

InfoSphere Information Server

Troubleshooting performance problems in Information Analyzer Column Analysis version 8



© 2012 IBM Corporation

This presentation will explain how to troubleshoot performance problems in Information Analyzer Column Analysis for version 8 and higher.

Objectives

- Performance factors
- Collecting runtime data
- Analyzing data
- Troubleshooting and diagnosing
- Recommendations

The objective of this presentation is to explain the most important factors that impact Column Analysis performance. This presentation explains how to collect runtime data to identify bottlenecks, how to troubleshoot and diagnose the root causes, and how to resolve and prevent these problems. In this presentation, Information Analyzer is referred to as IA and Column Analysis is referred to as CA.

Scope

- CA consists of three steps
 - Pre-process
 - Main-process
 - Post-process
- iasServer*.log lists the time each step takes
 - Example:

```
-----
- | Job Execution Status -> ColumnAnalysisTask61084527 |
- | JOB TYPE           : com.ascential.investigate.ca.job.BaseProfileJob |
- | JOB PROCESSOR      : com.ascential.investigate.utils.jobs.osh.PXProcessor |
- |-----|
- | Data Stage> ANALYZERPROJECT@SAWCHUCK:31538, credentialMapping=false |
- | Scratch Data Source> iadb@sawchuck.swg.usma.ibm.com:50000 (DB2) |
- | [] pre-process completed (0 sec) |
- | [] main-process completed (610 sec) |
- | [] post-process completed (63 sec) |
- |-----|
- | Job run *SUCCESS* in 11 minute(s) 17 seconds. |
- |-----|
```

- Focuses on troubleshooting main-process

3

Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

CA consists of three main steps, which are started sequentially. The first step is **pre-process**, which starts CA and prepares the execution of the next steps. The next one is the **main-process**, which runs the bulk of the data extraction, data analysis and creation of results. The last step is the **post-process**, which runs additional processes, such as Enhanced Data Classification, or EDC, depending on the options selected for CA.

The time each step takes can be obtained from the file iasServer*.log which is located under WebSphere/AppServer/profiles/<IS_HOME_PROFILE>. This slide displays an example of how the times are shown in this file.

This presentation will focus only on analyzing performance problems that affect the main-process, which is the longest process and the most likely step to experience performance issues. If you observe performance problems involving the pre-processing step, engage IBM Support for assistance. If you observe problems with the post-processing step, check the section “Troubleshooting: XMETA with EDC on”.

Performance factors

- Important factors for CA performance

Factor	Description
Network Connectivity/ Latency	Speed at which the different servers involved in CA communicate with each other.
Database Performance	Speed at which the source Database and the IADB Database process data. Note: If EDC is ON then XMETA can also be a factor
Engine Performance	Speed at which the Engine analyzes the data
Data Attributes	Particular attributes of the Data being analyzed
Patch Level	Level of Patches installed

4

Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

The performance of the CA main-process can be affected by several factors. This slide displays the most common and important factors.

The first factor, **Network Connectivity/Latency**, refers to the speed at which the different servers involved in CA communicate with each other. Depending on your configuration, there can be up to four different servers working together: the engine, the XMETA repository, the IADB repository, and the Source Database Server.

The next factor, **Database Performance**, refers to the speed at which the Source Database and the IADB Database process data. It is important to note that if EDC is used, the performance of the XMETA database can also be a factor.

The next factor, **Engine Performance**, refers to the speed at which the engine analyzes data.

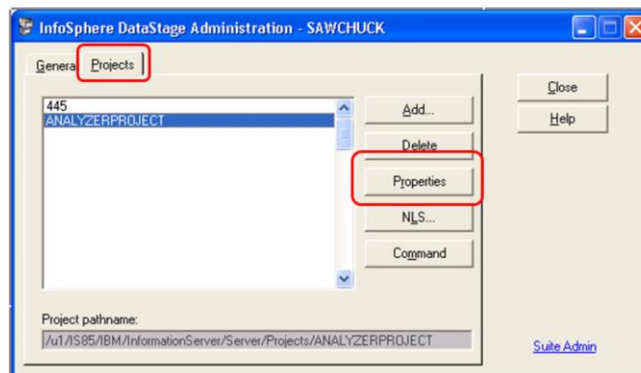
The next factor, **Data Attributes**, refers to the particular attributes of the data being analyzed. Different data types or large data volumes can take longer times to be analyzed.

The last factor is the **Patch Level**. If you are experiencing performance problems always make sure that you are running the latest Fix Pack and Rollup Patch level available for all your tiers.

This presentation will show how to diagnose and identify which of these factors can be impacting the performance of CA, with the exception of the Patch Level factor, as this varies over time and per version.

