

Workload XM 2.2.1

Workload XM Cluster Optimization

Date published: 2020-12-04

Date modified: 2021-09-21

CLOUDERA

<https://docs.cloudera.com/>

Legal Notice

© Cloudera Inc. 2024. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

Working with Workload XM.....	5
Specifying a time range.....	5
Analyzing Your Workload Cluster Costs with Workload XM Cost Centers.....	5
Creating a Workload XM Cost Center.....	5
Displaying Your Job Costs Associated with a Cost Center Cluster.....	6
Assigning Uncategorized Resources to a Cost Center.....	7
Triggering an Action across Jobs and Queries.....	8
Creating an Auto Action Event.....	8
Managing your Auto Actions.....	9
Trigger Email Notification Example.....	10
Troubleshooting an Abnormal Job Duration.....	11
Troubleshooting Failed Jobs.....	15
Determining the Cause of Slow and Failed Queries.....	18
Classifying Workloads for Analysis with Workload Views.....	20
Working with Auto Generated Workload Views.....	20
Defining Workload Views Manually.....	22
Troubleshooting with the Job Comparison Feature.....	27
Identifying File Size Storage Issues.....	32
Displaying File Size Metadata.....	33
Displaying the Metadata of a Table.....	34
Assigning Access Roles in Workload XM.....	35
Understanding the Workload XM Access Roles.....	35
Understanding a Workload XM Cluster Policy.....	37
Configuring a Default Systems Administrator for Workload XM.....	38
Assigning Workload XM Access Roles.....	39
Assigning a Workload XM System Admin Access Role.....	39
Assigning a Workload XM Cluster Admin Access Role.....	39

Assigning a Workload XM Cluster User Access Role.....	40
Managing Your Workload XM Access Roles.....	41

Purging HDFS Data..... 41

Understanding the Purge Date used by the Purge Event.....	42
Workload XM Purge Event Parameters.....	42
Configuring the Workload XM Purge Event.....	44
Manually Executing a Workload XM Purge Event.....	44
Managing your Workload XM Purge Event.....	45

Working with Workload XM

Tasks for identifying and troubleshooting job and query abnormalities and failures, optimizing workloads, and improving job performance with Workload XM.

Specifying a time range

Choose a time period in which your workload results are displayed in Workload XM for analysis and troubleshooting.

About this task

Describes how to specify a time period for displaying current or historical data about your cluster and the jobs performed in that cluster.

By default, Workload XM displays workload data for the last 24 hours. If there is no data available during that time, Workload XM displays the nearest date range that is available.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. In the Clusters page, select the cluster required for analysis.
3. From the time-range list in the Cluster Summary page, do one of the following:
 - For a predefined period, select one of the default periods of time that meets your requirements.
 - For an exact date and time range, select Customize and then either, enter the date and time range using the YYYY/MM/DD HH:MM:SS format for the beginning and the ending time period, or in the calendar element, select the beginning and ending time period.
4. Click Ok, which clears any existing workload data from the chart and table components for the existing period of time.

Results

All charts and tables in Workload XM are updated to reflect the workload data for the chosen time period.

Analyzing Your Workload Cluster Costs with Workload XM Cost Centers

Define customized cost centers based on user or pool resource criteria and CPU and memory consumption with the Chargeback feature. Once defined Workload XM visually displays a Workload cluster's current and historical costs. With these cost insights you can then plan and forecast budgets and future workload environments and/or justify current user groups and resources.

Creating a Workload XM Cost Center

Create Workload XM cost centers that enable you to display your current and historical workload cluster and resource costs that can be used for planning, budgeting, and forecasting future workload environments.

About this task

Describes how to configure your Workload XM Chargeback settings, which define your cost centers and the unit costs of your resources, and create a Workload XM cost center.



Note: To avoid cost duplication, resources must only be assigned one cost center.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. From the Workload XM Navigation side-bar, select Chargeback.
3. Globally define your cost center criteria and memory and CPU costs by clicking Chargeback Setup.
4. From the Setup page, do the following:

- a. From the Select Chargeback criteria section, select your cluster's chargeback criteria.



Note: Cost centers are associated with a specific criteria. If you later change the Chargeback criteria setting the cost centers associated with the previous selection are hidden. You can revert back to these cost centers by reselecting their Chargeback criteria.

To revert back to previous cost centers, from the Actions list on the Chargeback page, select Chargeback Settings and then reselect their criteria option.

- b. From the Cluster list of the Cluster Selection section, select the clusters required for your cost centers. Where, the cost calculations use resource utilization for each of your chosen clusters.
 - c. In the CPU field of the Unit cost section, enter the amount, in dollars, for each CPU core hour.
 - d. In the Memory field of the Unit cost section, enter the amount, in dollars, for each Gigabyte hour.
 - e. Click Complete Setup.
5. From the Chargeback page, create a new cost center by clicking Create a Cost Center.
 - a. In the Name field, enter a unique name for your cost center.
 - b. (Optional) In the Description field, enter a meaningful description for the cost center.
 - c. Depending on the Chargeback criteria value you selected when you configured your Chargeback settings, do one of the following:
 - If you selected Pool, in the Add Pools field, enter one or multiple resource pools.
 - If you selected User, in the Add Users field, enter one or multiple users.
 - d. Click Create.

Results

Once you have configured your Chargeback settings and created a cost center you can view your job costs associated with a cost center cluster.

Displaying Your Job Costs Associated with a Cost Center Cluster

Steps for displaying your Workload cluster jobs associated with a cost center cluster.

About this task

Describes how to view your workload costs associated with a cluster.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. From the Workload XM Navigation side-bar, select Chargeback.
3. In the Chargeback page, select a cost center.

Your cost center page opens displaying the costs, and the CPU and memory usage associated with the cost center.
4. To view more details about the pool, user, and job costs for a specific cluster in the cost center, from the Cluster column, locate the cluster and then either click its name or click the greater-than arrow (>) at the end of its row.

Assigning Uncategorized Resources to a Cost Center

Steps for moving unassigned resources into an existing or a new Workload XM cost center.

About this task

Describes how to locate and move uncategorized resources into an existing or a new Workload XM cost center.



Note: To avoid cost duplication, resources must only be assigned one cost center.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. From the Workload XM Navigation side-bar, select Chargeback.
3. In the Chargeback page, select a cost center and then a cluster.
4. From the Overview tab, scroll down and click inside the Uncategorized section.

The Uncategorized page opens.
5. Select the required uncategorized resource tab.
6. From either the Pools, Users, or Clusters page, select the check boxes of the resources you require for your cost center.

The Assign Cost Center button becomes visible.
7. Click Assign Cost Center.
8. From the Select Cost Center list, do one of the following:
 - a. To add the uncategorized resource/s in a new cost center, select New Cost Center and then click Create a new cost center.
 - b. To add the uncategorized resource/s in an existing cost center, select an existing cost center and then click Assign to Cost Center.
9. (Optional) Repeat steps 4-8 until all your uncategorized resources are placed in your Workload XM cost centers.

Triggering an Action across Jobs and Queries

You can trigger actions, in real-time, across jobs and queries with Workload XM auto action events that are defined by you. When a job or query matches your action's criteria and its conditions exist, the auto action event is triggered. For example, memory exhaustion can cause nodes to crash or jobs to fail. Knowing when available memory is falling below a specific threshold enables you to take steps before a potential problem occurs. With the Auto Actions feature, you can create an action that informs you through an email when a job is consuming too much memory so that you can take steps to alleviate a potential problem.

Creating an Auto Action Event

The steps to create a Workload XM auto action event, which is triggered when a job or query meets the action's criteria and conditions. For example, when a job uses too much memory it may cause other jobs to fail or increase a job's runtime. You can create an action that informs you when a job is consuming too much memory.

About this task

Describes how to create a Workload XM Auto Action.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. From the Workload XM Navigation side-bar, select Auto Actions.
3. Do one of the following:
 - If no other auto actions exist, click Auto Actions Setup.
 - If other auto actions exist, click Create an Auto Action.

The Auto Actions Create page opens.

4. In the Auto Action Name field, enter a unique name that is easily identifiable.
5. From the Engine list, do nothing.



Note: At this time only Spark Engines are supported for the Auto Action feature.

6. (Optional) Define the criteria for the auto action by doing the following:
 - a. From the Criteria list, choose between Pool and User.
 - b. From the Operator list, choose between ANY OF and NONE OF.
 - c. In the Value field, enter the value for this filter.



Tip: You can define multiple AND filters for the Criteria by clicking the plus sign.



Note: An Auto Action does not require the Criteria filter:

- When included, only those jobs on the selected engine that meet the criteria conditions are tested by the Trigger.
- When not included, all jobs on the selected engine are tested by the Trigger.

7. Define the trigger for the auto action by doing the following:

- a. From the Metric list, choose between Memory Allocated (MB) and Application Name.
- b. From the Operator list, select an operator.
- c. In the Value field, enter the value for this trigger condition.



Tip: You can define multiple OR conditions for the trigger by clicking the plus sign.

- d. From the Action options, do nothing.



Note: At this time only email notifications are supported for the Auto Action feature.

- e. In the Emails field, enter the email address that you use to log into Workload XM.
- f. In the Subject field, enter the subject for the email that distinguishes the subject matter from other auto action emails.
- g. Click Create, which creates the action and adds it on the Auto Actions home page.

The Auto Actions page displays the action's configuration and status details in the following entry fields:

- Status, contains the current state of the action, as either Enabled or Disabled.
- Name, contains the unique name you entered for the auto action.
- Action, contains either the command that is to be executed or the function that is to be performed when the action is triggered.
- Engine, contains the selected engine on which the action is to be performed.
- Triggers, contains the action's Trigger conditions.
- Criteria, contains the action's Criteria filters.

Results

When a job or query meets the auto action's criteria and conditions the action event is triggered.

Managing your Auto Actions

Steps for updating, deleting, and disabling an auto action, and viewing your actions in Cloudera Manager.

The following Auto Actions management tasks are performed in the Auto Actions page, which is accessed by selecting Auto Actions in the Workload XM Navigation side-bar.

Updating your Auto Action

In the Auto Action page, click the action's vertical ellipsis, and select Edit. Make your changes and then click Update.

Deleting an Auto Action

In the Auto Action page, click the action's vertical ellipsis, and select Delete. In the confirmation message, click OK to confirm. The action is permanently removed.



Note: Unless the action is no longer required, Cloudera recommends disabling the action, as you may require the action at another time.

Disabling an Auto Action

In the Auto Action page, click the action's vertical ellipsis, and select Disable. In the confirmation message, click OK to confirm. The action is no longer active and the Disabled state is displayed in the action's Status column on the Auto Actions page.

Viewing Workload XM Auto Actions in Cloudera Manager

You can display more details about your Workload XM actions in Cloudera Manager.

To view your Workload XM Auto Actions in Cloudera Manager:

- In a supported web browser log in to Cloudera Manager.
- From the Cloudera Manager Navigation side-bar, click Diagnostics and then select Events.
- In the Events page, search for Content that contains AUTOACTIONTRIGGER.
- Expand your action to display more information about the event.

Trigger Email Notification Example

An example of a Workload XM Auto Actions Notification email that was triggered when the job matched the action's criteria and condition.

The following email notification example was sent when the listed application met the action's criteria and the trigger conditions, which are also included in the email notification.

Auto Action Notification

Hello,

We found an application that violates the thresholds you set.

Application Details:

Id	: application_1619627211029_0045
Name	: AwesomeName2
User	: hdfs
Type	: SPARK
Pool	: default
Allocated MB	: 16,384

Auto Action Triggered:

Name	: FirstSatisfyAction
Criteria	: USER ANY hdfs AND USER ANY admin AND USER ANY hive
Trigger	: allocatedMb GREATER_THAN 1

Please review the application and take necessary action.

[View Job](#)

Troubleshooting an Abnormal Job Duration

Identify areas of risk from jobs running on your cluster that complete within an unusual time period.

About this task

Describes how to locate and troubleshoot an abnormal job duration.

Steps with examples are included that explain how to further investigate and troubleshoot the cause of an abnormal job duration.

