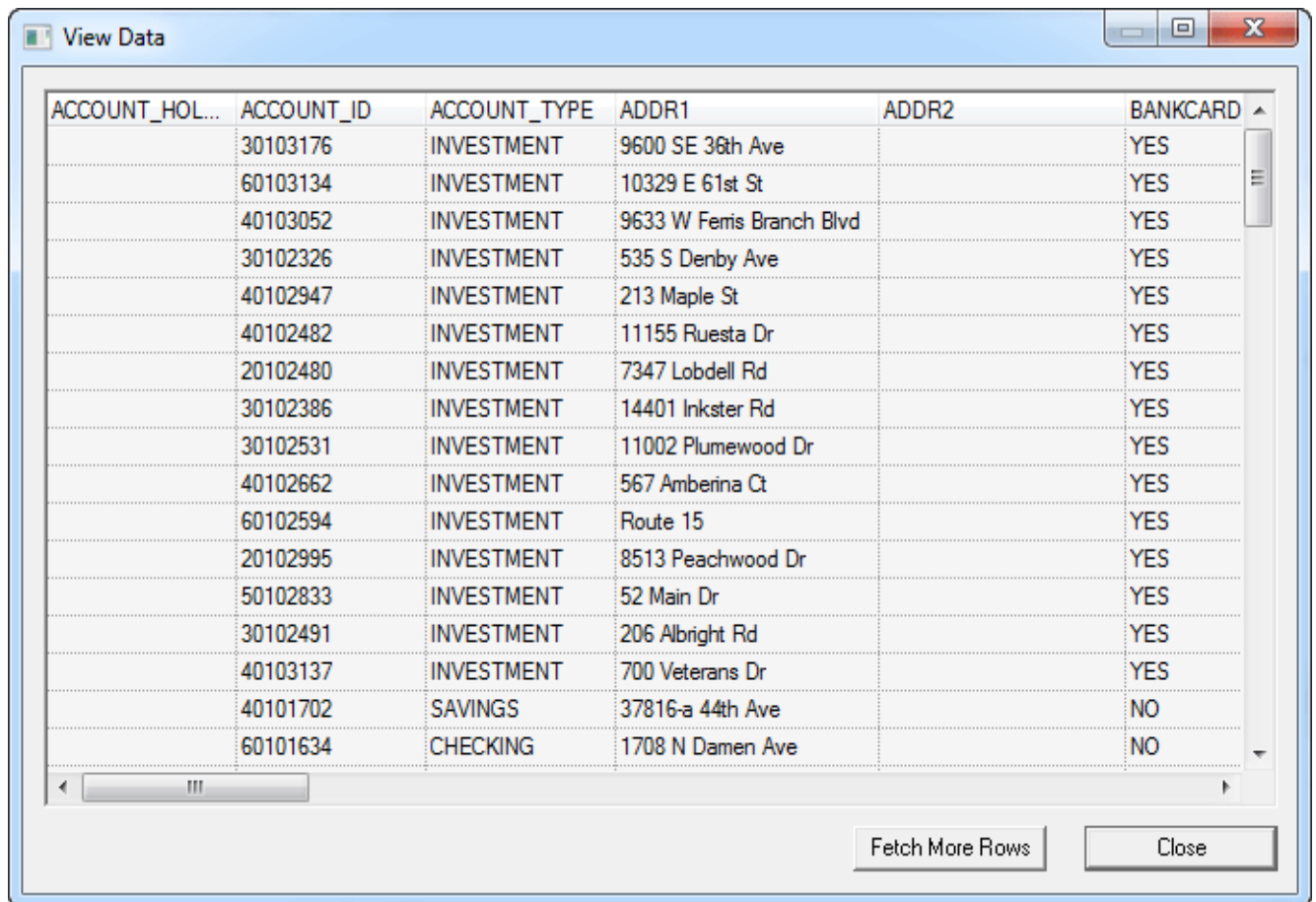
Run Validate Cancel Help

6. After the job runs, verify the data from the output stage:

- Click the **COMBINE\CST\INPUT** output database connector and click View Link\COMBINE\CST\INPUT Data



The data that is returned represents a combined customer input set for JKBANK1 & JK\BANK2 line-of-business customers in a non-standardized and un-cleansed form.



ACCOUNT_HOL...	ACCOUNT_ID	ACCOUNT_TYPE	ADDR1	ADDR2	BANKCARD
	30103176	INVESTMENT	9600 SE 36th Ave		YES
	60103134	INVESTMENT	10329 E 61st St		YES
	40103052	INVESTMENT	9633 W Ferns Branch Blvd		YES
	30102326	INVESTMENT	535 S Denby Ave		YES
	40102947	INVESTMENT	213 Maple St		YES
	40102482	INVESTMENT	11155 Ruesta Dr		YES
	20102480	INVESTMENT	7347 Lobdell Rd		YES
	30102386	INVESTMENT	14401 Inkster Rd		YES
	30102531	INVESTMENT	11002 Plumewood Dr		YES
	40102662	INVESTMENT	567 Amberina Ct		YES
	60102594	INVESTMENT	Route 15		YES
	20102995	INVESTMENT	8513 Peachwood Dr		YES
	50102833	INVESTMENT	52 Main Dr		YES
	30102491	INVESTMENT	206 Albright Rd		YES
	40103137	INVESTMENT	700 Veterans Dr		YES
	40101702	SAVINGS	37816-a 44th Ave		NO
	60101634	CHECKING	1708 N Damen Ave		NO

b. Click **Close** to close the window.

## Move data from Data Warehouse to Hadoop

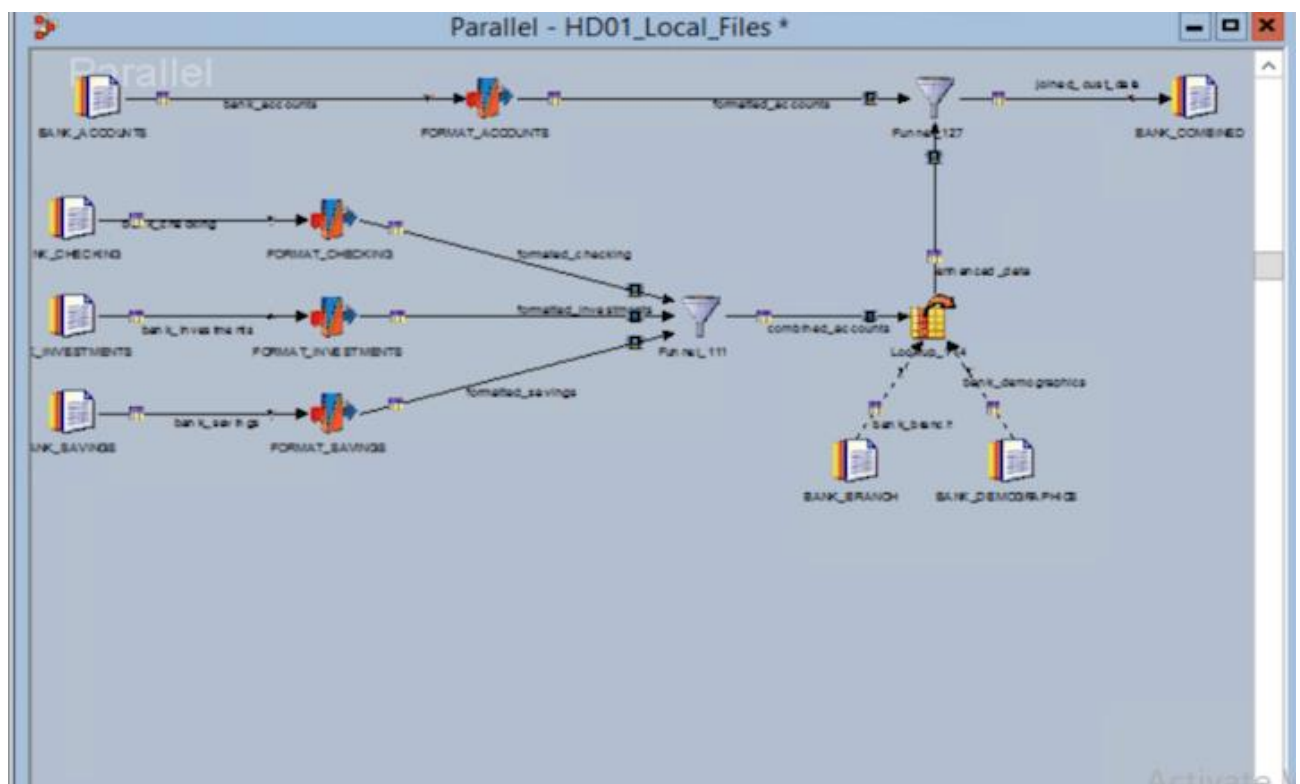
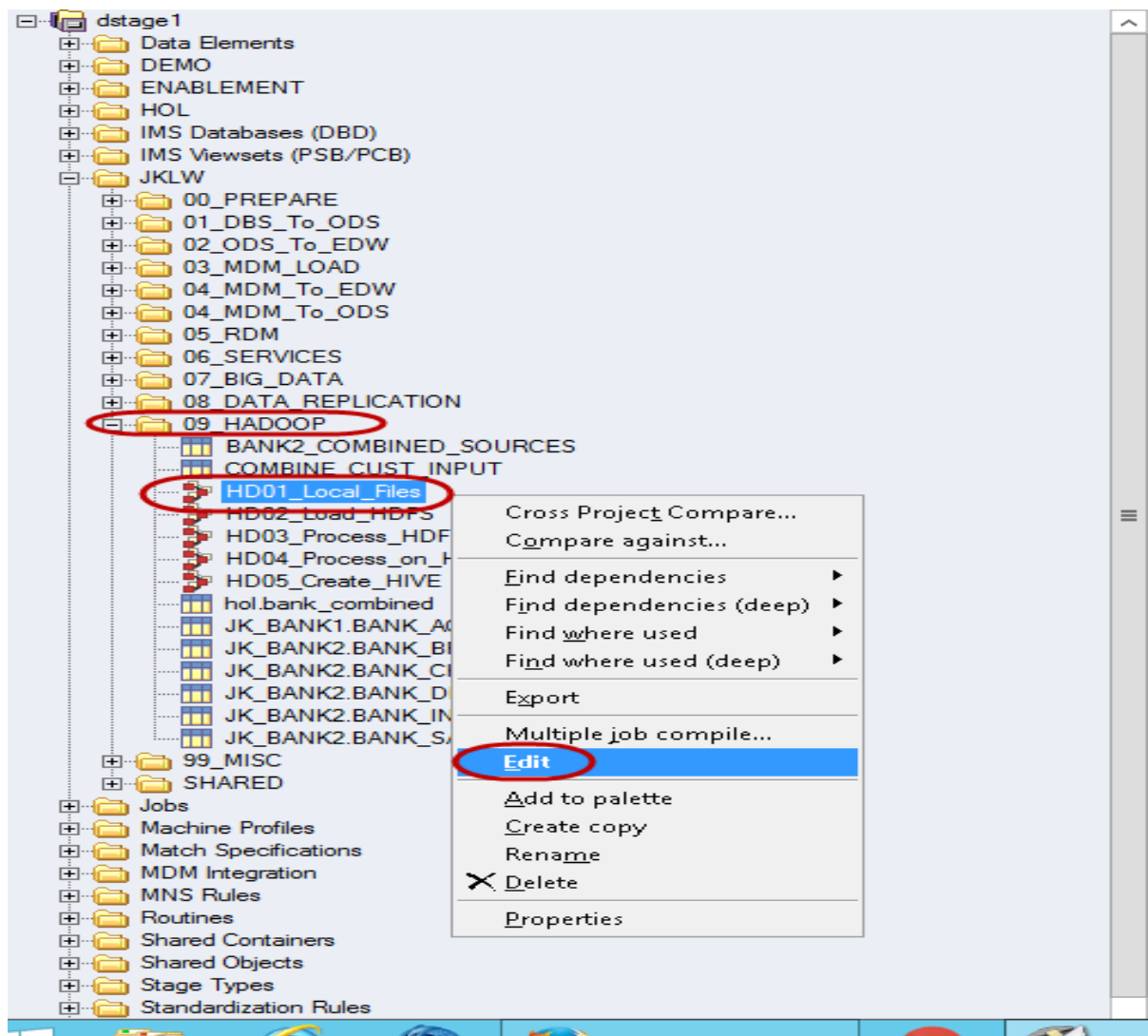
Explore how to use the scalability of the parallel engine and offload the processing from the enterprise data warehouse into a Hadoop cluster.