Collecting runtime data (1 of 13)

- Open DataStage® Administrator
- Login as DataStage administrator user
- Click Projects tab
 - Select ANALYZERPROJECT
- Click Properties



5

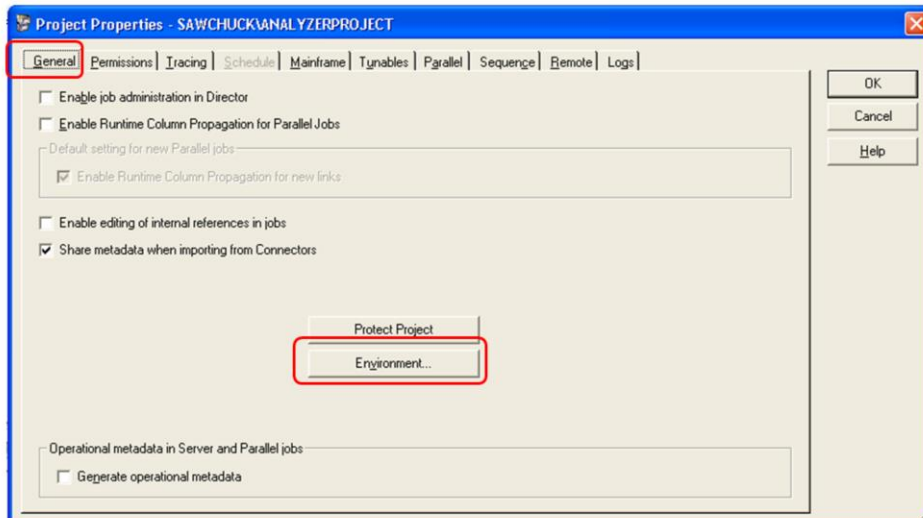
Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

To troubleshoot performance issues, you first need to collect runtime data. This is done by setting debug environment variables at the project level for the ANALYZERPROJECT. To do this, you will need access to DataStage Administrator and DataStage Director. DataStage Director is not included with Information Analyzer and will require a separate license to use it. Open DataStage Administrator, login with a DataStage administrator user, click the Projects tab, select the ANALYZERPROJECT and click Properties.

Collecting runtime data (2 of 13)

- Go to General tab
- Click Environment



6

Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

Go to the General tab and click the Environment Button to open the list of environment variables for this project.

Collecting runtime data (3 of 13)

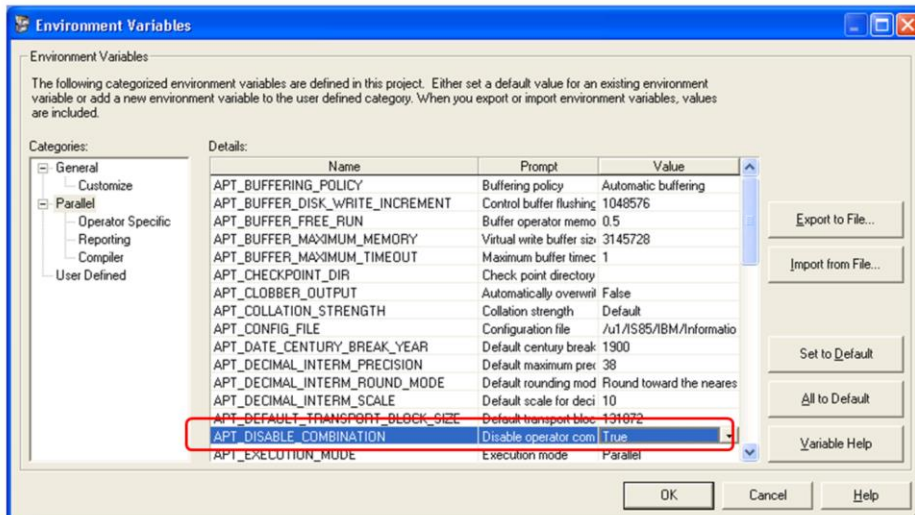
- Set debug environment variables to True
 - **APT_DISABLE_COMBINATION**
 - Located under Parallel category
 - Disables default optimization that CA does to improve performance by combining processes
 - **APT_PM_PLAYER_TIMING**
 - Located under Reporting category
 - Force each process to print processing time used
 - **APT_PM_SHOW_PIDS**
 - Located under Reporting category
 - Adds process ID or pid of each process to log
 - **APT_RECORD_COUNTS**
 - Located under Reporting category
 - Adds amount of records each operator is processing to log

There are four environment variables that should be set to True to collect runtime data relevant for performance issues. The first variable is **APT_DISABLE_COMBINATION** which is located under the Parallel Category. Setting this to true will disable the default optimization that CA does to improve performance by combining processes. By disabling this, you may experience some performance decrease but the logs will display more detailed information. Next, set **APT_PM_PLAYER_TIMING** which is located under the Reporting Category. This variable will force each process to print the processing time used to the log file. Next, set **APT_PM_SHOW_PIDS** which is located under the Reporting Category. This variable will add the process ID, or pid, of each process. This can be useful if there is a particular process that is taking significantly longer than others so you can trace this process at the OS level. The last environment variable to set to true is **APT_RECORD_COUNTS** which is also located under the Reporting Category. This variable will add the amount of records each operator is processing to the log file. This is useful to ensure that the data is evenly distributed on each operator.