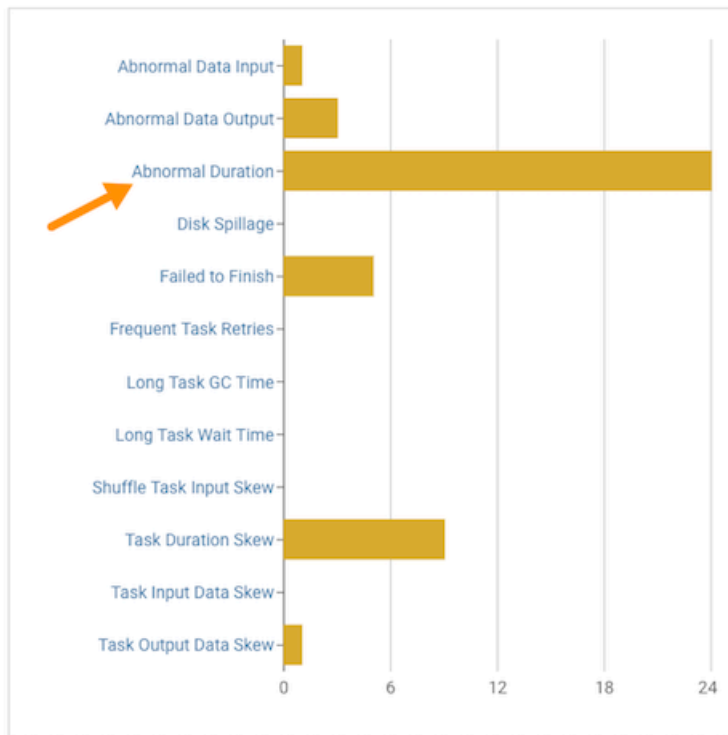
Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. In the Clusters page do one of the following:
 - In the Search field, enter the name of the cluster whose workloads you want to analyze.
 - From the Cluster Name column, locate and click on the name of the cluster whose workloads you want to analyze.
3. From the time-range list in the Cluster Summary page, select a time period that meets your requirements.
4. From the Engine column in the Usage Analysis chart, click an engine whose jobs you wish to analyze.

5. Display the number of jobs with an abnormal duration that executed within the selected time period by clicking the Abnormal Duration health check bar in the Suboptimal Jobs chart widget.

The Job page opens, listing all the jobs that have triggered the Abnormal Duration Health check.

Suboptimal Jobs



Tip: Any jobs that fall outside of their baseline are counted. You can hover over each bar within the chart to view how many jobs triggered each health check.

- To specify a specific amount of time in which the job either ran less than or more than the Health check rule, from the Duration list, either select a predefined time duration or select Customize and enter the minimum or maximum time period.

Jobs

Choose a duration range from the **Duration** list, or choose **Customize** to enter a custom maximum-minimum duration range.

Type	Job	Status	Start Time	Duration	Health Issue	Execution ID
HV	ins_from_...	Succeeded	04/19/2018 1:08 AM PDT		Abnormal Data Output Abnormal Data Input Abnormal Duration	
HV	insert over...	Succeeded	04/18/2018 10:48 PM P...	12m 10s	Task Output Data Skew Abnormal Data Output Task Input Data Skew Abnormal Data Input Abnormal Duration Data Processing Speed Skew	
HV	insert over...	Succeeded	04/18/2018 10:27 PM P...	16m 34s	Task Output Data Skew Abnormal Data Output Task Input Data Skew Abnormal Data Input	

- To view more details about a job, from the Job column, select a job's name and then click the Health Checks tab. The Baseline Health checks are displayed.
- To display more information about the job's duration, from the Baseline column, select Duration. For example, as shown in the following image the job finished much slower than the baseline:

Jobs

Log Analysis

Click to further investigate this job.

Baseline

Hide Normal Stages

Duration

Input Size

Output Size

Skew

Task Duration

Task Input Data

Task Output Data

Shuffle Input

Data Processing Speed

Resources

Task Wait Time

Log Analysis

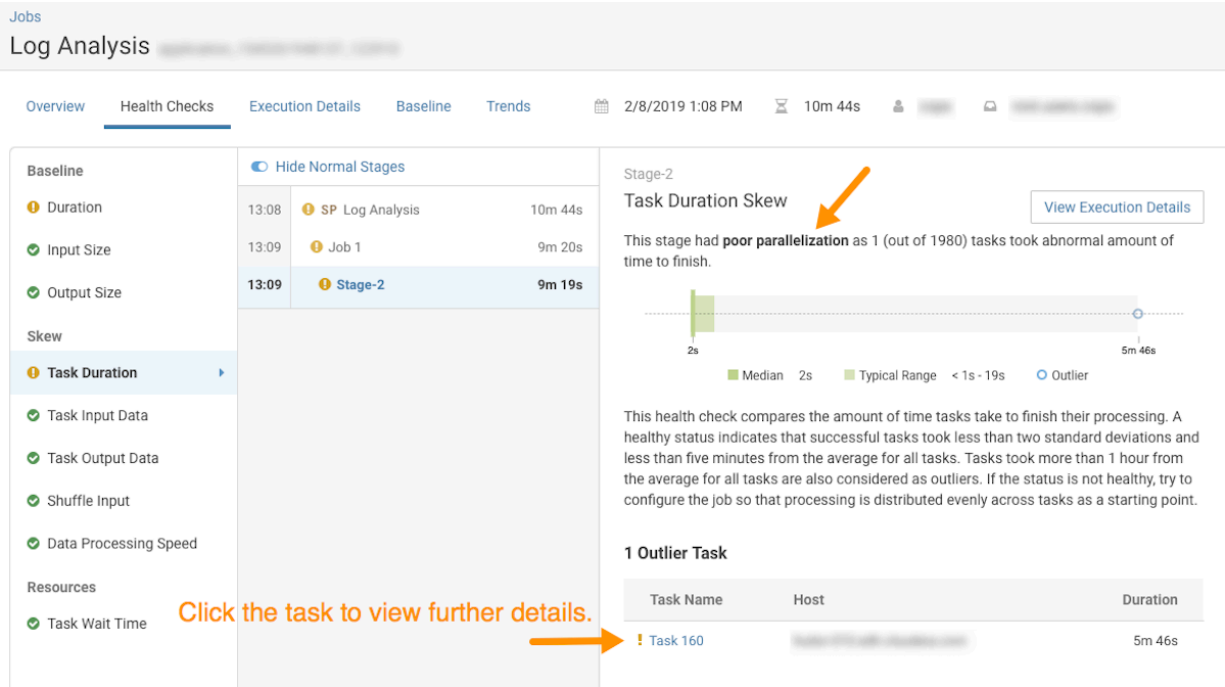
Abnormal Duration

Finished in 10m 44s, slower than the median duration 2m 8s. View all metrics.

Start Time	Execution ID	Duration
2/8/2019 1:08 PM		11m
2/8/2019 1:08 PM		8m
2/8/2019 1:08 PM		6m

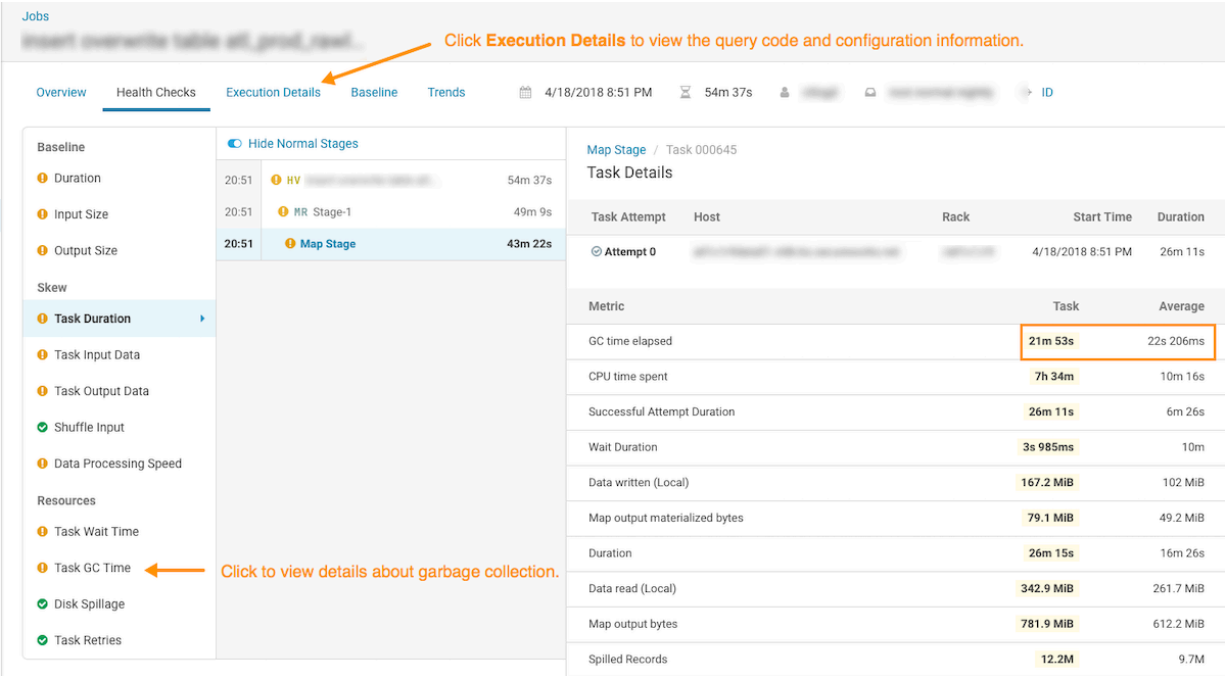
9. To display more information about the length of time the processing tasks took within a job, from the Baseline column, select Task Duration.

For example, as shown in the following image, a particular task took longer to complete than expected:



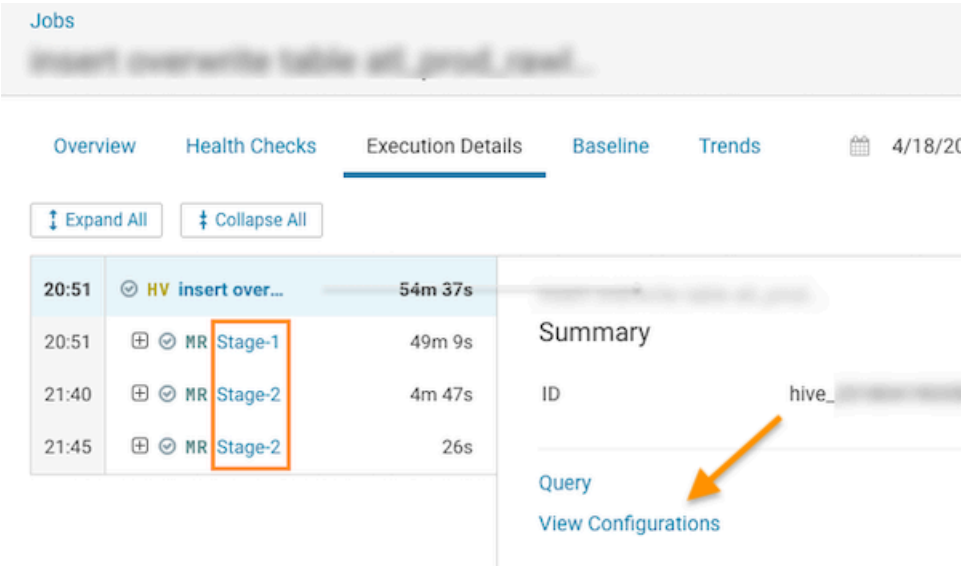
10. To display more information about the abnormal task, click the task, which opens the Task Details panel.

In the following example, the Task Details show that the abnormal task took significantly more time to complete the garbage collection process than the average:

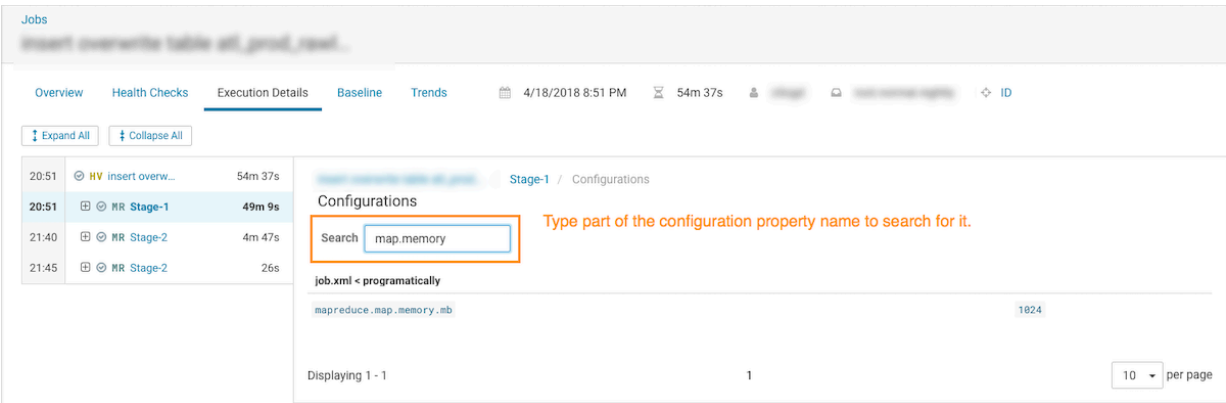


11. To display more information about the garbage collection for this abnormal task example, from the Baseline column, select Task GC Time.