In this task, you complete these steps:

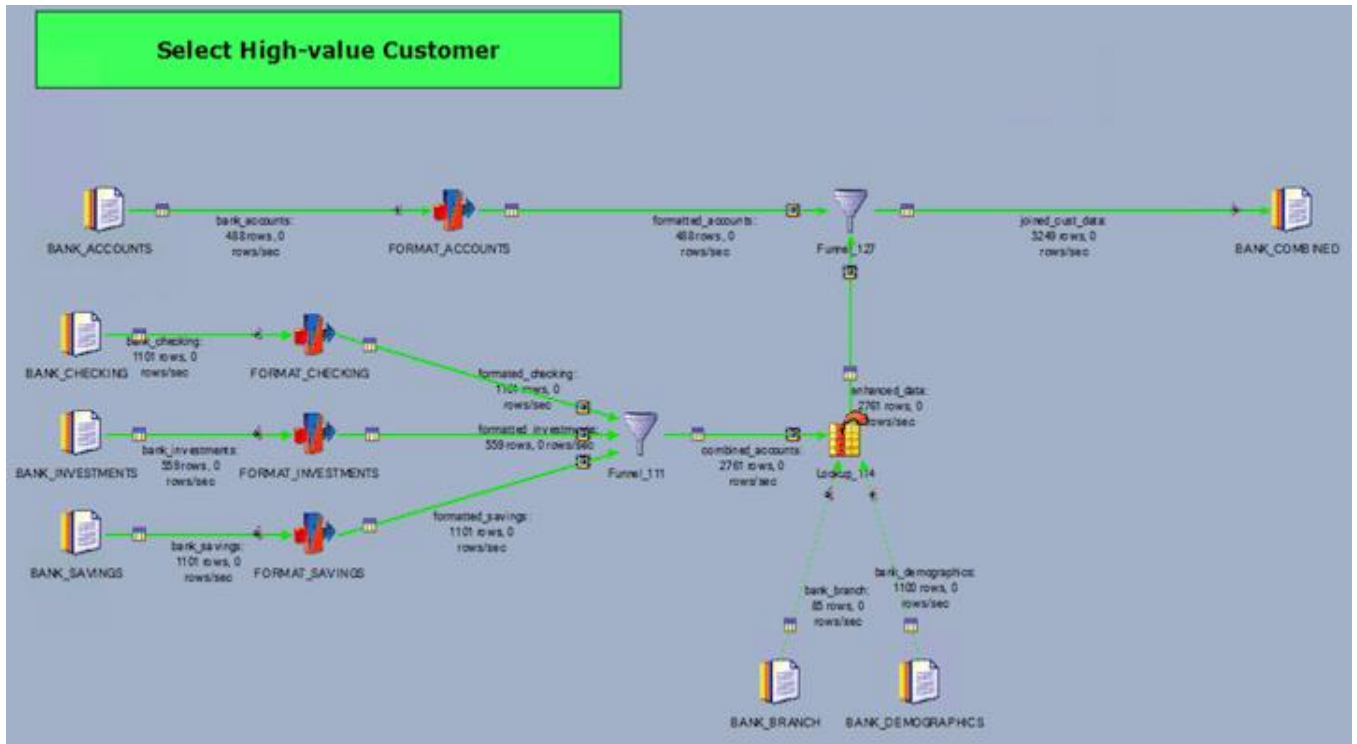
- Run a file processing job natively on the Linux host system
- Push the data to the Hadoop cluster

In this job, files that are based on Linux are ingested on the DataStage conductor node as the landing zone. After the files are ingested, DataStage processes the data files in a traditional DataStage configuration.

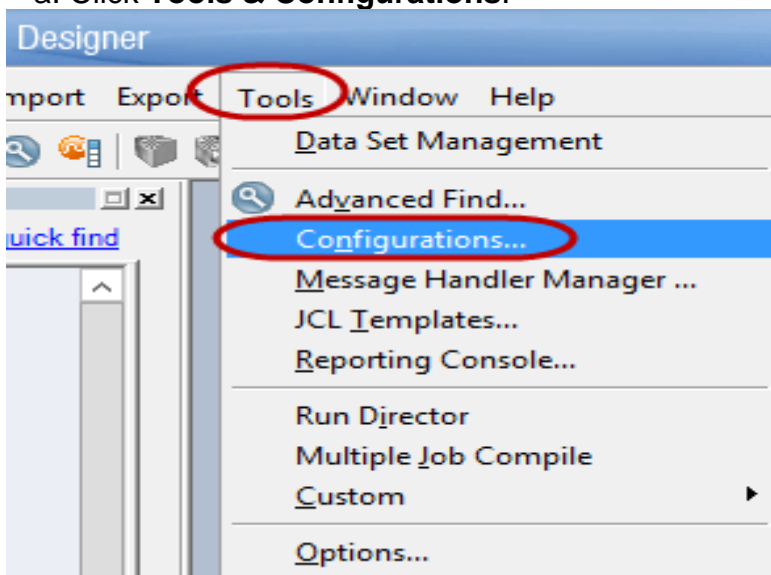
1. Click the **HD01\Local\Files** job and click **Edit**.



2. Compile the job by clicking the **Compile** icon on the toolbar.
3. Run the job by clicking the **Run** icon on the toolbar.
4. In the Job Run Options window, you can select the configuration file. For this demo, use the `default.apf` configuration file. Click **Run**.  
When the job is finished, all the job links turn green and show the number of rows on each link.

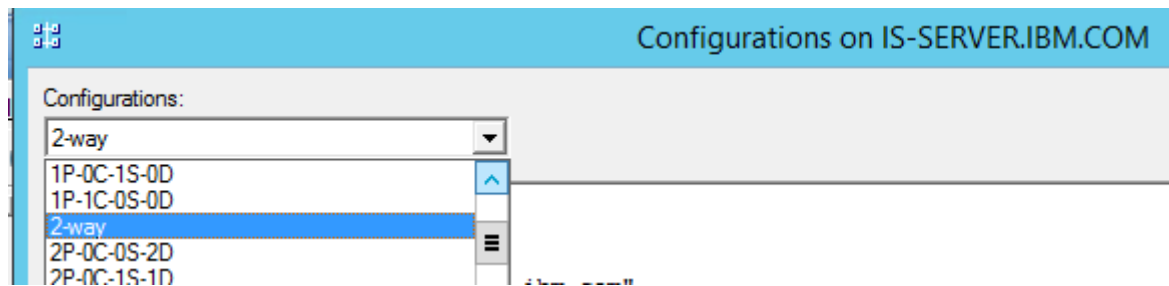


5. Review the default configuration file:
  - a. Click **Tools & Configurations**.



- b. From the Configurations list, select 2-way.