The next slide displays how to modify these environment variables.

Collecting runtime data (4 of 13)

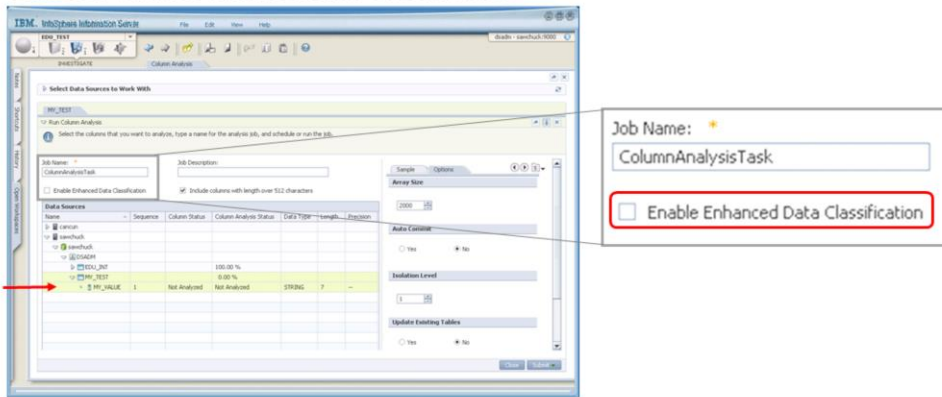
- Set each environment variable
- Save changes



To set these environment variables, find the variable you want to change then click the dropdown box to the right of the variable name and set the value to True. After having updated all the variables, click OK to save the changes. Click OK again to get back to the list of projects. It is not necessary to restart any component for the changes to take effect. All new CA submissions will use this configuration. Remember to set these variables back to False after you are done collecting runtime data.

Collecting runtime data (5 of 13)

- Select columns to analyze
- Uncheck "Enable Enhanced Data Classification"



9

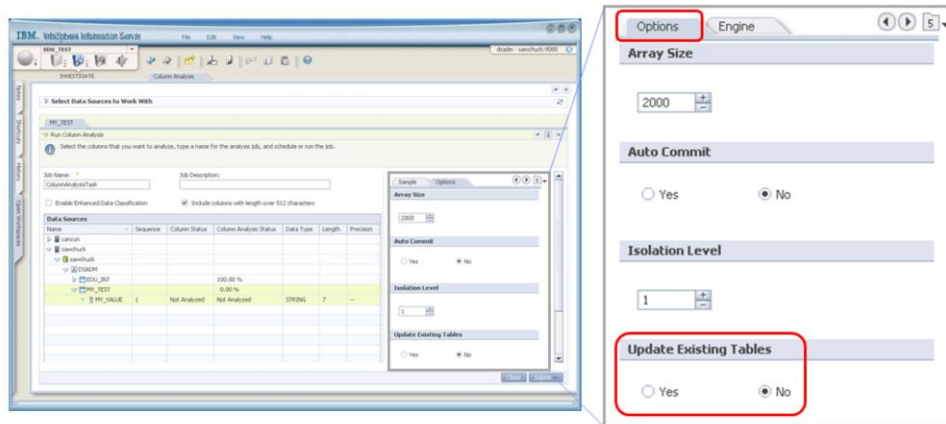
Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

Open the IA Client and select the column or columns you want to analyze. On the central panel, uncheck "Enable Enhanced Data Classification" if you see this option available. Older versions of IA did not have this option. Enabling this feature will add a process during the post-process step which will increase the time it takes for CA to complete. To isolate problems related only to the main process, disable this feature. If you have reasons to believe that this feature is causing performance problems then skip to the slide "Troubleshooting: XMETA when EDC on".

Collecting runtime data (6 of 13)

- Set "Update Existing Tables" to No



10

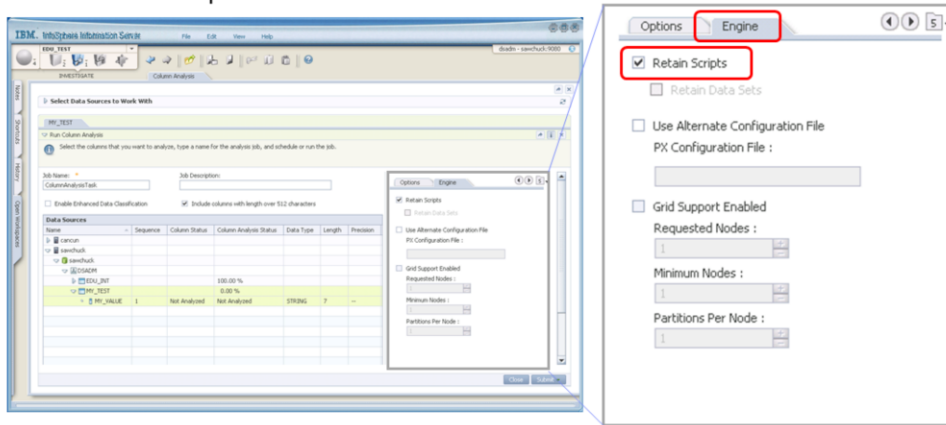
Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

On the right panel go to Options and set "Update Existing Tables" to "No". When troubleshooting performance issues, it is recommended not to update existing tables because database updates are slower than inserts. Setting this to No, forces CA to insert results into empty new tables which is a faster process. Set "Update Existing Tables" to "Yes" only if you want to keep previous results and you are not concerned about performance.