12. In the Task GC Time page, click the Execution Details tab and then click one of the MapReduce stages:



13. In the Summary panel, click View Configurations and then locate the configuration for the garbage collection by entering the MapReduce memory configuration property name in the Search field:



The configuration setting for the garbage collection is 1024. This value could be causing the mapper JVM to have insufficient memory and triggering too many garbage collection processes. Increasing the value will improve cluster performance and remove this task as a potential risk.

Troubleshooting Failed Jobs

Steps for troubleshooting incomplete jobs running on your cluster.

About this task

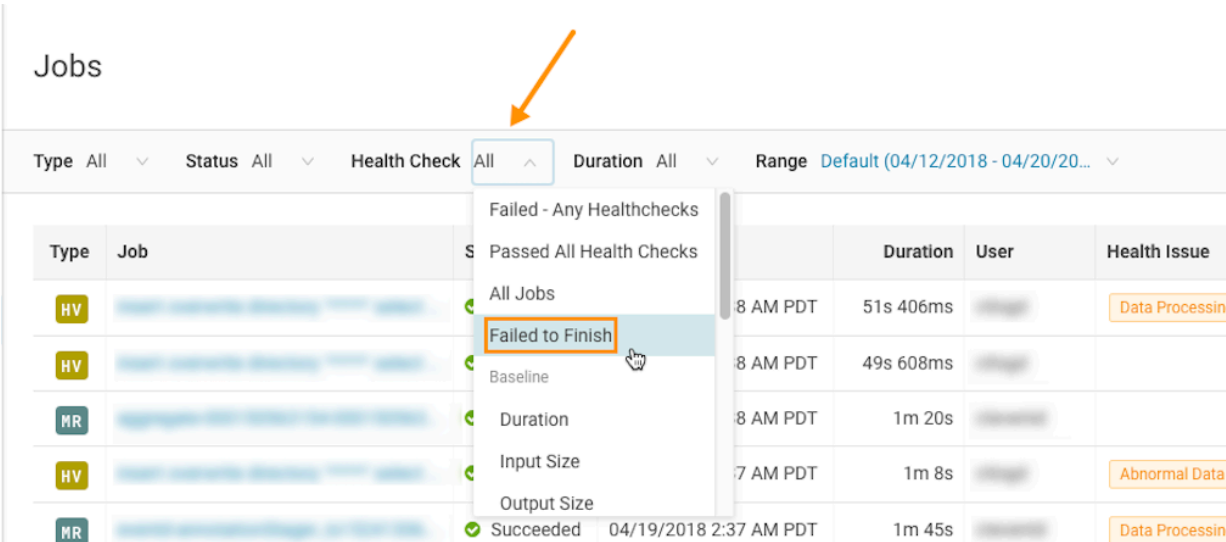
Describes how to locate and troubleshoot jobs that have failed to complete.

Steps with examples are included that describe how to further investigate and troubleshoot the root cause of an uncompleted job.

Procedure

- 1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
- 2. In the Clusters page do one of the following:
 - In the Search field, enter the name of the cluster whose workloads you want to analyze.
 - From the Cluster Name column, locate and click on the name of the cluster whose workloads you want to analyze.
- 3. From the time-range list in the Cluster Summary page, select a time period that meets your requirements.
- 4. From the Trend widget, select the tab of an engine whose jobs you want to analyze and then click its Total Jobs value.

The engine's Jobs page opens.
- 5. From the Health Check list, select Failed to Finish, which filters the list to display a list of jobs that did not complete.



- 6. To view more details about why a job failed to complete, from the Job column, select a job's name and then click the Health Checks tab.

The Baseline Health checks are displayed.

7. From the Health Checks panel, select the Failed to Finish health check.

For example, as shown in the following image, the failure occurred in the Map Stage of the job process:

Jobs

insert overwrite table concurrent_error...

Overview Health Checks Execution Details Baseline Trends

4/16/2018 3:39 AM 3m 4s

Health Check	Time	Status	Duration
Failed to Finish	03:39	Failed	3m 4s
Baseline	03:39	Failed	3m 4s
Duration	03:39	Failed	2m 52s
Input Size			
Output Size			

Click Map Stage and then click Execution Details.

Map Stage
Failed to Finish
Operation failed to finish.
This health check determines whether a job succeeded or failed.
Check Execution Details for more information.

View Execution Details

8. To display more information about the Map Stage process, click Map Stage and then from the Map Stage panel, click Execution Details.
9. To display all the failed tasks, in the Summary panel, click the Failed value:

Jobs

insert overwrite table concurrent_error...

Overview Health Checks Execution Details Baseline Trends

4/16/2018 3:39 AM 3m 4s

Expand All Collapse All

Time	Status	Duration
03:39	Failed	3m 4s
03:39	Failed	3m 4s
03:39	Failed	2m 52s

Stage-1 / Map Stage

Summary

Map Tasks

Completed 11 / 12

Failed 1

Average Map Time 1m 30s

Click the number of Failed tasks.

10. To display the reason for a task's failure, select and expand its error message.

For example, as shown in the following image, the task was not completed because it was stopped. To understand what triggered the Task KILL is received. Killing attempt! error message and to further troubleshoot the root cause, open the associated log file by clicking Logs.

The screenshot shows the 'Execution Details' tab in the Workload XM web UI. On the left, a list of tasks is shown with their status, name, and duration. The 'Map Stage' task is highlighted in red, indicating a failure. On the right, the detailed view of the 'Map Stage' is shown, including a table of attempts. The first attempt, 'Attempt 0', is shown with the error message 'Task KILL is received. Killing attempt!'. An orange arrow points to the 'Logs' link, with a text label 'Click to view log file.' above it.

Determining the Cause of Slow and Failed Queries

Identifying the cause of slow query run times and queries that fail to complete.

About this task

Describes how to determine the cause of slow and failed queries.

Steps with examples are included that explain how to further investigate and troubleshoot the cause of a slow and failed query.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. In the Clusters page do one of the following:
 - In the Search field, enter the name of the cluster whose workloads you want to analyze.
 - From the Cluster Name column, locate and click on the name of the cluster whose workloads you want to analyze.
3. From the time-range list in the Cluster Summary page, select a time period that meets your requirements.
4. From the Trend widget, select the tab of an engine whose jobs you wish to analyze and then click its Total Jobs value.
The engine's Jobs page opens.

- From the Health Check list in the Jobs page, select Task Wait Time, which filters the list to display a list of jobs with longer than average wait times before the process is executed.

The screenshot shows the 'Jobs' page with a table of job details. The 'Health Check' dropdown is open, and 'Task Wait Time' is selected. The table columns include Type, Job, Status, Start Time, User, Health Issue, and Execution ID. The 'Health Issue' column shows 'Data Processing Speed Skew' and 'Abnormal Data Input'.

Type	Job	Status	Start Time	User	Health Issue	Execution ID
HV	insert over...	✓ Succeeded	04/19/2018 2:37 AM PDT		Data Processing Speed Skew	hive_...
HV	insert over...	✓ Succeeded	04/19/2018 2:37 AM PDT			hive_...
MR	aggregate...	✓ Succeeded	04/19/2018 2:37 AM PDT			job_...
HV	insert over...	✓ Succeeded	04/19/2018 2:37 AM PDT		Abnormal Data Input	hive_...
MR		✓ Succeeded	04/19/2018 2:37 AM PDT		Data Processing Speed Skew	job_...

- To view more details, from the Job column, select a job's name and then click the Health Checks tab. The Baseline Health checks are displayed.

- From the Health Checks panel, select the Task Wait Time health check.

For example, as shown in the following image, the long wait time occurred in the Map Stage of the job process due to insufficient resources:

The screenshot shows the 'Health Checks' panel for a job. The 'Task Wait Time' health check is selected. The panel displays a timeline of stages, a histogram of task wait times, and a list of top 20 outlier tasks. The 'Map Stage' is highlighted, and the histogram shows a long tail of wait times. The text indicates that the stage had resource starvation as 344 (out of 32992) tasks suffered.

Baseline

- ✓ Duration
- ✓ Input Size
- ✓ Output Size
- Skew
- ✓ Task Duration
- ! Task Input Data
- ! Task Output Data
- ✓ Shuffle Input
- ! Data Processing Speed
- Resources
- ! Task Wait Time
- ✓ Task GC Time
- ✓ Disk Spillage

Map Stage

Long Task Wait Time

This stage had **resource starvation** as 344 (out of 32992) tasks suffered

5s 3m 33s 7m 12s 10m

Median 7m 12s Typical Range 3m 33s - 10m

This health check determines if some tasks took too long to start a success indicates that successful tasks took less than 15 minutes and less than 4. Sufficient resources cut the run time of the job by lowering the maximum healthy, try giving more resources to the job by running it in resource pool more nodes to the cluster as a starting point.

Top 20 Outlier Tasks

Task Name	Host
! Task 032985	
! Task 032975	
! Task 032688	

8. To display more information about the Map Stage tasks that are experiencing longer than average wait times before they can execute, click one of the tasks listed under Outlier Tasks.

In the following example, the Task Details show that the task's wait time is above average. When comparing the Wait Duration value with the Successful Attempt Duration value, the task when it does finish has a significantly better than average time. This indicates that insufficient resources are allocated for this job.

The screenshot displays the 'Task Details' for a 'Map Stage' task (Task 032953). The interface includes a sidebar with navigation tabs: Overview, Health Checks, Execution Details, Baseline, and Trends. The 'Task Details' section shows a table with metrics for 'Wait Duration' and 'Successful Attempt Duration', both highlighted with orange boxes. The 'Wait Duration' is 15m 20s (7m 9s average) and the 'Successful Attempt Duration' is 45s 10ms (59s 287ms average).

Metric	Task	Average
Wait Duration	15m 20s	7m 9s
Duration	16m 5s	8m 8s
RECORDS_OUT_0	22K	12.8K
Input split bytes	551 B	797.73 B
Data written (HDFS)	6.5 MiB	5.2 MiB
Successful Attempt Duration	45s 10ms	59s 287ms
Map input records	6.1M	8M

Classifying Workloads for Analysis with Workload Views

The Workload View feature enables you to analyze workloads with much finer granularity. For example, you can analyze how queries that access a particular database or that use a specific resource pool are performing against your SLAs. Or you can examine how all the queries are performing on your cluster that are sent by a specific user.

Working with Auto Generated Workload Views

Steps for using the Workload XM default workload views.

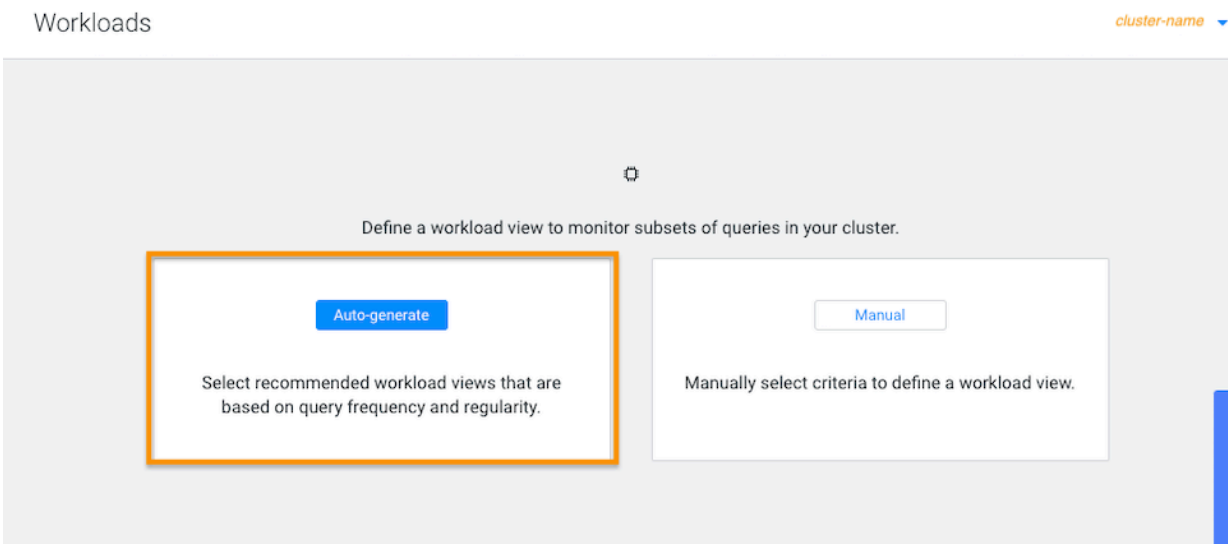
About this task

Describes how to use the workload views that Workload XM automatically generates.