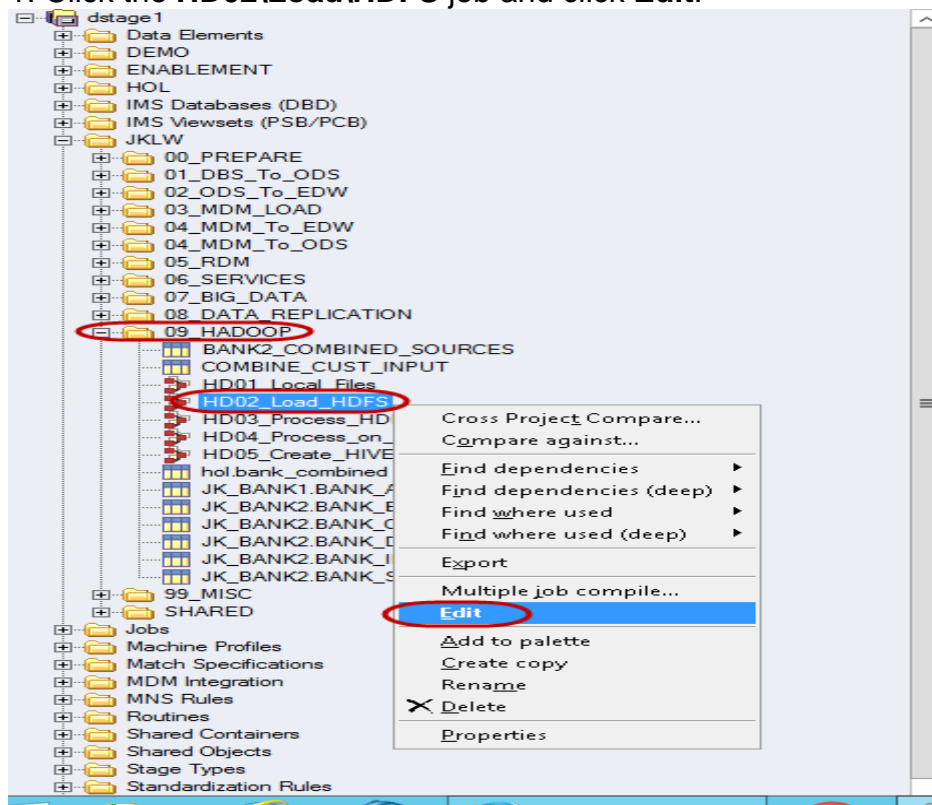
This file is a traditional configuration file, where the nodes and the location of the resources are named. A 2-node configuration file is shown, but you can have a 1-node configuration file.



c. Click Close.

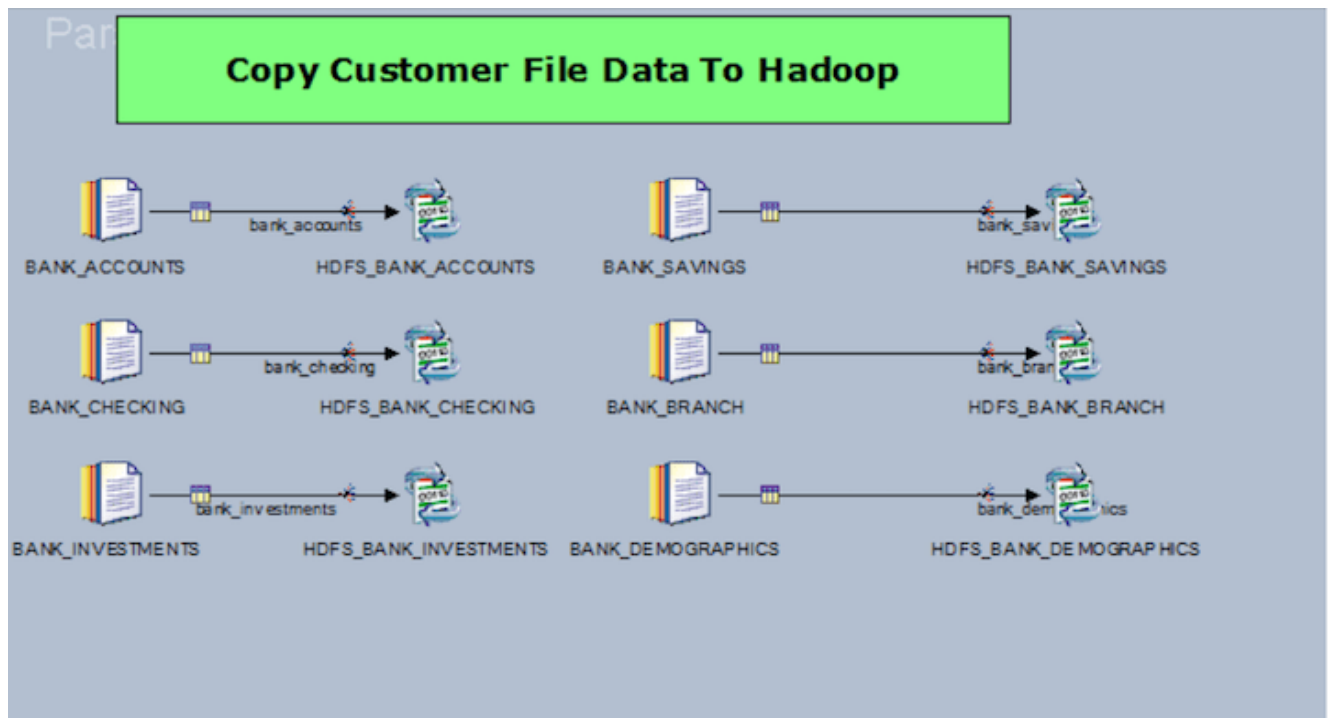
Load data into the Hadoop HDFS

1. Click the **HD02Load\HDFS** job and click **Edit**.

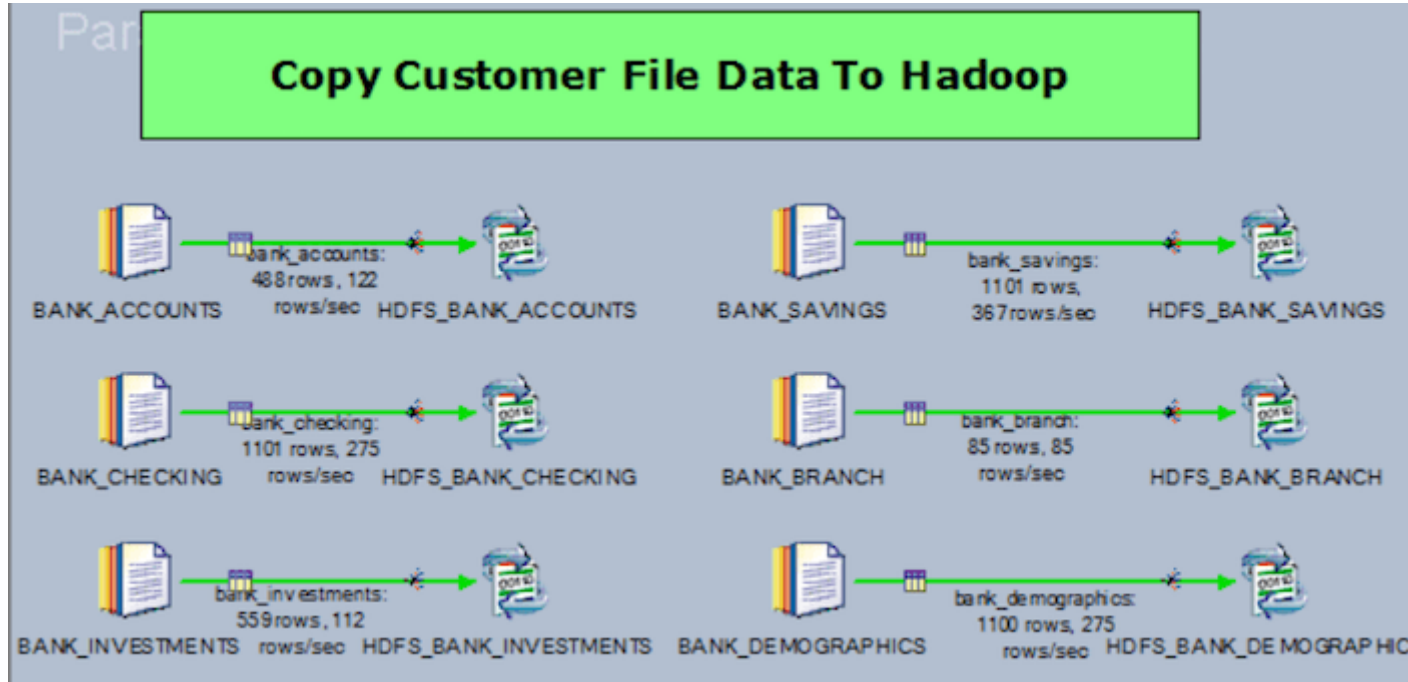


The HD02\_Load\_HDFS job loads your files onto HDFS for processing later.