Collecting runtime data (7 of 13)

- Check “Retain Scripts”



11

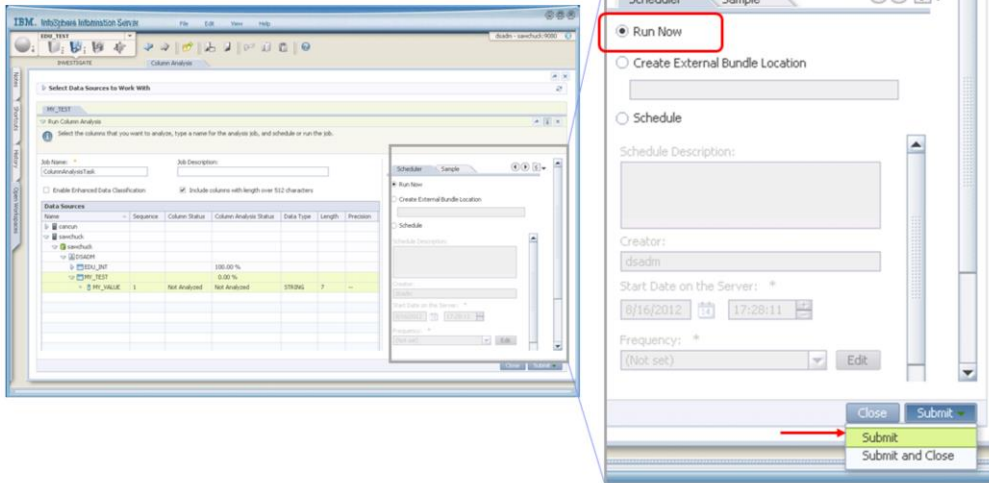
Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

Next, click the Engine tab and check the “Retain Scripts” box. This will keep additional files that can be used later to troubleshoot problems.

Collecting runtime data (8 of 13)

- Submit Column Analysis



12

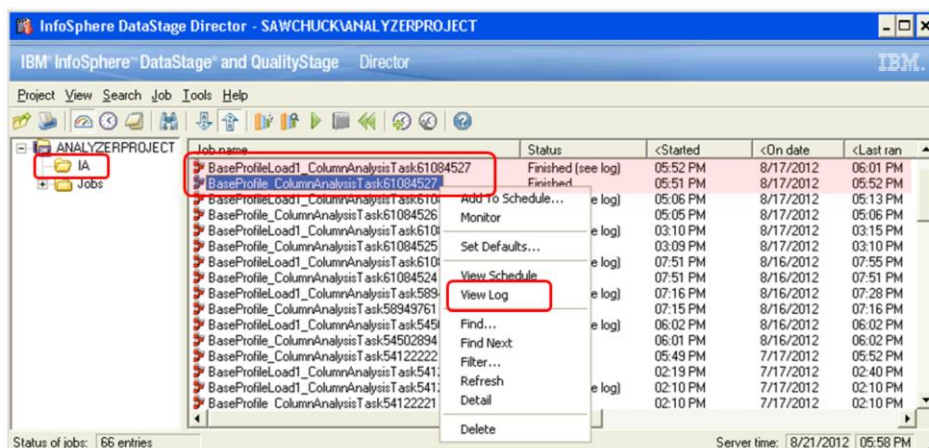
Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

Next, submit the Column Analysis.

Collecting runtime data (9 of 13)

- Open ANALYZERPROJECT project in Director
- Find Jobs created by Column Analysis
- Open logs



13

Troubleshooting performance problems in Information Analyzer Column Analysis version 8