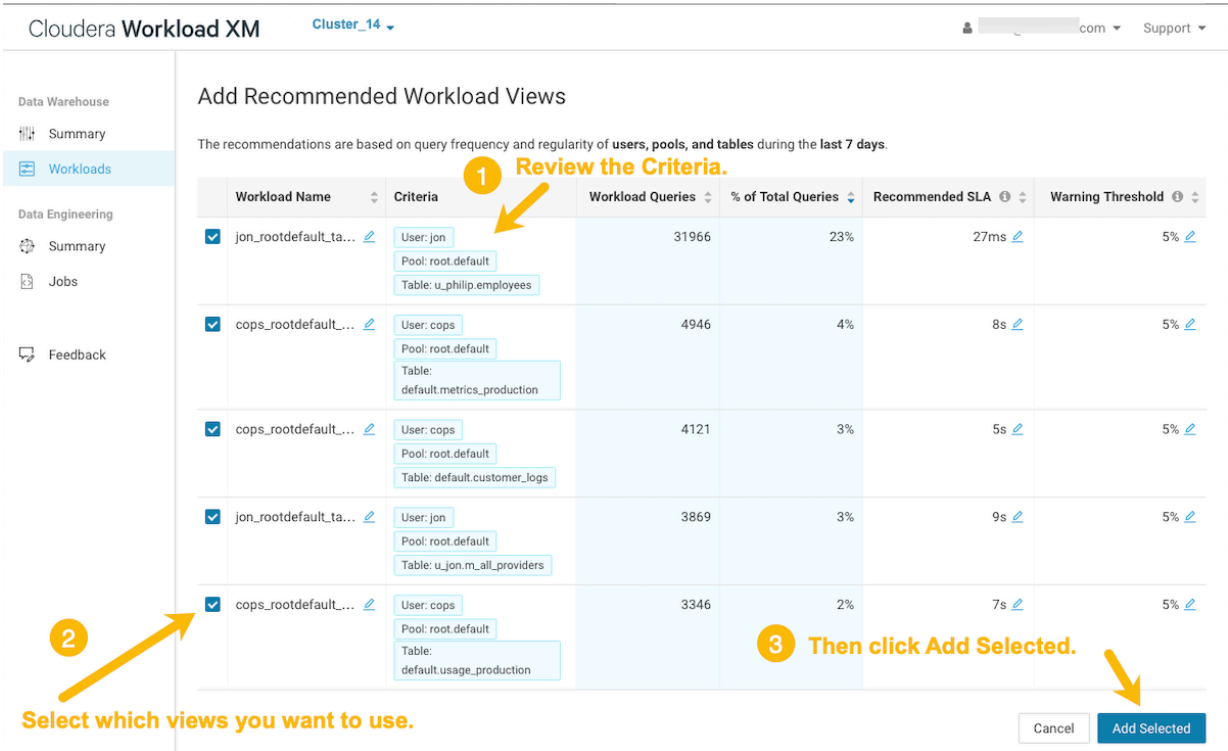
Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. In the Clusters page do one of the following:
 - In the Search field, enter the name of the cluster whose workloads you want to analyze.
 - From the Cluster Name column, locate and click on the name of the cluster whose workloads you want to analyze.
3. From the time-range list in the Cluster Summary page, select a time period that meets your requirements.

- 4. From the navigation panel, select Workloads.
- 5. In the Workloads page, click Auto-generate:



- 6. From the Criteria column, examine the criteria that are used to create the workload views, select the workload views required, and then click Add Selected:



The workload views you selected are saved and displayed on the Workloads page.

7. To verify your workload views, from the navigation panel, select Workloads and then on the Workload page locate the workload view you just added. When verified, click the workload to view its details:

The screenshot shows the Cloudera Workload XM interface for Cluster_14. The left navigation panel includes 'Data Warehouse' (Summary, Workloads), 'Data Engineering' (Summary, Jobs), and 'Feedback'. The main area is titled 'Data Warehouse Workloads' and contains a table of workloads. An orange arrow points to the workload named 'jon_rootdefault_tables_23_3869' with the text 'Click the workload name to view its details.'

Status	Workload	Criteria	SLA	Warning Thresh...	Missed SLA %	Failure %	Total Queries	Action
❌	MonitorSid	User: ANY OF ss, idalgic, mulyadi	30s	10%	23%	1%	2648	Actions ▾
✅	Clusterstats_Usage	Database: clusterstats	20s	20%	16%	1%	9884	Actions ▾
❌	Heavy_Users	User: ANY OF shashi, stephent, mulyadi, mstephenson	5s	10%	15%	13%	30263	Actions ▾
❌	M002	User: ANY OF mulyadi, vidya, brad, GABOR.SUDAR, abhishektalluri, abreshars, mkohs, raman, joydeep, ss, vmm	1m	1%	15%	1%	2648	Actions ▾
❌	jon_rootdefault_tables_23_3869	User: jon Pool: root.default Table: u_jon_m_all_providers	9s	5%	15%	0%	6016	Actions ▾
✅	ms	User: mstephenson	10s	20%	14%	2%	1147	Actions ▾
❌	cops_user	User: cops	5s	10%	12%	0%	2047472	Actions ▾
❌	DiagBundlesAnalysis	User: cops Statement Type: QUERY	5s	10%	12%	0%	1689191	Actions ▾

Defining Workload Views Manually

Steps for manually defining your workload views.

About this task

This task describes how to manually define your Workload Views.

To view this feature in action, watch the following video:

[Video: Classifying Workloads to Gain Insights](#)

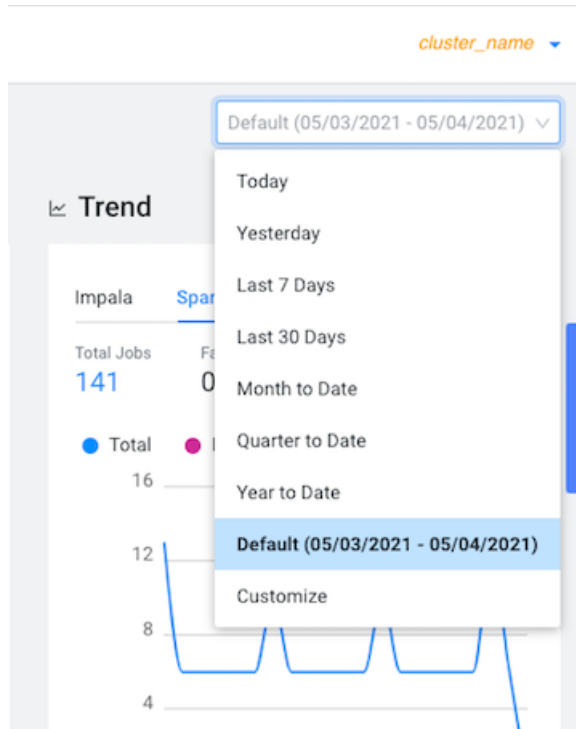
Figure 1: Video: Classifying Workloads to Gain Insights

For better video quality, click YouTube in the lower right corner of the video player to watch this video on YouTube.com.

Procedure

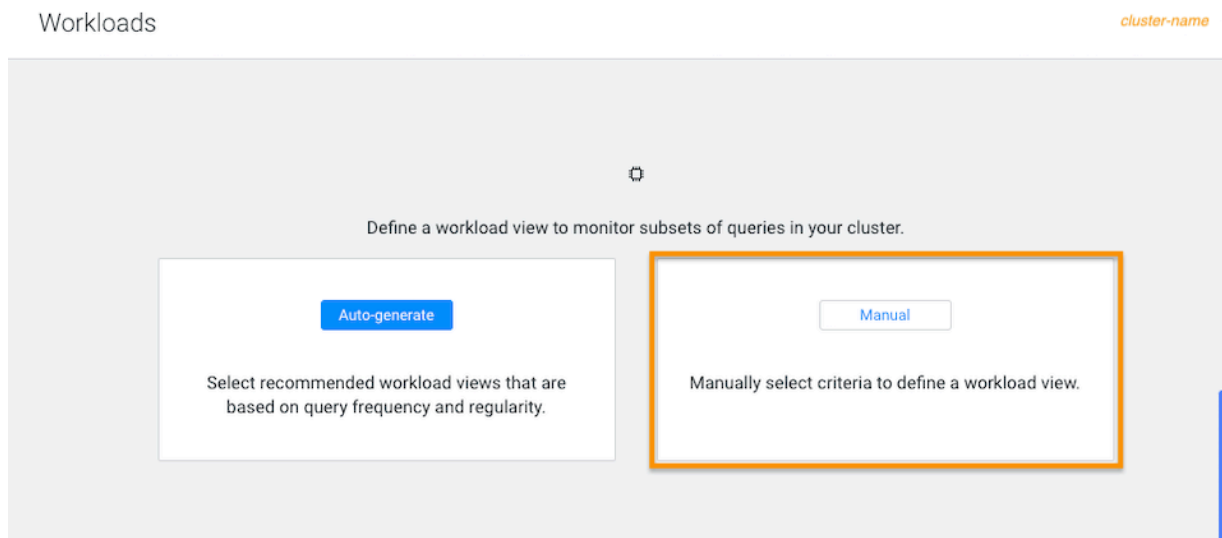
- In a supported browser, log in to the Workload XM web UI by doing the following:
 - In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - Click Log in.
- In the Search field of the Clusters page, enter the name of the cluster whose workloads you want to analyze.

3. From the time-range list in the Cluster Summary page, select a time period that meets your requirements.



4. From the navigation panel, select Workloads.

5. In the Workloads page, click Manual:



The Define Workload View widget opens, where you define a set of criteria that enables you to analyze a specific set of queries.

For example, as shown in the image below, you can list the total amount of failed queries, as a percentage, from a specific database that are subject to a fifteen second SLA.

Where, as defined by the criteria condition, Workload XM will monitor all query jobs from the applog database. When the total query execution time exceeds 15 seconds, as defined by the SLA condition, for 100 percent of these queries, as defined by the Warning Threshold, the workload is flagged with a failed state:

The screenshot shows the 'Define Data Warehouse Workload View' form. The form has the following fields and values:

- Name:** applog_db_under_15s
- Criteria:** Database = applog
- SLA:** 15s
- Warning Threshold:** 100 % queries missed SLA

Below the form, there is a 'Preview' button. Orange arrows point to the 'Name', 'Criteria', 'SLA', 'Warning Threshold', and 'Preview' fields.

- 6. (Optional) To display a summary of the queries matching your criteria, click Preview. As shown in the following example:

Preview

Cluster default date range is in the past, metrics reflect the status of the period.

Date range of queries in the workload. → 03/08/2018 - 05/03/2018

Summary of queries that match the criteria. ←

Total Queries	Missed SLA %
188	29%

- 7. When you are satisfied with the results, click Save.
The Workloads page opens and your workload view appears in the Workload column.

Cloudera Workload XMCluster_14.comSupport

Data WarehouseSummaryWorkloadsData EngineeringSummaryJobsFeedback

Data Warehouse Workloads

Status AllSearch workloads

Define New

Status	Workload	Criteria	SLA	Warning Thresh...	Missed SLA %	Failure %	Total Queries	Action
✓	applog_db_under_1...	Database = applog	15s	100%	29%	4%	188	Actions
✗	invest-db	Database = _int_prod	300ms	20%	40%	9%	30204	Actions
✗	invest-db-2	Database = _int_prod Statement Type = DDL	300ms	20%	32%	3%	21151	Actions

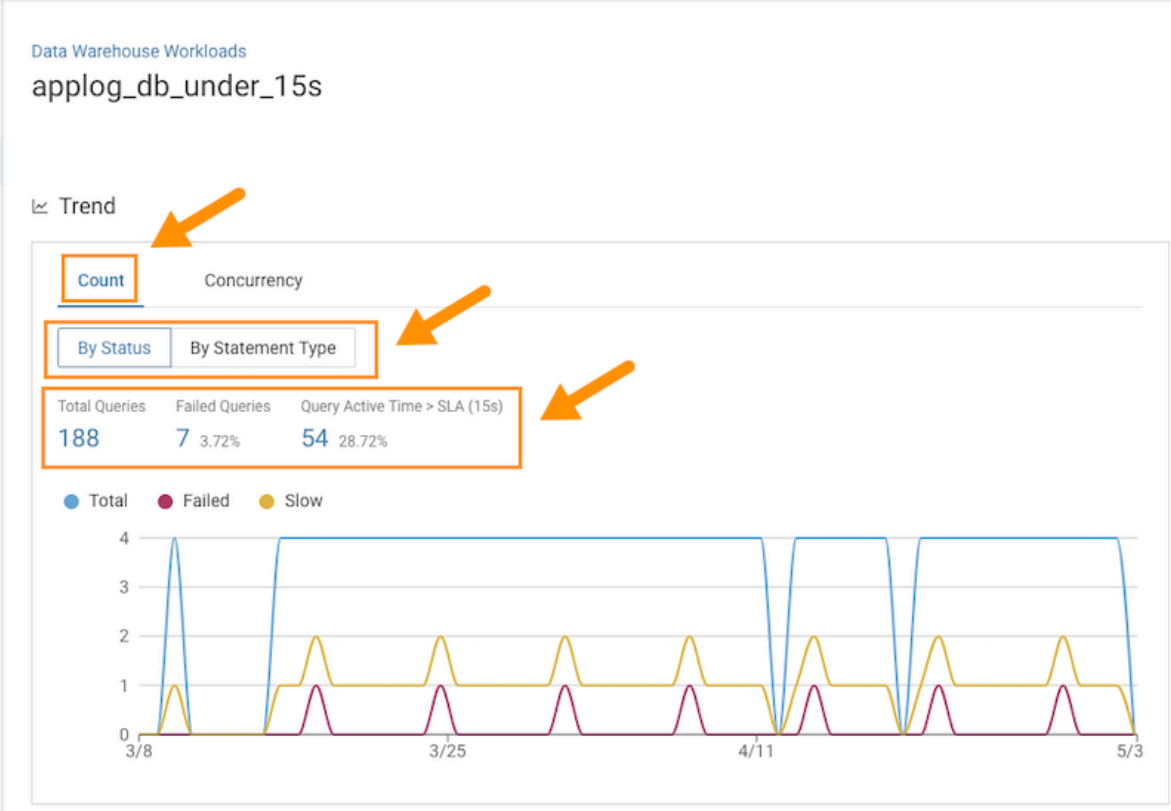


Tip: When you have a long list of Workload views, sorting the Workload column alphabetically in ascending or descending order by clicking the up or down arrows, helps locate the workload.