2. Compile the job by clicking the Compile icon on the toolbar.
  3. Run the job by clicking the Run icon on the toolbar.
  4. In the Job Run Options window, you can select the configuration file. Use the `default.appt` configuration file. Click **Run**.
- When the job is finished, all the job links turn green and show the number of rows on each link.



Review a File Connector property

1. Click the **HDFS\BANK\ACCOUNTS** file connector.

The File Connector window shows how you assigned the connectivity attributes for the big data server. You can also see the file attributes for where and how you write the data. Your host and file path might be different.

Stage | Input |

Stage name  
HDFS\_BANK\_ACCOUNTS

General | Properties | Advanced |

File system	WebHDFS
Use custom URL	No
Use SSL (HTTPS)	No
Use Kerberos	No
Use keytab	No
Host *	is-server.ibm.com
Port	
Service principal	
User name *	dsadm
Password	
Keytab *	
Custom URL *	

Usage

Write mode	Write single file
File name *	/klw/HDFS_BANK_ACCOUNTS.TXT
If file exists	Overwrite file
Split file on key changes	No
Maximum file size	0
Force sequential	No
Cleanup on failure	Yes
File format	Comma-separated value (CSV)
Avro format properties	
ORC Settings	
Delimited format properties	
Encoding	
Include byte order mark	No
First row is header	Yes
Include data types	Yes
Field delimiter	,
Row delimiter	<NL>
Escape character	
Quotation mark	None
Null value	----
Field format	

OK Cancel Help

When you're finished reviewing the information, click **OK** to exit the properties window.

2. Review the data in the HDFS files system. The Ambari web console is open and the user admin is logged in.

3. In the Ambari Console, click **File View**.

Ambari - hdp

https://169.55.181.195:8080/#/main/dashboard/metrics

Ambari hdp 0 ops 0 alerts

Dashboard Services Hosts Alerts Admin admin

Metrics Heatmaps Config History

Metric Actions Last 1 hour

NameNode Heap 21%

HDFS Disk Usage 10%

NameNode CPU WIO n/a

DataNodes Live 1/1

0.21 ms

YARN Queue Manager

Files View

Hive View

SmartSense View

Tez View

Zeppelin View

HDFS

YARN

MapReduce2

Tez







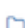
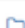
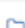
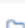
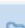
Hive

HBase




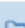

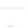
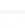
Pig

Sqoop

4. In the list of directories, click **user**.



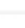
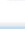
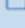

Name	Size	Last Modified	Owner	Group
 app-logs	--	2017-09-07 19:27	yarn	hadoop
 apps	--	2017-06-26 00:40	hdfs	hdfs
 ats	--	2017-06-25 07:23	yarn	hadoop
 biginsights	--	2017-06-26 01:08	hdfs	hdfs
 hdp	--	2017-06-25 07:23	hdfs	hdfs
 mapred	--	2017-06-25 07:23	mapred	hdfs
 mr-history	--	2017-06-25 07:23	mapred	hadoop
 spark-history	--	2017-08-04 15:49	spark	hadoop
 spark2-history	--	2018-03-26 15:22	spark	hadoop
 tmp	--	2018-03-24 15:19	hdfs	hdfs
 user	--	2017-08-18 10:16	hdfs	hdfs

5. Click the **dsadm** directory.

 ambari-qa	--	2017-08-04 13:39	ambari-qa	hdfs
 as_user	--	2017-07-04 17:00	as_user	hdfs
 biadmin	--	2017-08-18 09:54	admin	hdfs
 bigsql	--	2017-06-26 17:41	bigsql	hdfs
 dsadm	--	2018-03-08 22:19	dsadm	dsta
 hbase	--	2017-06-25 07:23	hbase	hdfs
 hcat	--	2017-06-26 00:40	hcat	hdfs
 hive	--	2017-08-04 12:51	hive	hdfs

6. Click the **jklw** folder.

The files that you loaded into the HDFS file system are in the folder.

 hdfs1_yarn.txt	0.1 kB	2018-03-08 22:19	hdfs	dsta
 hdfs1_yarn_debug.txt	0.1 kB	2017-09-28 13:01	hdfs	dsta
 insurance	--	2018-02-28 11:26	dsadm	dsta
 jklw	--	2018-03-26 12:48	dsadm	dsta
 my_hdfs_yarn_txt	693.9 kB	2018-03-08 22:19	hdfs	dsta
 sample	--	2018-03-07 19:01	dsadm	dsta

For the purposes of this demo, you look at the HDFS\_BANK\_ACCOUNTS.TXT file.

7. Click **Open** on the top menu bar.

Ambari hdp 0 ops 0 alerts Dashboard Services Hosts Alerts

/ > user > dsadm > jklw 1 Files, 0 Folders selected

Open Rename Permissions Delete Copy Move Download concatenate

Name	Size	Last Modified	Owner
HDFS_BANK_ACCOUNTS.TXT	117.4 kB	2018-03-28 20:04	dsadm
HDFS_BANK_BRANCH.TXT	4.2 kB	2018-03-28 20:04	dsadm

8. Review the data in the file. When you're finished, click **Cancel**.

File Preview

/user/dsadm/jklw/HDFS\_BANK\_ACCOUNTS.TXT

```

"RECORD_ID:VarChar(255)","SS_NUM:VarChar(255)","NAME:VarChar(255)","ADDR1:Var
"AV000065","458774796","Etta A Metheny","9600 SE 36th Ave","","FORT MEADE","F
"AV000026","403055607","Dwayne A Marquez","10329 E 61st St","","HOLLYWOOD","F
"AV0000221","401527635","Helen C Boe","535 S Denby Ave","","PENSACOLA","FL","
"AV0000842","734470487","Rodeny I Wigle","213 Maple St","","SEQUATCHIE","TX",
"AV0000692","580374025","Donovan U Marchetti","750 Greenlee Rd","","MOUNT JUL
"AV0000733","412271519","Alan G Hanley","27858 Palmetto Ridge Dr","","AUSTIN"
"AV0000982","570815245","Vera V Eisenberg","317 Road M56 St","","CLEVELAND",
"AV0000679","448876507","Butch A Shabashevich","650 Jenkins Rd","","MOUNT PLE
"AV0000323","445288231","Brendan Tracey","102 S 264th St","","FARGO","NC","72
"AV0000931","700128845","Robert 'E' Gooden","11 Kathleen Way","","POLKTON","N
"AV0000791","664775365","Gael W Rubel","Hwy 41 N","","CHINA GROVE","NC","7532
"AA00000048","528074458","Kaye X Riley","1747 Le Flore Dr","","BANKS","AL","3
"AA00000202","999-99-9999","Eric X Hirahara","6337 99th St E","","BANKS","AL"
"AA00000137","419029160","Moshe D Mercadante","900 4th St NE","","GOSHEN","AL
"AA00000082","264246606","Mack W Perreault","39737 Road 274","","MONTG","AL",
"AA00000118","217661742","Daryle F Anchondo","9704 Snra Hrdn Sprgs Rd","","MO

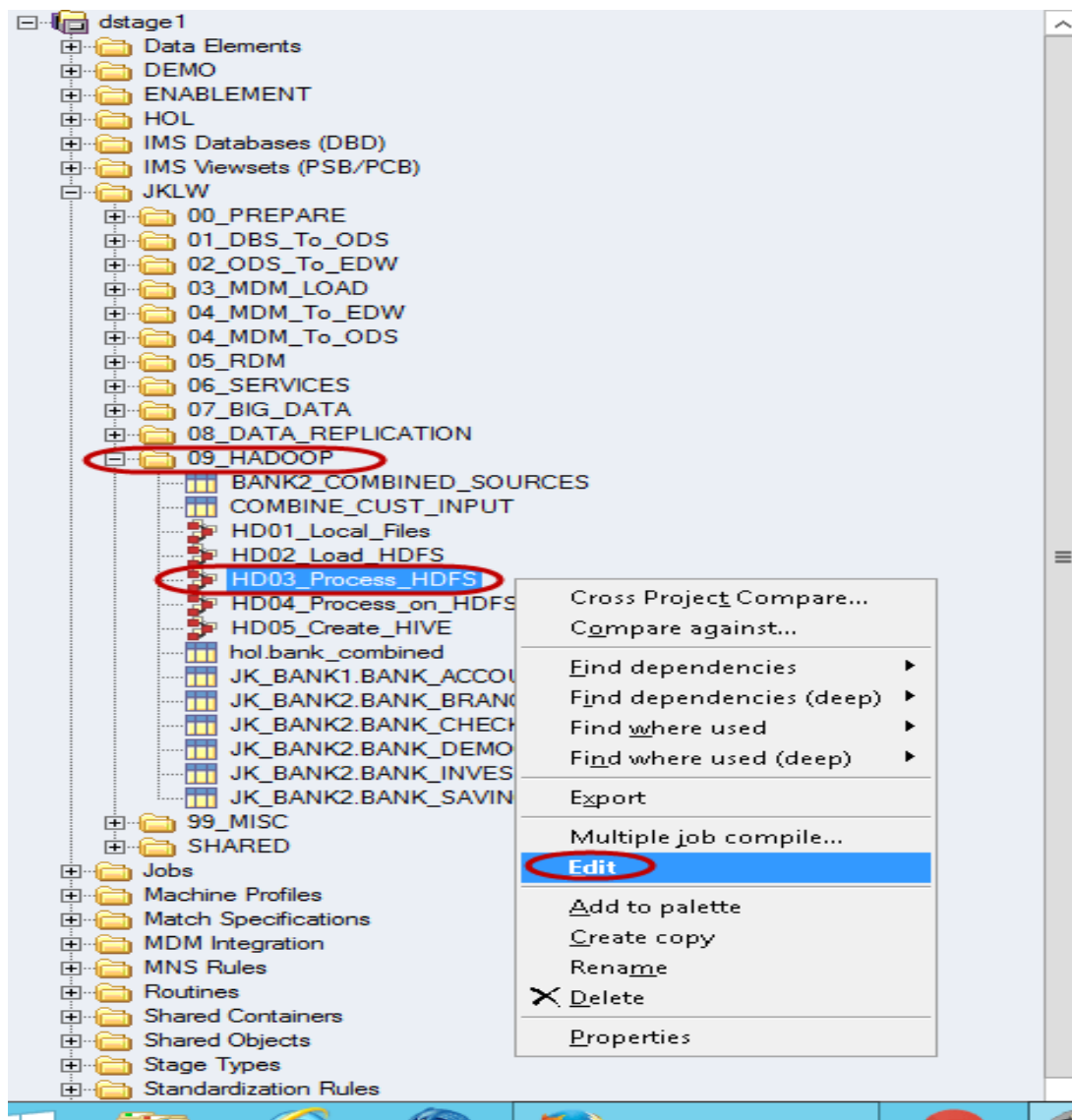
```