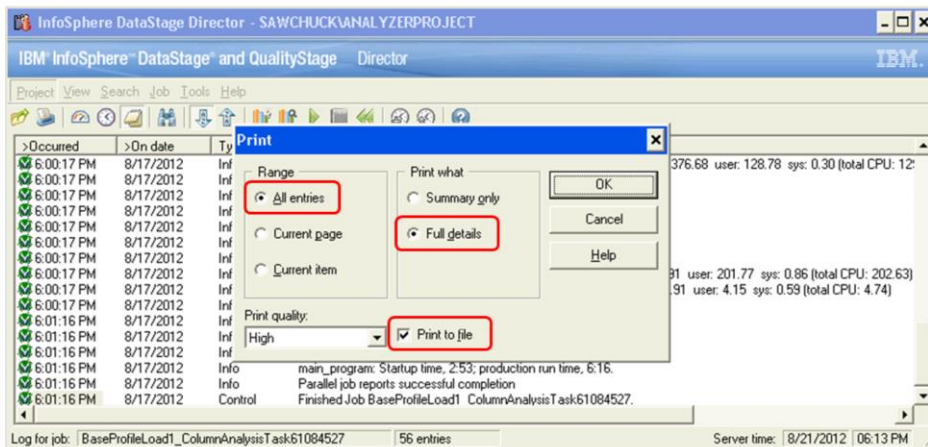
© 2012 IBM Corporation

Next, open the DataStage Director and connect to the ANALYZERPROJECT. Select the IA folder and sort the Jobs by the “On Date” Column. Use the timestamps to identify the most recent jobs. CA creates two types of jobs, one with the prefix “BaseProfile” to extract the data from source database and a second type with the prefix “BaseProfileLoad” to analyze the data and insert the results into the IADB database. The number of actual jobs that CA creates, depends on the number of columns submitted for analysis. Ten columns or fewer will only create two jobs, one BaseProfile and one BaseProfileLoad job. More than ten columns will create additional jobs.

Right click the “BaseProfile” log and click “View Log” to open the log.

Collecting runtime data (10 of 13)

- Save log to a file



14

Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

The next step is to save the log to a file. Click “Project” in the Menu bar, select the options “All entries”, “Full details”, and “Print to file”. Click OK and save the file. This will create a text file with the contents of this log. Repeat this process for the “BaseProfileLoad” log.

Collecting runtime data (11 of 13)

- Open log files in text editor
- Identify lines that include elapsed times

```
Item #: 118
Event ID: 117
Timestamp: 2012-08-17 17:06:21
Type: Info
User Name: dsadm
Message Id: IIS-DSEE-TFPM-00325
Message: fd_compute,3: Operator completed. status: APT_Statusok elapsed: 37.50 user: 1.44
sys: 0.09 (total CPU: 1.53)
```

Operator, Node

```
Item #: 119
Event ID: 118
Timestamp: 2012-08-17 17:06:21
Type: Info
User Name: dsadm
Message Id: IIS-DSEE-TFPM-00326
Message: fd_compute,3: Heap growth during runLocally(): 0 bytes
```

Total CPU time

```
Item #: 120
Event ID: 119
Timestamp: 2012-08-17 17:06:21
Type: Info
User Name: dsadm
Message Id: IIS-DSEE-TFPM-00325
Message: fd_compute,3: Operator completed. status: APT_Statusok elapsed: 37.49 user: 1.88
sys: 0.06 (total CPU: 1.94)
```

Next, identify the processing time associated with each operator. An operator represents a part of the job that is responsible for a specific task, for example, to extract data. Open the log files in a text editor and identify the lines that include the string “Operator completed”. These lines will include the operator name, the node and the “total CPU” time, expressed in seconds, for the operator. Write down all “total CPU” times shown in the log and add them together by operator. The order in which these values appear in the log is not relevant for this analysis.

Collecting runtime data (12 of 13)

- Run time versus Startup time

```
Item #: 123
Event ID: 122
Timestamp: 2012-08-17 17:06:23
Type: Info
User Name: dsadm
Message Id: IIS-DSEE-TCOS-00026
Message: main_program: Startup time, 0:07; production run time, 0:37.

Item #: 124
Event ID: 123
Timestamp: 2012-08-17 17:06:23
Type: Info
User Name: dsadm
Message Id: IIS-DSTAGE-RUN-I-0124
Message: Parallel job reports successful completion

Item #: 125
Event ID: 124
Timestamp: 2012-08-17 17:06:23
Type: Control
User Name: dsadm
Message Id: IIS-DSTAGE-RUN-I-0077
Message: Finished Job BaseProfile_ColumnAnalysisTask61084526.

End of report.
```

Towards the end of the log you will see the **Startup and Production Run Time** values. The **Startup time** refers to the time it takes for DataStage to start all the processes of the job. The larger the number of nodes used in the configuration file, the greater the Startup time will be. The **production run time** is the elapsed time from the moment the first operator starts to run until the last operator finishes running. Production run time includes idle times. If you add these two times together you will get the amount of time the job took to complete. The processing times that each operator reports and that was discussed in the previous slide occur during the production run time part of the job. However, keep in mind that processing times are reported individually by each operator and operators can run in parallel. Because of this, the total processing time of all operators is typically larger than the production run time shown by the log. However, the opposite is also possible because “Total CPU” times do not include idle times between processes and in busy systems the Production Run time can be larger.

As a result of this, Total CPU times by operators should only be used as a reference to understand how the processing time is being distributed by step and not as a way of calculating the actual time it took the job to complete. To calculate that, add the Startup time plus the production run time.