- 8. (Optional) To view more information about the workloads using the view's formula, open the workload's details page by clicking the name of the workload view in the Workload column, which visually displays the view's details as widgets that you can use to further analyze the results.

The following examples, display how this group of queries are meeting the Workload view's SLA in the Trend chart, where:

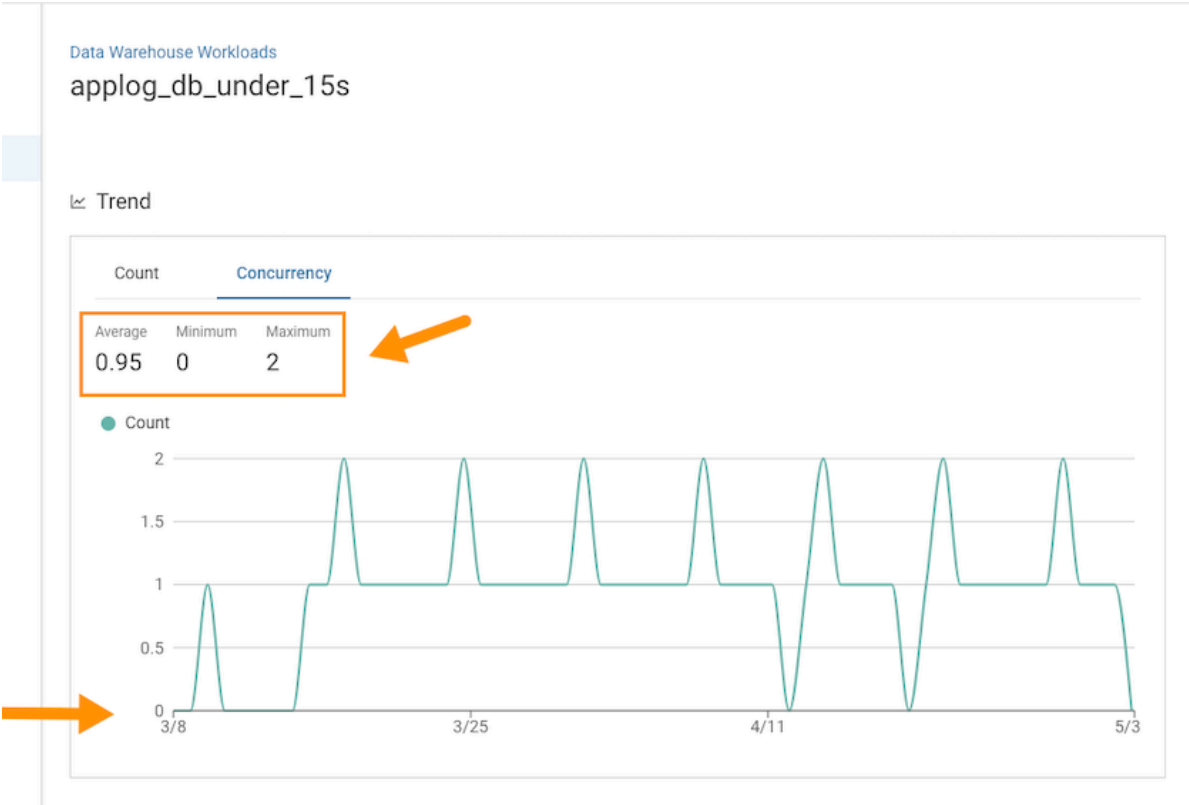
- The Count tab, displays the number of executing queries, either By Status or By Statement Type. To view further details, click the value under Total Queries, Failed Queries, and Query Active Time.



- The Concurrency tab, displays the number of queries executing concurrently.

In the following example, the maximum concurrency for this view is 2. This indicates that for the queries monitored by this view, only two queries accessed the same data at the same time during the specified time

period. The graph view displays how the concurrency fluctuates over the date range specified for the workload view.




Troubleshooting with the Job Comparison Feature

Steps for comparing two different runs of the same job, which is especially useful when you notice unexpected changes. For example, when you have a job that consistently completes within a specific amount of time and then it starts taking longer, comparing two runs of the same job enables you to analyze the differences so that you can troubleshoot the cause.

About this task

Describes how to compare any two runs of a job using the Job Comparison tool.

Steps with examples are included that help explain how to further investigate and troubleshoot.

 **Note:** When a job is flagged as slow, you can select the slow job from the Slow Jobs widget in the job's engine page and then in the details page, click Compare with Previous Run. The job is compared with its last run and the results are displayed the Job Comparison page for you to analyze.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. In the Search field of the Clusters page, enter the name of the cluster whose workloads you want to analyze.
3. From the time-range list in the Cluster Summary page, select a time period that meets your requirements.

- In the Trend widget, select the tab of an engine whose jobs you want to analyze and then click its Total Jobs value. The engine's Jobs page opens.

- Examine the list of jobs that have executed during the selected time period:

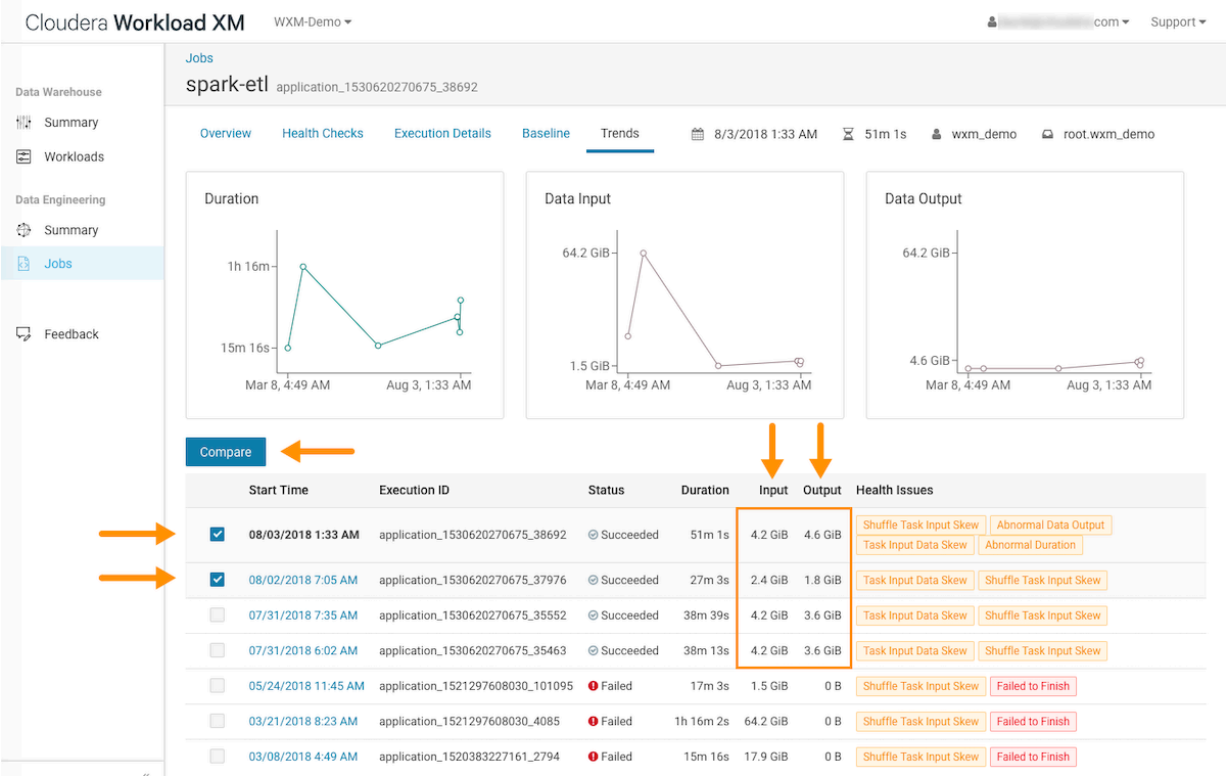
For example, as shown in the following image, the last three runs of the spark-etl job have taken significantly longer to complete than usual. Where, on August 2, the duration was 27 minutes, but on August 3, the duration almost doubled to 51 minutes:

User All ▾ Type All ▾ Status All ▾ Health Check All ▾ Duration All ▾ Range Default (03/08/2018 - 09/05/2018) ▾							
Type	Job	Status	Start Time	Duration	User	Health Issue	Execution ID
SP	Target Acc...	Failed	09/05/2018 3:04 AM PDT	18m 50s	wxm_demo	Failed to Finish	application_1527
SP	spark-etl	Succeeded	08/06/2018 9:46 AM PDT	51m 42s	wxm_demo	Shuffle Task Input Skew Abnormal Data Output Task Input Data Skew Abnormal Data Input Abnormal Duration	application_1530
SP	spark-etl	Succeeded	08/03/2018 8:37 AM PDT	51m 16s	wxm_demo	Shuffle Task Input Skew Abnormal Data Output Task Input Data Skew Abnormal Duration	application_1530
SP	spark-etl	Succeeded	08/03/2018 1:33 AM PDT	51m 1s	wxm_demo	Shuffle Task Input Skew Abnormal Data Output Task Input Data Skew Abnormal Duration	application_1530
SP	spark-etl	Succeeded	08/02/2018 7:05 AM PDT	27m 3s	wxm_demo	Task Input Data Skew Shuffle Task Input Skew	application_1530
SP	spark-etl	Succeeded	07/31/2018 7:35 AM PDT	38m 39s	wxm_demo	Task Input Data Skew Shuffle Task Input Skew	application_1530
SP	spark-etl	Succeeded	07/31/2018 6:02 AM PDT	38m 13s	wxm_demo	Task Input Data Skew Shuffle Task Input Skew	application_1530

- List and display details of all the runs of a specific job, by selecting one of the job runs and then in the Jobs details page, click the Trends tab.