Cancel Download

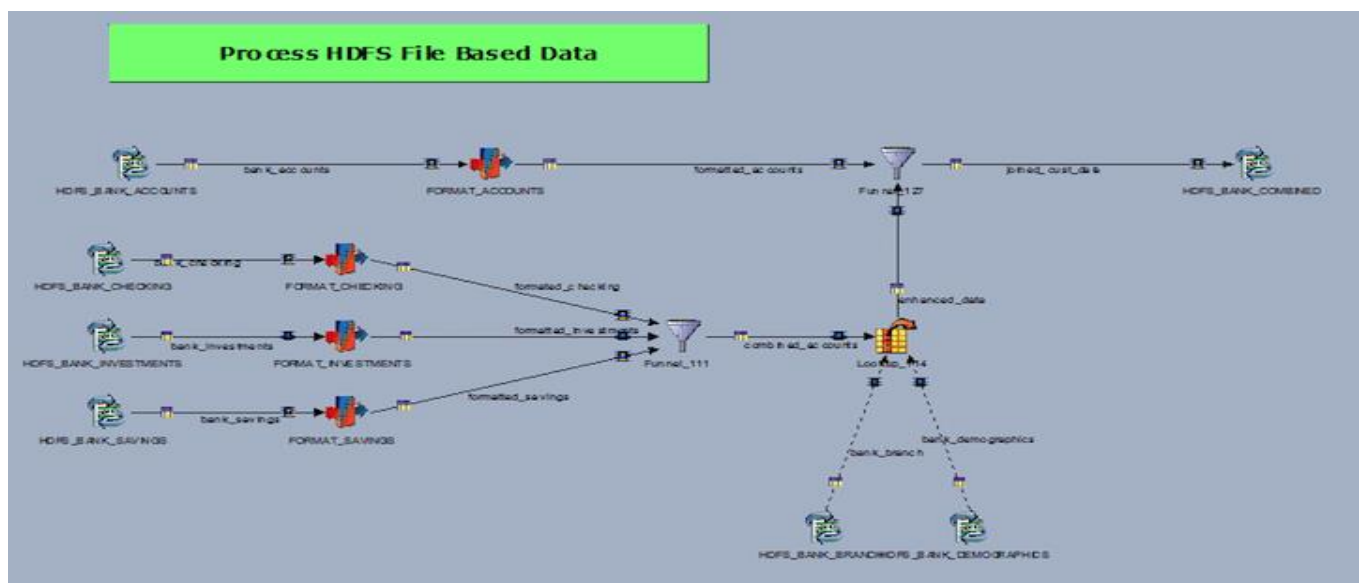
Process the Hadoop files

Process the HDFS-file-based data in a traditional DataStage configuration.

1. Click the **HD03\Process\HDFS** job and then click Edit.



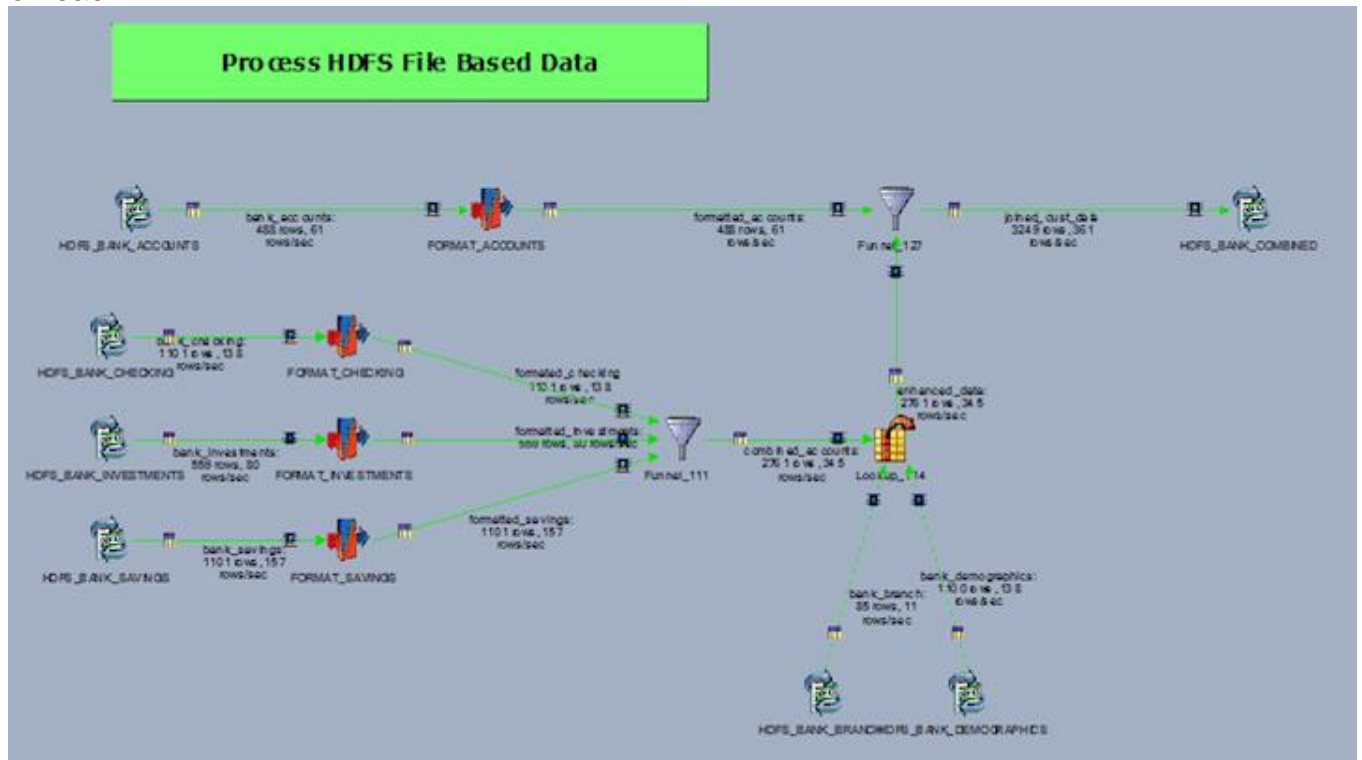
The HD03\Load\HDFS job reads data from your Hadoop HDFS file system, processes the data, and writes the consolidated customer data back to the Hadoop HDFS file system.





2. Compile the job by clicking the **Compile** icon on the toolbar.
3. Run the job by clicking the **Run** icon on the toolbar.
4. In the Job Run Options window, you can select the configuration file. Use the `default.apt` configuration file. Click **Run**.

When the job is finished, all the job links turn green and show the number of rows on each link.



5. Return to the HDFS File view by clicking the Maximize Window icon.

📄 HDFS_BANK_ACCOUNTS.TXT	117.4 kB	2018-03-28 20:04
📄 HDFS_BANK_BRANCH.TXT	4.2 kB	2018-03-28 20:04
📄 HDFS_BANK_CHECKING.TXT	174.3 kB	2018-03-28 20:04
📄 HDFS_BANK_COMBINED.TXT	832.6 kB	2018-03-28 13:06
📄 HDFS_BANK_DEMOGRAPHICS.TXT	137.0 kB	2018-03-28 20:04
📄 HDFS_BANK_INVESTMENTS.TXT	83.6 kB	2018-03-28 20:04
📄 HDFS_BANK_SAVINGS.TXT	161.4 kB	2018-03-28 20:04

Notice the time that the HDFS\BANK\COMBINED\TEXT file was written onto the Hadoop file system.