The next slide displays how to use the total processing times by operator to identify bottlenecks.

Collecting runtime data (13 of 13)

- Manually accumulate total values by operator and then by step

Job	Operator	Total CPU	Phase		
			Extract	Analyze	Insert
BaseProfile	pxbridge	18.6	x		
	fd_compute	27.4		x	
BaseProfileLoad	fd_extract	3.4		x	
	ca_properties	129.0		x	
	generator	4.7		x	
	pxbridge	202.6			x
			18.6	164.7	202.6

17

Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

For the purpose of this presentation, the production run times of the operators is classified in three phases. The first phase, **Extract**, includes operators that connect to the source database and extract the data to be analyzed. The next phase, **Analyze**, includes operators that look at the extracted data and analyze it to find data formats, data types, patterns, null ability, and so on. Finally, the **Insert** phase includes operators that take the results generated during the Analyze step and inserts them into the IADB repository.

Using the total processing times by operator, you can now calculate a metric to estimate how much processing time took each phase. The table displayed on this slide shows how to assign each operator's processing time to a specific phase. For example, the "pxbridge" operator in the "BaseProfile" job belongs to the Extract phase.

Use this table to compare different executions of CA and understand how processing times per Phase behave after you change certain variables such as the number of nodes, the source database, or the data. You should repeat the steps in this exercise a few times changing one variable at a time to better understand what are the common bottlenecks.

Troubleshooting

- Important performance factors by step

Factors	Extract	Analyze	Insert
Network Latency: Source DB – Engine	Yes		
Source DB Performance	Yes		
Network Latency: Between Engines (MPP)		Yes	
Data Attributes		Yes	
Network Latency: Engine – IADB			Yes
IADB Performance			Yes
Engine Performance	Yes	Yes	Yes

18

Troubleshooting performance problems in Information Analyzer Column Analysis version 8

© 2012 IBM Corporation

Once the bottlenecks are identified, the next step is to troubleshoot the possible causes. The table displayed on this slide shows the main performance factors for each phase. For example, if the “Insert” phase is taking too long then check three things: the “Network Latency” between the Engine and IADB, the IADB Performance, and the Engine Performance.

Be aware that a large value for a phase, in particular for the Analyze or Insert phases, does not necessarily mean that you are experiencing a performance problem. Analyzing very large tables with high cardinality can take a lot of time to complete. Before drawing any conclusions about performance, run and compare multiple executions to make sure that you are observing a consistent behavior.

The next slides display how to troubleshoot each factor.

Troubleshooting: Source database performance (1 of 3)

- Find job number
 - Open BaseProfile log
 - Identify library path variable

Example:

```
LIBPATH=/opt/IBM/InformationServer/Server/Projects/ANALYZERPROJECT/RT_BP427.O:/opt/IBM/InformationServer/Server/DSComponents/lib:/opt/IBM/InformationServer/Server/DSComponents/bin:/opt/IBM/InformationServer/Server/DSParallel:/opt/IBM/InformationServer/Server/PXEngine/user_lib:/opt/IBM/InformationServer/Server/PXEngine/lib:/opt/IBM/InformationServer/Server/Projects/ANALYZERPROJECT/buildop:/usr/lib/lib_64:/opt/IBM/InformationServer/Server/branded_odbc/lib:/opt/IBM/InformationServer/Server/DSEngine/lib:/opt/IBM/InformationServer/Server/DSEngine/uvdlls:/opt/IBM/InformationServer/ASBNode/apps/jre/lib/ppc64:/opt/IBM/InformationServer/ASBNode/apps/jre/lib/ppc64/j9vm
```

–Find RT_BP folder

```
LIBPATH=/u1/IS85/IBM/InformationServer/Server/Projects/ANALYZERPROJECT/RT_BP427.O
```

– Write down job number

To confirm if there is a performance problem on the source database you need to test the SQL query that IA is running to extract data. To find this SQL query, you first need to identify the job number. The easiest way to do this is to open the BaseProfile log in a text editor and identify the line that includes the library path variable. Environment variable values are listed in the second entry of each log. This variable can be named LIBPATH, LD_LIBRARY_PATH or SHLIB_PATH depending on your operating system. Inspect the value of this variable and identify the path that includes the RT_BP folder. This typically appears at the beginning of the entry. Write down the number after the prefix RT_BP. This is the job number. In this example, the number is 427.