In the following example, notice how the amount of data changes in the Input and Output columns. Where, on August 2, the job processed 2.4 GB of data, but on August 3, the job processed 4.2 GB, which is almost twice as

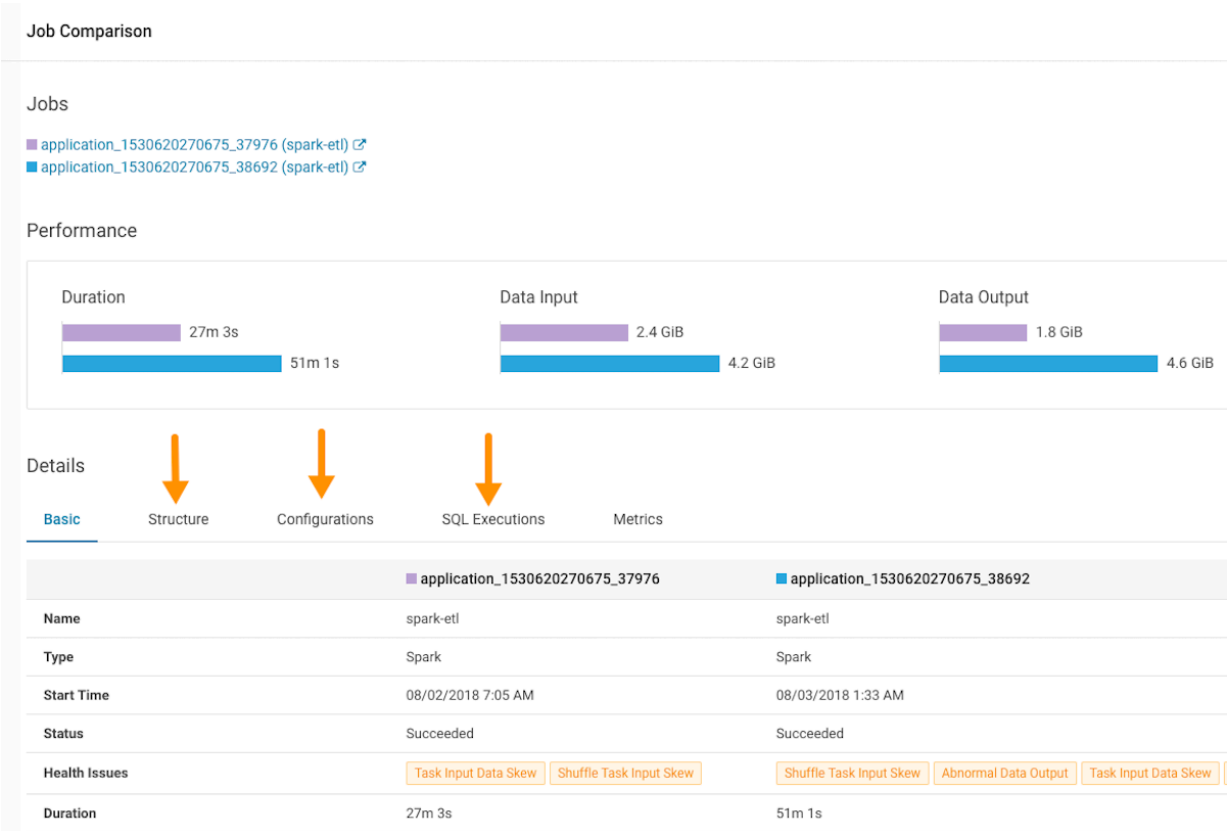
much data. The Job Comparison tool will enable you to examine both runs to determine why the amount of data changed:



7. To compare two job runs, select the check boxes adjacent to the job runs you require, in this case the runs for August 2 and August 3 are selected, and then click Compare.

The Job Comparison page opens displaying more details about each job.

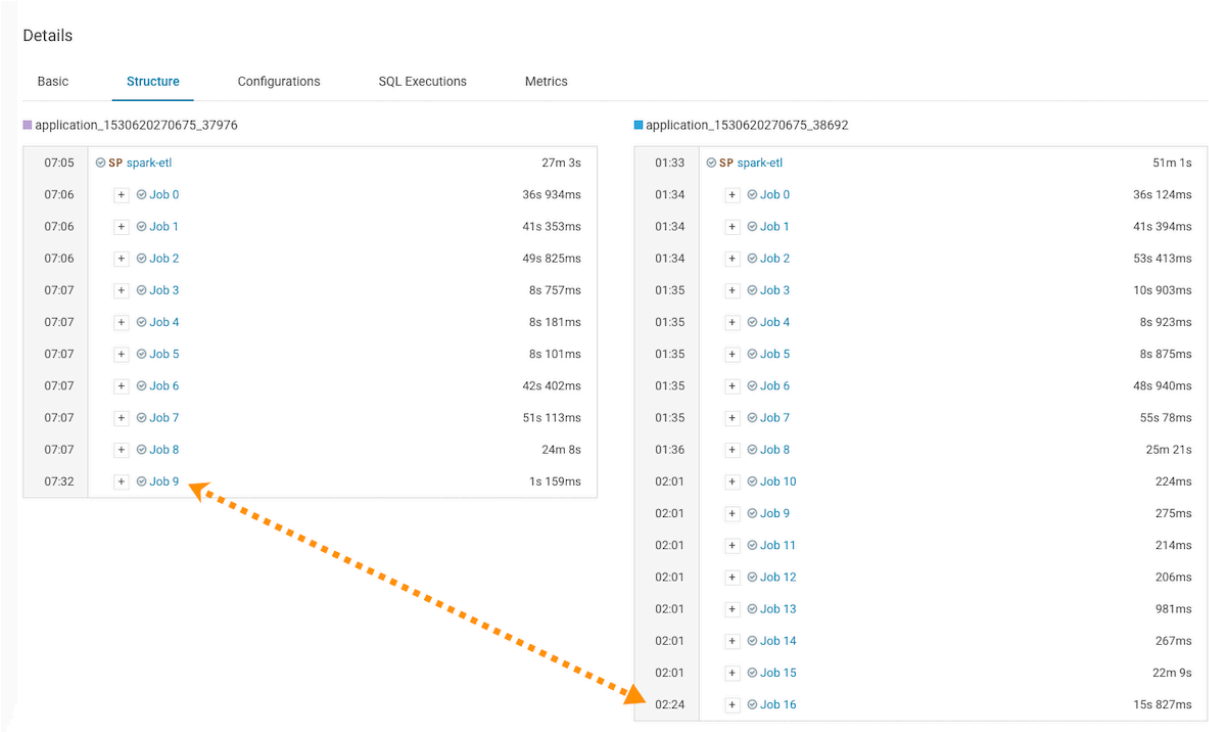
For this example's comparison, the tabs that contain more information about the job runs are the Structure, Configurations, and the SQL Executions tabs:



Note: The SQL Executions tab is only available for Spark jobs.

8. Display the sub-jobs executed for both of your selected job runs by selecting the Structure tab.

For example, as shown in the following image, the job that took 27 minutes only executed 9 sub-jobs and the job that took 51 minutes, almost twice as much time, executed 16 sub-jobs, almost twice as many. Clicking any of the listed sub-jobs displays more details.

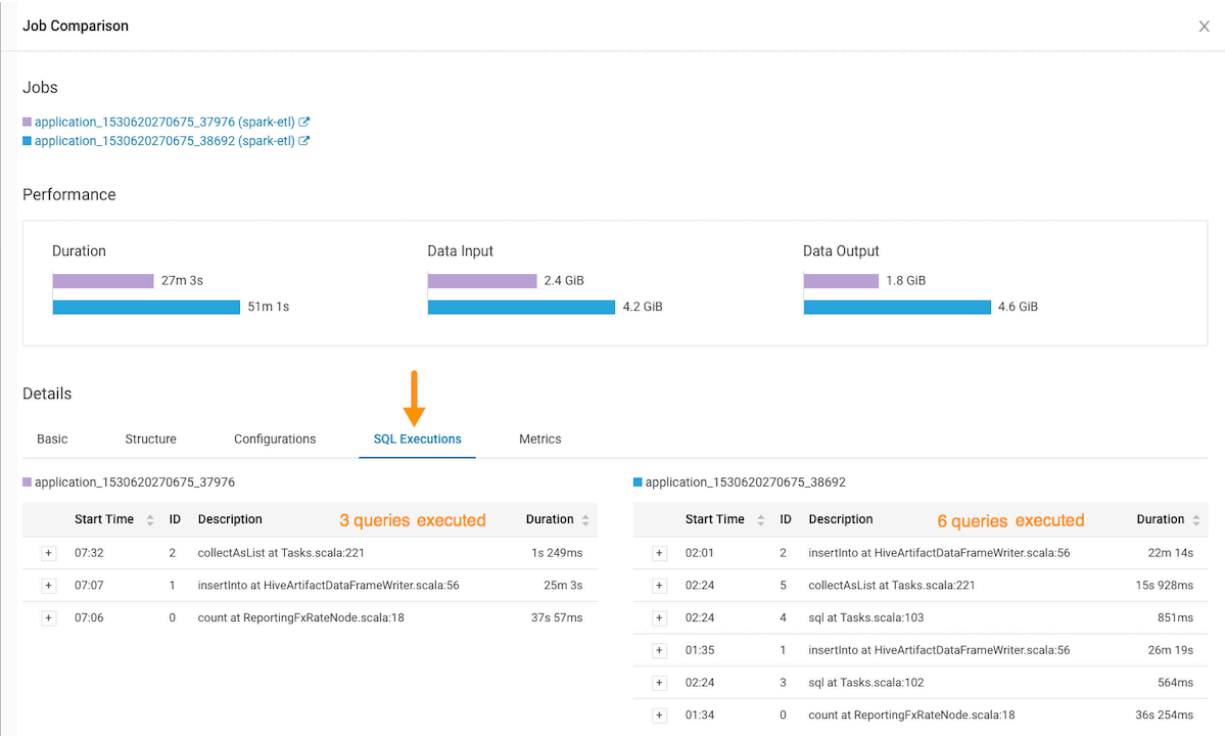


9. Display the jobs configurations by selecting the Configurations tab.

For our example, the configurations between the two runs of this job were identical, so a configuration change probably did not cause the anomaly.

10. Display the number of SQL executions by selecting the SQL Executions tab.

For our example, as shown by the following image, twice as many Spark queries executed for the job that took the longest duration.




Results

The analysis from the Job Comparison tells us that either the Spark SQL code was changed by the Job Developer or that the data on which the code ran triggered more of the Spark queries in the job. The Workload XM Job Comparison tool helped narrow the number of causes that produced the anomaly. For our example, the change in job duration appears to be expected so no further troubleshooting is required.


Identifying File Size Storage Issues

Data stored in small files or partitions may create performance issues. The File size reporting feature helps you identify data that is stored inefficiently in small files or partitions.

 **Important:** At this time the Workload XM File Size Report feature is only supported on CDH Workload clusters, version 6.3 to version 7.0, with Cloudera Navigator enabled. CDP Workload clusters are not supported.

A table's data maybe stored in a large number of files, perhaps millions of files. For example, the first time you run an Impala query, Impala also loads the metadata for each file, which can cause processing delays. In addition, every time you change a query, refresh the metadata, or add a new file or partition, Impala reloads the metadata. This puts pressure on the NameNode, which stores each file's metadata. For more information about the problems caused by small files and what you can do to fix those problems, see [Handling Small Files on Hadoop with Hive and Impala](#) on the Cloudera Engineering Blog.

The Workload XM file size reporting enables you to identify tables that have a large number of files or partitions. For example, for queries that run slowly or for Impala cluster crashes, you can view a table's metadata to determine whether a large number of files or partitions are causing the problem.

 **Note:** Before you can view the file size metadata in Workload XM, you must enable file size reporting in Cloudera Manager. Once enabled, the file size metadata is saved in HDFS, which is then forwarded to Workload XM by Telemetry Publisher.

Displaying File Size Metadata

Steps for displaying a table's File Size report and the metadata about the table's file's size distribution.

About this task

Describes how to open a table's File Size report and the metadata of a file.



Important: At this time the Workload XM File Size Report feature is only supported on CDH Workload clusters, version 6.3 to version 7.0, with Cloudera Navigator enabled. CDP Workload clusters are not supported.

Procedure

1. In a supported browser, log in to the Workload XM web UI by doing the following:
 - a) In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - b) When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - c) Click Log in.
2. In the Search field of the Clusters page, enter the name of the cluster whose workloads you want to analyze.
3. From the navigation panel under Data Warehouse, select File Size Report.
4. In the File Size Report page, either search for a specific table, or locate the table by sorting the tables by the number of files, the number of partitions, or the table size.

For example, the File Size Reports shows that the Animantarx table has 7 million files and 913 partitions.

Cloudera Workload XM

MyCluster ▼

triceratops@cloudera.com ▼

Support ▼

Data Warehouse

Summary

Workloads

File Size Report

Data Engineering

Summary

Jobs

Feedback

Table File Size Report

As of Mon, Apr 15, 2019 5:13 PM

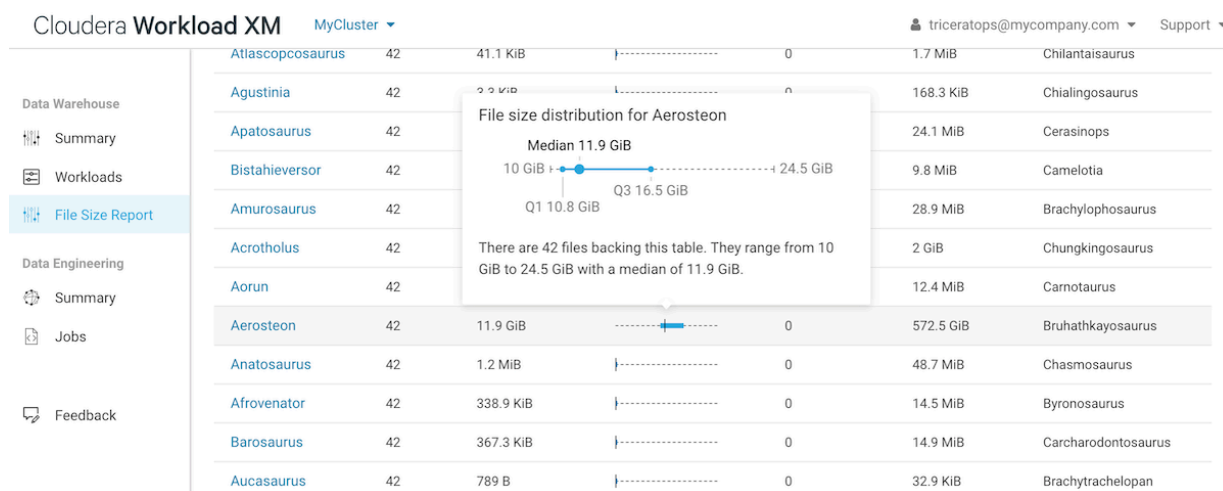
Search

Table and Db Name

Table	Files	Median File Size	File Size Distribution	Partitions	Table Size	Database
Animantarx	7M	36.7 KiB	-----	913	229.6 GiB	Carnotaurus
Bonapartenykus	3.1M	1 MiB	-----	397.3K	3.3 TiB	Bruhathkayosaurus
Balaor	1.7M	469 KiB	-----	1K	1.7 TiB	Chasmosaurus
Alwalkeria	595.3K	2.5 MiB	-----	1.7K	1.4 TiB	Cetiosaurus
Atlasaurus	401.8K	1.2 KiB	-----	4	477.6 MiB	Chilantaisaurus
Angolatitan	358.9K	168 KiB	-----	7.1K	455.9 GiB	Cerasinops
Anatosaurus	346.9K	1.9 KiB	-----	5.1K	27.3 GiB	Byronosaurus

- To display details about the table's file size distribution, select a table name.