6. Click the **refresh icon** at the top of the File view to reload the file information.

<span>Open</span> <span>Rename</span> <span>Permissions</span> <span>Delete</span> <span>Copy</span> <span>Move</span> <span>Download</span>		
Name	Size	Last Modified
↩		
📄 HDFS_BANK_ACCOUNTS.TXT	117.4 kB	2018-03-28 20:04
📄 HDFS_BANK_BRANCH.TXT	4.2 kB	2018-03-28 20:04
📄 HDFS_BANK_CHECKING.TXT	174.3 kB	2018-03-28 20:04
📄 HDFS_BANK_COMBINED.TXT	832.6 kB	2018-03-28 20:15
📄 HDFS_BANK_DEMOGRAPHICS.TXT	137.0 kB	2018-03-28 20:04
📄 HDFS_BANK_INVESTMENTS.TXT	83.6 kB	2018-03-28 20:04

The file time stamp is updated after the job ran.

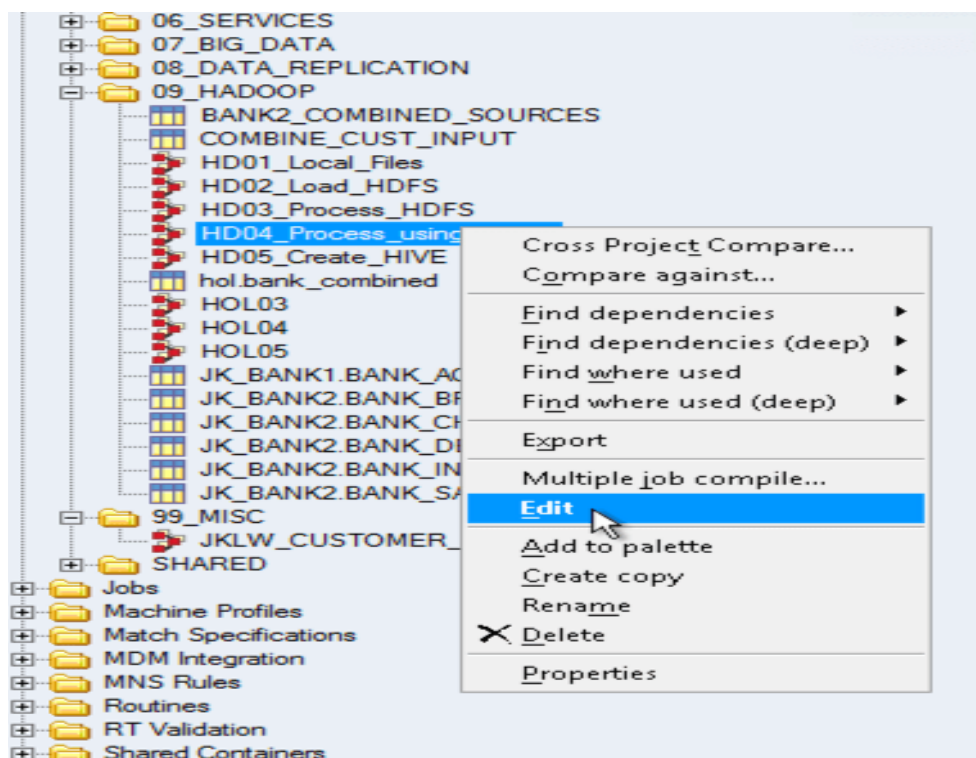
Close the File view by clicking the highlighted file.

7. Return to DataStage and close the job.

## Run ETL processing inside Hadoop by using YARN

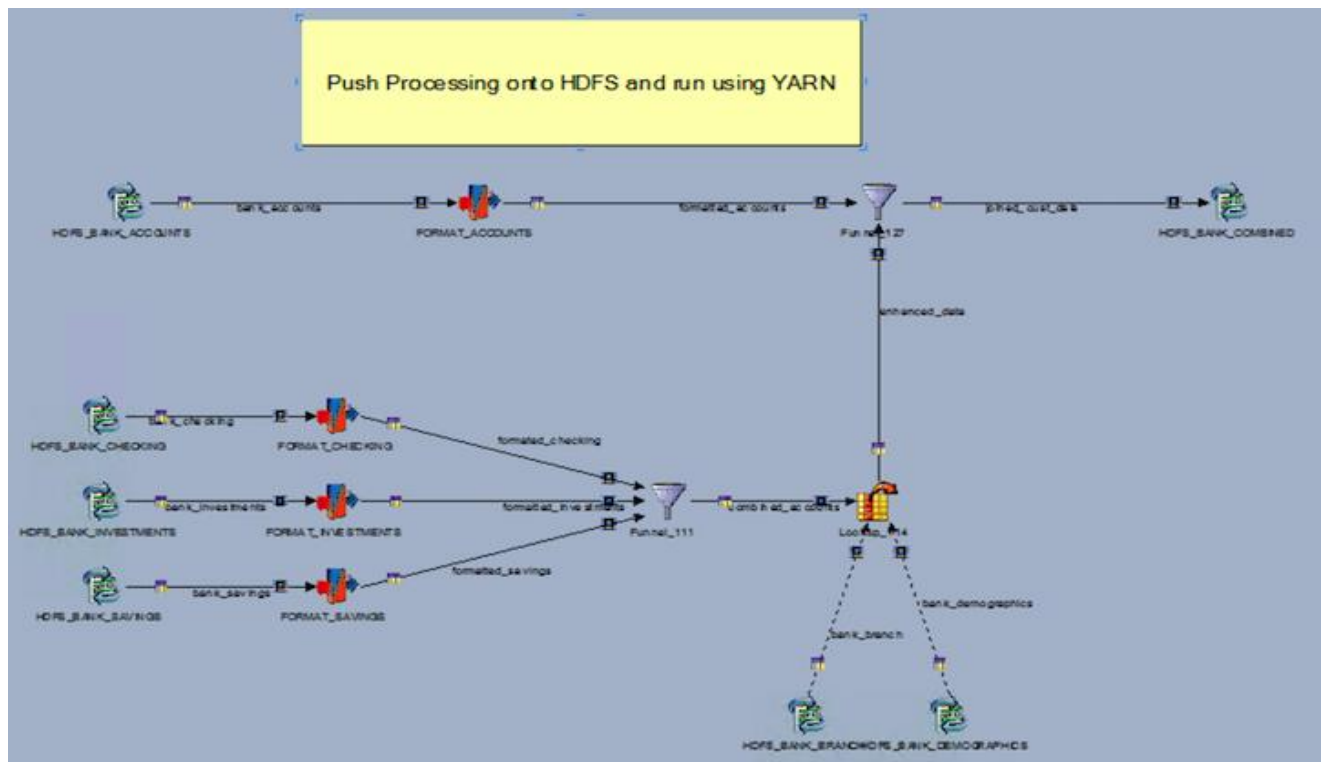
Process the HDFS-file-based data by pushing the processing to the Hadoop cluster and allowing the processing to run by using YARN.

1. Right-click the HD04\Process\using\YARN job and click **Edit**.



The HD04\Process\on\HDFS\using\YARN job reads data from your Hadoop HDFS file system, processes the data, and writes the consolidated customer data back to the Hadoop HDFS file system entirely in the Hadoop cluster.





2. Compile the job by clicking the **Compile** icon on the toolbar.

3. Run the job by clicking the **Run** icon on the toolbar.

In the Job Run Options window, you can see that you have two properties.

The first property points to the `yarnconfig.cfg` file. This file indicates to DataStage which settings it needs to communicate with and run the process on YARN.

The second property points to the configuration file so that DataStage detects how many nodes to run and where the resources are.

4. Before you click **Run**, learn about the properties:

- \* On the server, the `yarnconfig.cfg` file contains a number of settings to tune how DataStage processes are run with YARN. Two important settings are `APT\_YARN\_MODE` and `APT\_YARN\_USE\_HDFS`.

- \* The `APT\_YARN\_MODE` setting tells the engine where to run. A value of `false` tells the engine to run normally on the host system. A value of `true` tells the engine to hand the process to YARN for processing on the Hadoop cluster.

- \* The `APT\_YARN\_USE\_HDFS` setting tells the engine whether the data resources are being written to the local file system or the HDFS file system. A setting of `false` indicates that the data is written to the local file system. A setting of `true` indicates that the data is being written to the HDFS file system.

- \* You can use many other settings to further tune the interaction. Each setting is described in the `yarnconfig.cfg` file.

```

## Licensed Materials - Property of IBM
## (c) Copyright IBM Corp. 2015
#
# DataStage PX Yarn Configuration
# =====
# Lines in this file are either comments, introduced by a # sign like this,
# or of the form "key=value". Key lines may be commented out below.
#
# IMPORTANT:
# Ensure when making changes to this file that it is saved with the encoding set to
# UTF-8. Please be aware if the encoding isn't set to UTF-8 this may produce undesired
# behaviour.

APT_YARN_MODE=true
# If defined and set to 1 or true runs the given PX job on
# the local Hadoop install in YARN mode.

APT_YARN_CONTAINER_VCORES=0
# Defines the number of virtual cores that the containers will request to run
# PX Section Leader and Player processes in.
# The default is 0 which means "Don't set it".

APT_YARN_CONTAINER_SIZE=64
# Defines the size in MBs of the containers that will be requested to run
# PX Section Leader and Player processes in.
# The default is 64MB if not set.

APT_YARN_CONTAINER_SIZE_AUTO=false
# When defined will automatically use the estimated container size for the largest partition
# as the container size if it is larger than the set container size.
# It accepts a value of true or false.

APT_YARN_BASE_PROCESS_SIZE=16
# Defines the base process size in MB to be used when estimating the size of containers
# to request from YARN. The default value for this is 16MB.

APT_YARN_ALLOCATION_TIMEOUT=180
# Specifies the amount of time in seconds to wait for allocations of containers

```

\* When parallel jobs run on Hadoop, the jobs request a set of containers from YARN. The containers represent the resources that the job was allocated. Each resource has a designated amount of virtual CPU and memory for each container. The number of containers that are requested is equal to the number of logical nodes that are defined in the `APT\_CONFIG\_FILE` file.