Troubleshooting: Source database performance (2 of 3)

- Find SQL statement

- Open terminal session

- cd to ANALYZERPROJECT/RT_SCxxx

```
cd /opt/IBM/InformationServer/Server/Projects/ANALYZERPROJECT/RT_SC427
```

- Open OshSript.osh in text editor

- Find <SQL> section

```
# OSH / orchestrate script for Job BaseProfile_ColumnAnalysisTask581
  compiled at 19:15:54 16 AUG 2012
pxbridge
-XMLProperties '<?xml version='1.0' encoding='UTF-16'?>
<Properties version='1.0'><Common><Context
  type='int'>1</Context><Variant
  type='string'>9.1</Variant><MaximumRecords type='int'>-
1</MaximumRec
ords></Common>
<Usage><TableName type='string'><![CDATA[DSADM.MY_TEST]]></TableName>
<GenerateSQL type='bool'><![CDATA[0]]></GenerateSQL>
<EnableQuotedIDs type='bool'><![CDATA[1]]></EnableQuotedIDs>
<SQL><SelectStatement type='string'>
<![CDATA[select "MY_VALUE" from DSADM.MY_TEST]]>
</SelectStatement>
</SQL>
```

Next, open a terminal session and change directories to the RT_SC folder that ends with the job number you obtained in the previous slide. In this example, it is RT_SC427. Open the OshScript.osh file in a text editor and find the SQL tag. This will contain the query that CA is using to extract data.

Troubleshooting: Source database performance (3 of 3)

- Test SQL statement outside of IA with database client
 - If time is similar to Extract Phase
 - Performance issue in source database
 - Contact DBA
 - If time is better than Extract Phase
 - Run SQL statement using ODBC example program
 - Refer to [Tech Note 1434177](http://www-01.ibm.com/support/docview.wss?uid=swg21434177) for instructions
 - If example program reproduces issue
 - Review “Troubleshooting: Network connectivity” section

Run the SQL statement outside of IA using a database client application and compare the times it takes to complete the query. If you obtain a time similar to what you got in the Extract phase, then the performance problem is coming from the source database. Contact your DBA to discuss the performance of this source database. If you obtain a significantly better time using a database client application, then run the SQL using the “example” program from the Engine Server. Refer to Tech Note 1434177 for instructions on how to do this. This will use the same type of connectivity that IA is using to extract data. If the SQL statement executed with the example program reproduces the slow performance you get in IA, review the section “Troubleshooting: Network connectivity”.

Troubleshooting: IADB performance

- Check
 - Array size
 - IA patches
 - IADB creation
 - ISA Lite Health Check
 - Database errors
 - See “Troubleshooting: Network latency” section

Next, this presentation discusses troubleshooting IADB performance. IADB performance can be impacted by several factors. First, try to increase the array size located in the Options tab of CA. Higher values for this parameter can improve the performance of inserting rows to IADB. Try doubling the value and compare results. If you observe an improvement, you can continue doubling value until you find an optimal setting. If the array size is set too high, performance will begin to decrease.

Next, make sure that you are running the latest patches for the version of IA that you are using. New patches may introduce improvements in performance.

Next, confirm that the IADB database was created using the scripts provided with the installation media. If IADB was created without following the documented procedure, then IADB needs to be re-created.

Next, check that there are no errors for IA shown in the ISA Lite Health Check. For more details about ISA Lite, see [Tech Note 4022700](#).

Next, check database logs to confirm that there are no errors, such as lack of space. Check with your DBA to obtain this information.

Finally, if the IADB is on a different machine than the engine, check the section “Troubleshooting: Network latency” to measure the latency between the engine and the IADB.

Troubleshooting: XMETA performance (when using EDC)

- EDC adds extended data class to XMETA
- CA with EDC takes longer time to complete
- XMETA can be a bottleneck when EDC is enabled and database statistics are not updated

Another performance issue may be caused by EDC. EDC is a feature that identifies the extended data class of each distinct value of a column. When EDC is enabled, CA will take longer time to complete; this is the expected behavior. However, some customers have reported significant performance problems after adding EDC. These problems were caused by XMETA not having updated database statistics. If you are having performance problems that occur only when EDC is enabled, contact your DBA to make sure that RUNSTATS has been executed recently to update the statistics of the XMETA database.

Troubleshooting: Network latency (1 of 2)

- Use ping command
- Example

```
$ ping myServer.newco.com
PING myServer.newco.com : (9.3.2.1): 56 data bytes
64 bytes from 9.3.2.1: icmp_seq=0 ttl=255 time=0 ms
64 bytes from 9.3.2.1: icmp_seq=1 ttl=255 time=0 ms
64 bytes from 9.3.2.1: icmp_seq=2 ttl=255 time=0 ms
64 bytes from 9.3.2.1: icmp_seq=3 ttl=255 time=0 ms

---- myServer.newco.com PING Statistics----
 9 packets transmitted, 9 packets received, 0% packet loss
round-trip min/avg/max = 0/0/0 ms
```

- Check [Tech Note 1515972](http://www-01.ibm.com/support/docview.wss?uid=swg21515972):

<http://www-01.ibm.com/support/docview.wss?uid=swg21515972>

Network connectivity can also impact the performance of CA. To investigate if this is a problem, measure the speed of a data packet round-trip using the ping command. This slide displays an example using the ping. For a more accurate way to measure network latency, see Tech Note 1515972

Troubleshooting: Network latency (2 of 2)

- For problems during Extract phase
 - Measure latency between the engine and source database
- For problems during Insert phase
 - Measure latency between engine and IADB
- For problems during Analyze phase with MPP configuration
 - Measure latency between Engines
- Packet round-trip average should be close to 0ms for optimal performance

If you are experiencing problems during the Extract phase, measure latency between the engine and the Source Database.