For example, the following table's details window shows that the Aerosteon table uses 42 data files that range from 10 to 24.5 GiB and the graph displays the Q1 and Q3 file size distribution.



Displaying the Metadata of a Table

Steps for displaying a table's metadata that could be causing a query to run slowly.

About this task

Describes how to display the metadata of table used in your query, such as the table's file size distribution that could be causing your query statement to run slowly.

Procedure

- In a supported browser, log in to the Workload XM web UI by doing the following:
 - In the web browser URL field, enter the Workload XM URL that you were given by your system administrator and press Enter.
 - When the Workload XM Log in page opens, enter your Workload XM user name and password access credentials.
 - Click Log in.
- In the Search field of the Clusters page, enter the name of the cluster whose workloads you want to analyze.
- From the navigation panel under Data Warehouse, select Summary.

4. In the Queries page, select the query of interest and then select the HDFS Tables Scanned tab.

For example, the Duration column shows that the query took over six hours to run and the HDFS Tables Scanned section displays the metadata for the tables that were scanned.



Note: This is not the number of files accessed, but the total number of files that were in the table the last time a HDFS snapshot was taken before the query was run.

Cloudera Workload XM MyCluster ▾

triceratops@cloudera.com ▾ Support ▾

Data Warehouse

- Summary
- Workloads
- File Size Report

Data Engineering

- Summary
- Jobs

Feedback

Queries

1c4d1a6ca67e94f4:7977bd4d00000000

Summary Trend ✓ Succeeded 04/24/2019 2:05 AM CDT stegosaurus DML root.default Profile

Joins	Duration	Rows Produced	Aggregate CPU Usage	Aggregate Memory Usage	Peak Memory Usage
2	6h 34m 3s	571	41h 29m 23s	30.6 TiB × s	1.5 GiB

Basic Operators Hosts **HDFS Tables Scanned**

As of Mon, Apr 15, 2019 5:13 PM

! The table metadata below is taken from the last HDFS snapshot before the query was run.

Table	Files	Median File Size	File Size Distribution	Partitions	Table Size	Database
Aragosaurus	739	39.5 KiB	+	0	78.8 MiB	Brachytrachelopan

5. To display the file size distribution details for a table, click the Table name .

Assigning Access Roles in Workload XM

Workload XM supports cluster privilege role types that define who is entitled to access jobs and queries that are created by the user, who is entitled to create and administer cost centers and view cluster costs, and who is entitled to access and administer jobs and queries within either a specific cluster or across all clusters within the Workload XM environment.

Limiting the trust boundary for jobs, queries, cluster costs, and administrative management at the cluster level, enables more control over the security and access management of your Workload XM environment.

Understanding the Workload XM Access Roles

Describes the Workload XM access roles.



Important: Customers are responsible for managing and reviewing access credentials for their Workload XM accounts and activities. All user privileges and access rights should periodically be reviewed and monitored, including who should access Workload XM, its services, and components. For example, access rights should be reviewed when a user moves to another business unit.

Workload XM supports cluster privilege roles that define Workload XM users as a:

- System Admin
- Cluster Admin
- Cluster User

The following tables describe these cluster privilege roles, also known as access roles:

System Admin Access Role

An authentic Workload XM user who is assigned the System Admin access role has full access rights and system administrator privileges across all clusters within the Workload XM environment. Where they can view, edit, and create cost centers, view, edit, and create auto actions, and view all the jobs and queries in all the Workload clusters. These users have the least restrictive access permissions.

Table 1: System Admin

Resource	Actions
Access Management page	View and manage all the Workload XM cluster policies and user access from the Access Management page
Cluster	<ul style="list-style-type: none"> View all the workload clusters on the Clusters page Rename a workload cluster Delete a workload cluster
Workloads	<ul style="list-style-type: none"> Create workloads View all the workloads in a cluster Update all the workloads in a cluster Delete all the workloads in a cluster
Queries	View all the queries in all the clusters of the Workload XM environment
Jobs	View all the jobs in all the clusters of the Workload XM environment
Chargeback	<ul style="list-style-type: none"> Create cost centers Update cost centers List cost centers Delete cost centers View all the Chargeback related dashboards
Auto Actions	<ul style="list-style-type: none"> Create auto actions View auto actions Update auto actions Disable auto actions Delete auto actions Enable an auto action email

Cluster Admin Access Role

An authentic Workload XM user who is assigned the Cluster Admin access role has full access rights and cluster administrator privileges across an assigned cluster within the Workload XM environment. Where they can view all the jobs and queries in the assigned Workload cluster.

Table 2: Cluster Admin

Resource	Actions
Cluster	<ul style="list-style-type: none"> View the assigned Workload cluster on the Clusters page Rename the Workload cluster Delete the Workload cluster
Workloads	<ul style="list-style-type: none"> Create workloads View all workloads in the assigned cluster Update all workloads in the assigned cluster Delete all workloads in the assigned cluster
Queries	View all the queries in the assigned cluster
Jobs	View all the jobs in the assigned cluster

Cluster User Access Role

An authentic Workload XM user who is assigned the Cluster User access role has limited access rights across an assigned cluster within the Workload XM environment. Where they can view only those jobs and queries they created and executed in the assigned Workload cluster.

Table 3: Cluster User

Resource	Actions
Cluster	View their assigned cluster on the Clusters page.
Workloads	View their assigned workloads on the Workloads page
Queries	View their queries in the assigned cluster
Jobs	View their jobs in the assigned cluster

The Cluster User access role type has the most restricted access permissions, where the user may only view their own jobs and queries.

This access role further restricts the Cluster User to one cluster per policy. For users who are responsible for jobs and queries in more than one cluster they must also be assigned access rights to those clusters. You can either add them to the Cluster Policy for that cluster or include the pool that contains those workloads in the Cluster Policy in which they are assigned.

Also, for users who require access to jobs and queries executed by other users, you can create a Custom Policy as part of the Cluster Policy. This policy includes the user names of the users who execute those jobs and queries and/or the pool names in which they are executed.

For example, though user A and user B have been granted the same Cluster User role type their access to jobs and queries is different. This is due to the conditions of the Cluster Policy in which they are assigned. Where:

- The cluster policy that defines user A's Cluster User role type does not permit the user to view workloads within a pool or view other user workloads. In this case, user A is restricted to only view their own jobs and queries within their policy's assigned cluster.
- The cluster policy that defines user B's Cluster User role type contains a Custom Policy that permits the user to view workloads within a pool and view other user workloads. In this case, user B can view the jobs and queries executed by other users and the jobs and queries executed in the pool.

Understanding a Workload XM Cluster Policy

Describes the Workload XM Cluster Policy criteria that is used to assign Workload XM access roles to your users.

Access to your Workload jobs and queries is determined by a Workload XM Cluster Policy, which comprises two or more of the following conditions:

- One or more LDAP Group identifier account names.
- One or more user names. By default, Workload XM authenticates user access by checking that the user is a member of an LDAP group.
- A Workload XM access role type. The access role is assigned to the users that you provide in the Users field and/or the users who are part of the groups you provide in the Groups field and is defined by the conditions in the Cluster Policy.
- (Cluster User and Cluster Admin only) The cluster associated with the access role.
- (Cluster User only) A custom policy whose criteria is defined from the provided user names and/or the provided pools. A custom policy enables the user or users defined in the Cluster Policy to view the jobs and queries executed by other users and/or the jobs and queries executed in a pool.

Workload XM Cluster Policies are created, managed, and maintained from the Access Management page. Only users who have been granted the System Admin access role type can view and manage your Workload XM cluster policies.

Configuring a Default Systems Administrator for Workload XM

Pre-tasks that are required before you can start enabling role based access in Workload XM.

About this task

Describes how to enable role based access in Workload XM and configure a Workload XM default systems administrator.

Before you can assign access roles in Workload XM you must first enable role based access and configure a default systems administrator. Both tasks are completed in Cloudera Manager. Once configured, the default administrator (also known as a superuser) can log into the Workload XM UI and assign the System Admin access policy role to one or more users.

Procedure

1. In a supported web browser on the Workload XM on-premises cluster, log in to Cloudera Manager.
2. In Cloudera Manager, select Clusters, WXM, and then click the Configuration tab.
3. In the Configuration page, search for the Role Based Access enabled property and then select its WXM (Service-Wide) check box.
4. According to your requirements, do one of the following:
 - a. In the WXM (Service Wide) field of the WXM Default Super Users property, enter either the user name or the account name of a system administrator who is to be granted access to perform administration tasks in Workload XM. By default, admin.



Tip: If the WXM (Service Wide) field is not displayed, click the plus sign circle icon.

- b. In the WXM (Service Wide) field of the WXM Default Super Groups property, enter the group account name of your LDAP admin group. For example, admin_grp.

The following image shows the configuration properties:

5. Click Save Changes.
6. Navigate to the top of the Workload XM service page and from the Actions menu, restart the Workload XM service, by selecting Restart.

Assigning Workload XM Access Roles

Role based access to your Workload jobs and queries requires a Workload XM Cluster Policy that defines the conditions for the role based access type and assigns it to your users. You can have multiple Cluster Policies that define the access criteria for all of your workloads.

Assigning a Workload XM System Admin Access Role

Steps for assigning a System Admin access role to your Workload XM users.

About this task

Describes how to assign a Workload XM Role Based Access (RBAC) role for a system administrator. This access role has full access rights and system administrator privileges across all clusters within the Workload XM environment and can create your Workload XM Cluster Policies that define your access roles.



Note: Generally, only a user assigned the System Admin access role can create a Workload XM Cluster Policy. But until the first System Admin access role is assigned, a Cluster Policy can only be created by a default systems administrator, also known as a default super user.

Before you begin

This task assumes that you have:

- Enabled role based access in Cloudera Manager.
- Created a default systems administrator, also known as a default super user, in Cloudera Manager.

Procedure

1. In a supported browser, log in to Workload XM as the user with default system administrator privileges.
2. From the Workload XM Navigation side-bar, select Access Management.
3. In the Access Management page, click New Cluster Policy.
The Create Cluster Policy page opens.
4. Do one or more of the following:
 - a. In the Groups field, enter the name of the LDAP administration group account whose users will be assigned this cluster policy's access role.
 - b. In the Users field, enter the user name or user names who will be assigned this cluster policy's access role.
5. From the Assign Roles list, select System Admin.
6. Click Create.



Note: Workload XM will take at least 60 minutes to assign the access role to the user, users, and/or groups provided in the Cluster Policy.

Results

The Successfully created access policy message appears when the Cluster Policy is created and the policy is displayed in the Access Management's home page.

Assigning a Workload XM Cluster Admin Access Role

Steps for assigning a Cluster Admin access role to your Workload XM users.

About this task

Describes how to assign a Workload XM Role Based Access (RBAC) role for a cluster administrator.



Note: Only a user assigned the System Admin access role can create a Workload XM Cluster Policy.

Procedure

1. In a supported browser, log in to Workload XM as a user that has been granted the System Admin access role.
2. From the Workload XM Navigation side-bar, select Access Management.
3. In the Access Management page, click New Cluster Policy.
The Create Cluster Policy page opens.
4. Do one or more of the following:
 - a. In the Groups field, enter the name of the LDAP group account whose users will be assigned this cluster policy's access role.
 - b. In the Users field, enter the user name or user names who will be assigned this cluster policy's access role.
5. From the Assign Roles list, select Cluster Admin.
6. From the Cluster list, select the name of the cluster that will be assigned to this policy's access role.
7. Click Create.