\* These files can be configured in one of three ways: static, dynamic, and mixed.

\* A static file looks like a regular `APT\_CONFIG\_FILE`, except the resource disk can be either on the local file system or in HDFS depending on the `APT\_YARN\_USE\_HDFS` setting in the `yarnconfig.cfg` file.

```

{
  node "node1"
  {
    fastname "mymachine.domain.com"
    pools ""
    resource disk "/mydisk/tmp" {pools ""}
    resource scratchdisk "/myscratch/tmp" {pools ""}
  }
}

```

\* A dynamic configuration file uses the same format as the static Information Server configuration files, which assign fixed nodes to the job. However, a dynamic configuration file uses a fastname value of ``$host``, as opposed to a static configuration file that usually contains a host name. One node in the configuration file must contain the engine tier node, but this node can be defined with a conductor node pool if you don't want to run data processing on the conductor node.

```

{ node "node0"
  {
    fastname "the-engine-tier-machine.domain.com"
    pools "conductor"
    resource disk "/sandbox/bsmith/tmp" {pools ""}
    resource scratchdisk "/scratch" {}
  }
  node "node1"
  {
    fastname "$host"
    pools ""
    resource disk "/sandbox/bsmith/tmp" {pool ""}
    resource scratchdisk "/sandbox/bsmith/tmp" {pools ""}
    instances 30
  }
}

```

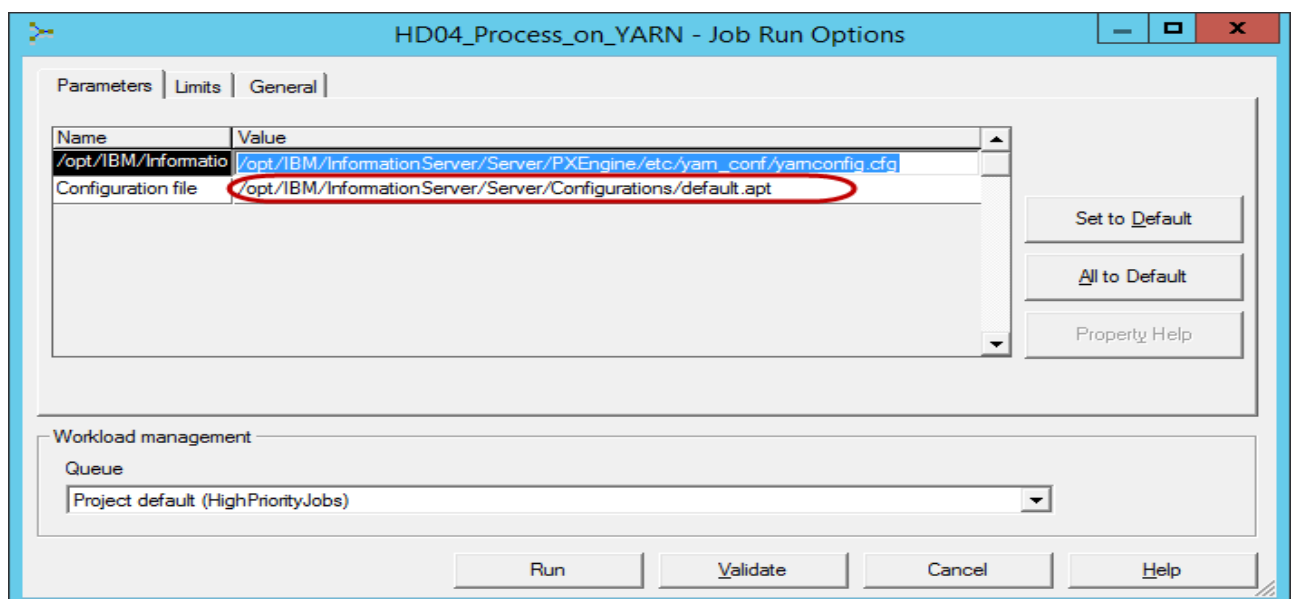
\* A mixed configuration file contains a mix of static host names (an actual host name) and dynamic host names (fastname ``$host``). The following configuration file specifies that the first 30 nodes are defined by YARN. Then, the next 10 nodes are run on machineA, and the 41st node runs on machineB:

```

{
  node "node0"
  {
    fastname "the-engine-tier-machine.domain.com"
    pools "conductor"
    resource disk "/sandbox/bsmith/tmp" {pools ""}
    resource scratchdisk "/scratch" {}
  }
  node "node1"
  {
    fastname "$host"
    pools ""
    resource disk "/mydisk1/tmp" {pools ""}
    resource scratchdisk "/myscratchdisk1/tmp" {pools ""}
    instances 30
  }
  node "node31"
  {
    fastname "machineA.domain.com"
    pools ""
    resource disk "/mydisk2/tmp" {pools ""}
    resource scratchdisk "/myscratchdisk2/tmp" {pools ""}
    instances 10
  }
  node "node41"
  {
    fastname "machineB.domain.com"
    pools ""
    resource disk "/mydisk3/tmp" {pools ""}
    resource scratchdisk "/myscratchdisk3/tmp" {pools ""}
  }
}

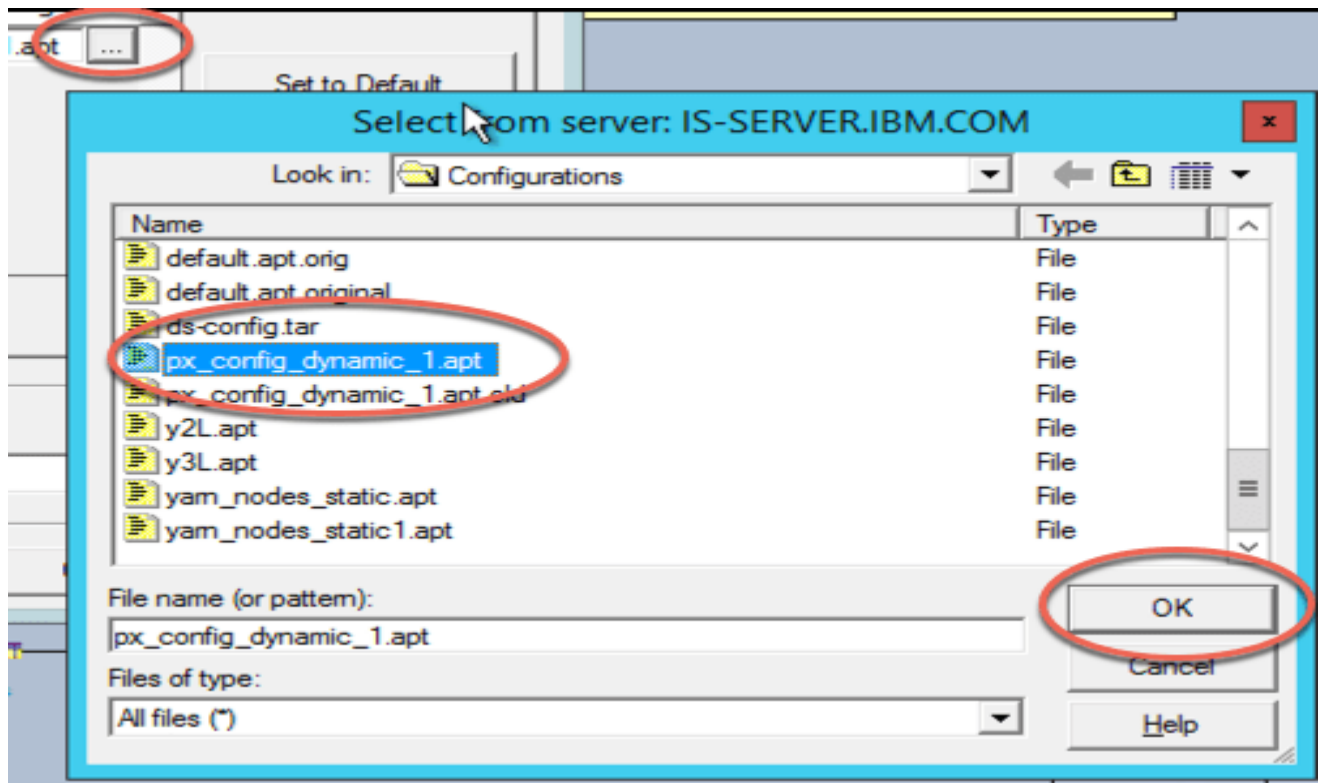
```

5. In the Job Run Options window, use a dynamic configuration file for your processing. Click the field that defines the configuration file to show the **Options** icon.