If you are experiencing problems during the Insert phase, measure latency between the engine and the IADB.

If you are running an MPP configuration and you are experiencing problems during the Analyze phase, measure latency between engines.

Ideally, a packet round-trip average should be close to 0ms for optimal performance. Larger values will impact the communication between the servers. If you measure high values for this metric, contact your Network administrator to discuss ways to improve this.

Troubleshooting: Data attributes

- Factors that may impact CA's time
 - High Cardinality (% of distinct values in column)
 - Number of rows
 - Column length
 - Date types and values
 - Virtual columns
- Consider
 - Checking patch level
 - Using data sampling
 - Increasing number of nodes (see "Troubleshooting: Engine performance" section)

The characteristics of the data you analyze also have a significant impact in the overall performance of CA. This is the expected behavior of the product. There are several data attributes that can impact CA's performance.

High cardinality, which is defined as the percent of distinct values in a column, may have an impact. The higher the cardinality of a column, the longer it will take to be analyzed. Number of rows may also have an impact on performance. The larger the number of rows of a table, the longer the data will take to be analyzed.

Column length is another contributing factor. The more bytes a column has the longer it will take to be analyzed. This impact is only noticeable when comparing columns of very different lengths.

Data types may also be a factor as certain values of data types, such as dates or strings that can be interpreted as dates, will take longer time to be analyzed.

Another factor that may cause performance problems are virtual columns. These are columns that do not exist in the source database and are defined in IA. These columns will require extract calculations and will increase the time to complete CA.

If you identify that one of these factors is associated with significantly slower times, make sure that you are using the latest patches available. Try to use data sampling to reduce the number of rows you are analyzing and increase the number nodes as explained in the next slide, Troubleshooting: Engine performance.

Troubleshooting: Engine performance

- Increase number of nodes of configuration file
 - Refer to

http://publib.boulder.ibm.com/infocenter/isinfsv/v8r7/topic/com.ibm.swg.im.is.ds.parjob.dev.doc/topics/c_deeref_The_Parallel_Engine_Configuration_File.html

- One node per two CPUs in SMP or one node per CPU in MPP
- Try with one, two, four and eight nodes and compare results
- As number of nodes increase, startup time increases

You can improve the performance of the engine by increasing the number of nodes in the configuration file used by IA. For instructions on how to do this, see the online documentation at the link displayed on this slide. The rule of thumb is to have one node per each two CPUs in an SMP configuration. In an MPP configuration, you can use one node per CPU. This is a suggested starting point. Modify the number of nodes and compare values to find what works best in your environment.

Run tests using one, two, four and eight nodes and compare the results.

Keep in mind that the more nodes you have, the greater the startup time will be. This time, it is not included in the processing times reported by operators, which refer to run time, so when tuning the number of nodes, take into account this time as well.

Trademarks, disclaimer, and copyright information

IBM, the IBM logo, ibm.com, DataStage, DB2, DB, IA, and InfoSphere are trademarks or registered trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of other IBM trademarks is available on the web at "[Copyright and trademark information](http://www.ibm.com/legal/copytrade.shtml)" at <http://www.ibm.com/legal/copytrade.shtml>

Other company, product, or service names may be trademarks or service marks of others.

THE INFORMATION CONTAINED IN THIS PRESENTATION IS PROVIDED FOR INFORMATIONAL PURPOSES ONLY. WHILE EFFORTS WERE MADE TO VERIFY THE COMPLETENESS AND ACCURACY OF THE INFORMATION CONTAINED IN THIS PRESENTATION, IT IS PROVIDED "AS IS" WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED. IN ADDITION, THIS INFORMATION IS BASED ON IBM'S CURRENT PRODUCT PLANS AND STRATEGY, WHICH ARE SUBJECT TO CHANGE BY IBM WITHOUT NOTICE. IBM SHALL NOT BE RESPONSIBLE FOR ANY DAMAGES ARISING OUT OF THE USE OF, OR OTHERWISE RELATED TO, THIS PRESENTATION OR ANY OTHER DOCUMENTATION. NOTHING CONTAINED IN THIS PRESENTATION IS INTENDED TO, NOR SHALL HAVE THE EFFECT OF, CREATING ANY WARRANTIES OR REPRESENTATIONS FROM IBM (OR ITS SUPPLIERS OR LICENSORS), OR ALTERING THE TERMS AND CONDITIONS OF ANY AGREEMENT OR LICENSE GOVERNING THE USE OF IBM PRODUCTS OR SOFTWARE.

© Copyright International Business Machines Corporation 2012. All rights reserved.