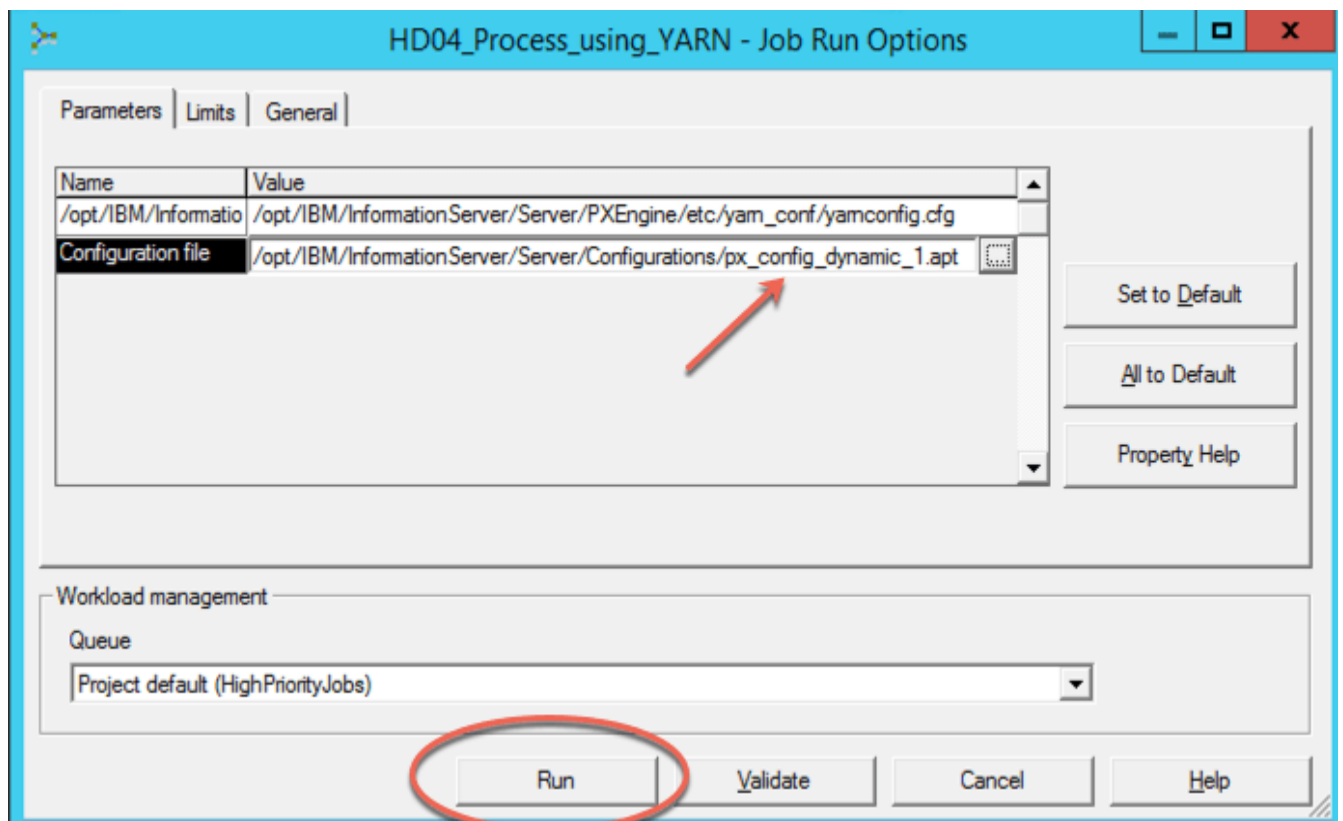


A file-browser window is displayed and the `px\_config\_dynamic\_1.apl` file is selected.

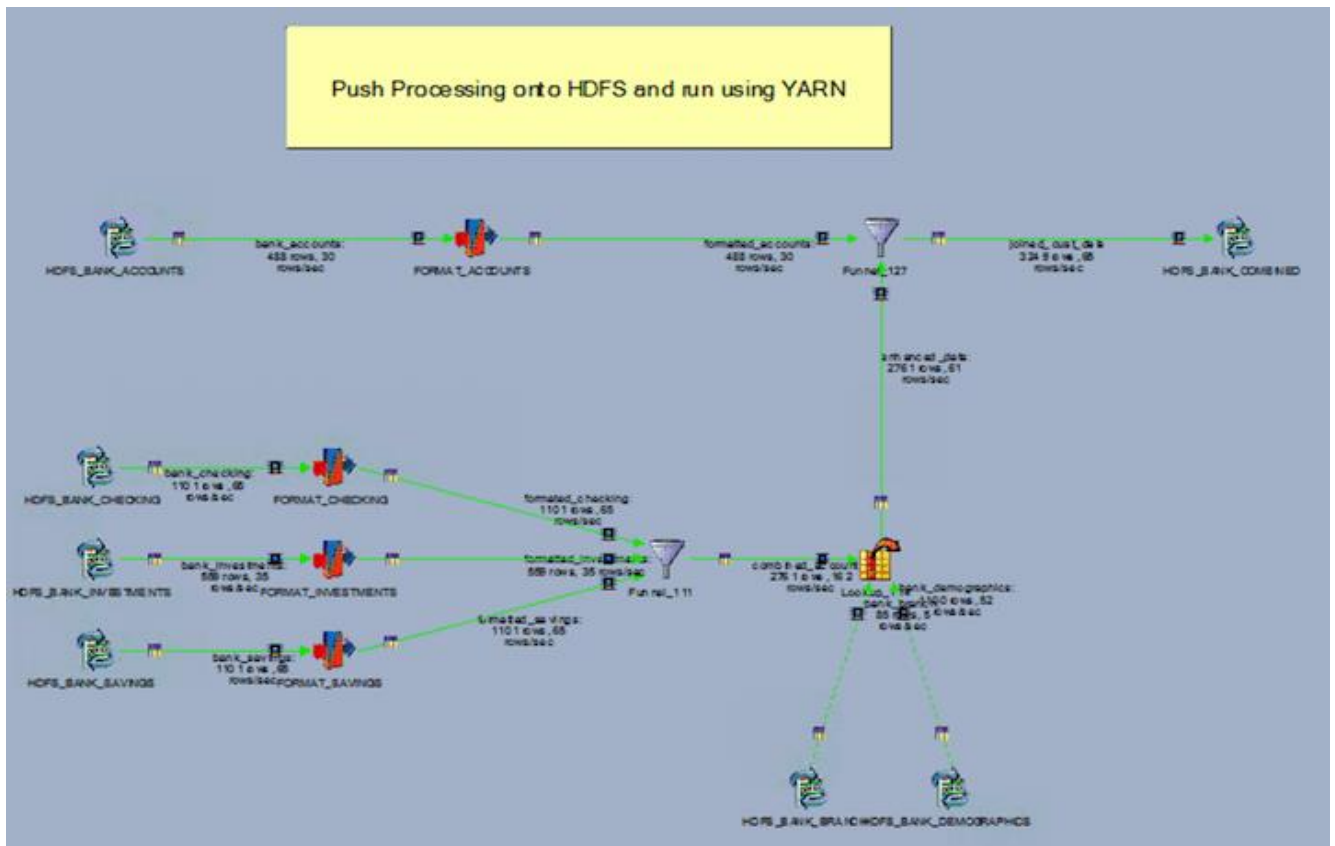
6. Click **OK**.



7. Ensure that you're using the correct configuration file and click **Run**.



When the job is finished, all the job links turn green and show the number of rows on each link.



8. Click the center of the lower part of the page to raise the job log window.
9. Look for these entries:

```
Starting Job HD04_Process_using_YARN. (...)
Environment variable settings: (...)
Parallel job initiated
OSH script (...)
Parallel job default NLS map UTF-8, default locale OFF
main_program: IBM InfoSphere DataStage Enterprise Edition 11.5.0.8169 (...)
main_program: The open files limit is 1024; raising to 4096.
main_program: conductor uname: -s=Linux; -r=2.6.32-642.13.1.el6.x86_64; -v=#1 SMP Wed Nov 23 16:03:01 EST 2016; -n=is-serve...
main_program: orchgeneral: loaded (...)
main_program: Parallel Engine running in YARN execution mode.
main_program: IPv6 isn't currently supported by Hadoop/YARN and the required environment variable APT_USE_IPV4 is not set. It ...
HDFS_BANK_ACCOUNTS: Accessing file via WebHDFS file system.
HDFS_BANK_CHECKING: Accessing file via WebHDFS file system.
HDFS_BANK_INVESTMENTS: Accessing file via WebHDFS file system.
HDFS_BANK_SAVINGS: Accessing file via WebHDFS file system.
HDFS_BANK_DEMOGRAPHICS: Accessing file via WebHDFS file system.
HDFS_BANK_BRANCH: Accessing file via WebHDFS file system.
HDFS_BANK_COMBINED: The connector was configured to run in parallel on 3 nodes, but the Read/Write mode is Read/Write sin...
HDFS_BANK_COMBINED: Accessing file via WebHDFS file system
main_program: APT configuration file: /opt/IBM/InformationServer/Server/Configurations/px_config_dynamic_1.apt (...)
```

## Summary

You ran a traditional Data Warehouse ETL job and moved data from Data Warehouse to Hadoop. Then, you ran a Data Warehouse ETL by using Hadoop data and ran ETL processing inside Hadoop by using YARN.