Note: Workload XM will take at least 60 minutes to assign the access role to the user, users, and/or groups provided in the Cluster Policy.

Results

The Successfully created access policy message appears when the Cluster Policy is created and the policy is displayed in the Access Management's home page.

Assigning a Workload XM Cluster User Access Role

Steps for assigning a Cluster User access role to your Workload XM users.

About this task

Describes how to assign a Workload XM Role Based Access (RBAC) role for a cluster user.



Note: Only a user assigned the System Admin access role can create a Workload XM Cluster Policy.

Procedure

1. In a supported browser, log in to Workload XM as a user that has been granted the System Admin access role.
2. From the Workload XM Navigation side-bar, select Access Management.
3. In the Access Management page, click New Cluster Policy.
The Create Cluster Policy page opens.
4. Do one or more of the following:
 - a. In the Groups field, enter the name of the LDAP group account whose users will be assigned this cluster policy's access role.
 - b. In the Users field, enter the user name or user names who will be assigned this cluster policy's access role.
5. From the Assign Roles list, select Cluster User.
6. From the Cluster list, select the name of the cluster that will be assigned to this policy's access role.
7. (Optional) Enable the user or users defined in this cluster policy to view executed workloads from other users or executed workloads from a pool by doing the following:
 - a. In the Users field, enter the user name or user names whose jobs and queries can be viewed by the user or users defined in this cluster policy.
 - b. In the Pools field, enter the pool name or pool names whose jobs and queries can be viewed by the user or users defined in this cluster policy.

8. Click Create.



Note: Workload XM will take at least 60 minutes to assign the access role to the user, users, and/or groups provided in the Cluster Policy.

Results

The Successfully created access policy message appears when the Cluster Policy is created and the policy is displayed in the Access Management's home page.

Managing Your Workload XM Access Roles

Describes how to manage your Workload XM cluster policies and access roles.

Information about your Workload XM Cluster Policies are displayed on the Access Management page, which are viewed and managed by the user with the System Admin access role.

Each row displays a Cluster Policy and its conditions, where:

- The Status column displays the state of the policy, as either Enabled or Disabled.
- The Clusters column displays the name of the cluster assigned to the Workload XM access role.
- The Role column displays the Workload XM access role type.
- The Groups column displays the LDAP group users who are assigned the Cluster Policy's access role.
- The Users column displays the user names who are assigned the Cluster Policy's access role.
- The Custom Policy column displays the user and pool filter conditions.
- The Last Updated column displays the date when the policy was last updated.
- The Actions column's vertical ellipses, when selected, lists the management tasks that can be performed.

The following management tasks are performed from the Access Management home page by a user with the System Admin access role, which is accessed by selecting Access Management from the Workload XM Navigation side-bar.

Updating a Cluster Policy

In the Access Management page, click the cluster policy's vertical ellipsis in the Actions column, and select Edit. In the Cluster Policy, make your changes and then click Update.

Deleting a Cluster Policy

In the Access Management page, click the cluster policy's vertical ellipsis in the Actions column, and select Delete. In the confirmation message, click OK to confirm the action. The policy is permanently removed.



Tip: Cloudera recommends disabling rather than deleting a Cluster Policy.

Disabling a Cluster Policy

In the Access Management page, click the cluster policy's vertical ellipsis in the Actions column, and select Disable. In the confirmation message, click OK to confirm the action. The Status column displays the state of the policy as Disabled.

Purging HDFS Data

Reduce bottlenecks between Telemetry Publisher and Workload XM, free up storage space, and increase job and query runtime efficiency by removing obsolete HDFS data that exceeds the maximum retention limit.



Note: Cloudera recommends performing regular purge events for HDFS files that are no longer required.

Understanding the Purge Date used by the Purge Event

Describes the Workload XM purge event's criteria that is based on the file's data group and the data group's retention limit and how the purge date is calculated.

The purge event's criteria is based on the maximum data retention policy, described in days, for the following HDFS data groups:

- Temporary data, when the retention period exceeds 8 days
- Staging data, when the retention period exceeds 31 days
- Detailed data, when the retention period exceeds 181 days
- Summarized data, when the retention period exceeds 731 days

The purge date is calculated by subtracting the retention days, specified by the maximum data retention period policy, from the current date and comparing the resultant date with the data's timestamp date. If the data's timestamp date is less than or equal to the resultant date the data is removed.

The data's timestamp date is determined by where the data resides:

- If the data resides in the cloudera-bus root directory, the timestamp date is extracted from the subdirectory name. For example, if the directory name is /cloudera-dbus/HiveAudit/2021030623. The timestamp date extracted by the purge event is 2021/03/06, using the YYYY/MM/DD date format.



Important: The purge event deletes files from the cloudera-dbus directory as follows:

- If the date is successfully extracted and is less than or equal to the resultant date, all the files in the directory are removed and are counted as one file by the maximum deletion limit.
- If the date is successfully extracted, is less than or equal to the resultant date, and a file or files are set in the blobstore.purger.paths.to.keep parameter, all the files except the file or files set in the blobstore.purger.paths.to.keep parameter are removed and each file that is removed is counted by the maximum deletion limit.
- If the data resides in a cloudera-sigma-olap-impala, cloudera-sigma-partial-pse, cloudera-sigma-pse-extended, or cloudera-sigma-sdx-payloads root directory, the timestamp date is extracted from the file's last modified time.



Obsolete data can be purged from the following HDFS root directories:

- cloudera-dbus
- cloudera-sigma-olap-impala
- cloudera-sigma-partial-pse
- cloudera-sigma-pse-extended
- cloudera-sigma-sdx-payloads

Workload XM Purge Event Parameters

Lists the Workload XM purge event parameter settings that enable you to set the event's execution time, frequency, and maximum purge duration. You can also exclude files and directories from being purged with the blobstore.purger.paths.to.keepparameter setting.

Table 4: Purge Event Parameters

Parameter	Description	Example
blobstore.purger.frequency	<p>The purge event's recurring schedule, based on one of the following values:</p> <ul style="list-style-type: none"> None. By default, the purge process is set to none. Daily. When this value is set for the first time, files are automatically deleted the next day at 1am. Weekly. By default, files are automatically deleted every Saturday at 1am. Monthly. When this value is set for the first time, files are automatically deleted the last Saturday of the month at 1am. Thereafter, files are deleted every 28th day. The monthly parameter uses the 28 day calendar format 	blobstore.purger.frequency = none
blobstore.purger.start.time	<p>The purge event's start time, based on the 24-hour time format. Where, 01:00 and 0:00 are valid time values, and 24:00, 1:0, and 01:0 are not valid time values</p> <p>By default, Workload XM schedules the purge process when it will cause the least amount of disruption to users.</p> <p> Note: Cloudera recommends scheduling a time during non-peak working hours or job execution hours.</p>	blobstore.purger.start.time = 01:00
blobstore.purger.paths.to.keep	<p>Lists the files and directories that are to be excluded from the purge event.</p> <p>Where each file and/or directory is separated by a comma and where:</p> <ul style="list-style-type: none"> a file value must use its full path, directory name, and file name. a directory value must use its full path and directory name. 	blobstore.purger.paths.to.keep=/cloudera-dbus/ImpalaQueryProfile/2021030217/7d2bcefa-8819-4fa1-be0c-4529ee4eb98f/cloudera-dbus/HiveAudit/cloudera-sigma-olap-impala/02f54999-b9a4-4dca-8237-d1b04775efb/cloudera-sigma-sdx-payloads/2bc85719-7a3e-4438-96a4-8fc0f77ff
blobstore.purger.delete.request.limit	<p>The maximum deletion limit.</p> <p>By default, the maximum number of files that can be deleted by the purge process is 500,000. This ensures that a purge cycle is not overloaded, does not introduce bugs, or takes up too much time.</p> <p>When the deletion limit is met, the purge process:</p> <ul style="list-style-type: none"> Stops processing for a daily scheduled value. Stops processing and restarts the next day for all other scheduled values. <p> Note: The purge event's maximum deletion limit calculates all the files in a dbus directory as one file. When you exclude a file or files that reside in the dbus directory from the purge process, the purge event's maximum deletion limit condition calculates all the files in the directory minus those files you have excluded.</p>	blobstore.purger.delete.request.limit=500000

Configuring the Workload XM Purge Event

Steps for scheduling and configuring a purge event.

About this task

Describes how to schedule and configure the Workload XM purge event.

Procedure

1. In a supported web browser, log in to Cloudera Manager as a user with full system administrator privileges.
2. From the Navigation panel, select Clusters and then WXM.
3. In the Status Summary panel of the WXM page, select Admin API Server.
4. Click the Configuration tab.
5. Search for the Admin API Server Advanced Configuration Snippet (Safety Valve) for the wxm-conf/sigmaadminapi.properites option.
6. In the text field enter your purge event's parameter settings, using the *Purge Event Parameters* table.

For example,

```
blobstore.purger.delete.request.limit=9990000
blobstore.purger.paths.to.keep=/cloudera-dbus/ImpalaQueryProfile/202103021
7/7d2bcefa-8819-4fa1-be0c-4529ee4eb98f,/cloudera-dbus/HiveAudit,/cloudera-
sigma-olap-impala/02f54999-b9a4-4dca-8237-d1b047755efb,/cloudera-sigma-sdx
-payloads/2bc85719-7a3e-4438-96a4-8fc0f77ff79e
blobstore.purger.frequency=daily
blobstore.purger.start.time = 0:00
```

7. Click Save Changes, which sets and schedules the purge process.
8. From the Actions menu, select Restart this Admin API Server.
9. In the Restart this Admin API Server message, confirm your changes by clicking Restart this Admin API Server.
10. When the Restart API Server step window displays Completed, click Close.

Manually Executing a Workload XM Purge Event

You can manually run your purge event immediately with a one-time operation, rather than scheduling a purge event.

About this task

Describes how to manually run a Workload XM purge event.

A one-time purge event is based on the maximum data retention policy using the Workload XM purge event's parameter values, without the frequency value.

Procedure

1. In a supported web browser, log in to Cloudera Manager as a user with full system administrator privileges.
2. From the Navigation panel, select Clusters and then WXM.
3. In the Status Summary panel of the WXM page, select Admin API Server.
4. Click the Configuration tab.
5. Search for the Admin API Server Advanced Configuration Snippet (Safety Valve) for the wxm-conf/sigmaadminapi.properites option.

6. In the text field enter your purge event's parameter settings, using the *Purge Event Parameters* table.

For example,

```
blobstore.purger.delete.request.limit=9990000
blobstore.purger.paths.to.keep=/cloudera-dbus/ImpalaQueryProfile/202103021
7/7d2bcefa-8819-4fa1-be0c-4529ee4eb98f,/cloudera-dbus/HiveAudit,/cloudera-
sigma-olap-impala/02f54999-b9a4-4dca-8237-d1b047755efb,/cloudera-sigma-sdx
-payloads/2bc85719-7a3e-4438-96a4-8fc0f77ff79e
blobstore.purger.frequency=none
blobstore.purger.start.time = 0:00
```

7. Click Save Changes.
8. From the Actions menu, select Restart this Admin API Server.
9. In the Restart this Admin API Server message, confirm your changes by clicking Restart this Admin API Server.
10. When the Restart API Server step window displays Completed, click Close.
11. When a manual purge event run is required, do the following:
 - a) Log in to Cloudera Manager.
 - b) From the Navigation panel, select Clusters and then WXM.
 - c) From the Actions menu, select Purge HDFS Bucket Data.
 - d) In the Purge HDFS Bucket Data confirmation message, confirm the purge event by clicking Purge HDFS Bucket Data.
 - e) When the Purge HDFS Bucket Data window displays Completed, click Close.

Managing your Workload XM Purge Event

Steps for updating, stopping, and troubleshooting your Workload XM Purge event.

The following management tasks can be performed:

Updating your Workload XM Purging Event

To update your purge event:

1. In a supported web browser, log in to Cloudera Manager as a user with full system administrator privileges.
2. From the Navigation panel, select Clusters and then WXM.
3. From the Status Summary panel, select Admin API Server.
4. Click the Configuration tab.
5. Search for the Admin API Server Advanced Configuration Snippet (Safety Valve) for the wxm-conf/sigmaadminapi.properites option field.
6. In the text field, change the required values.
7. Click Save Changes.
8. From the Actions menu, select Restart this Admin API Server.
9. In the Restart this Admin API Server message, confirm your changes by clicking Restart this Admin API Server.
10. When the Restart API Server step window displays Completed, click Close.

Stopping the Workload XM Purge Event

You can stop a recurring purge event or stop a scheduled purge event whilst still running.

- To stop a recurring purge event:
 1. In a supported web browser, log in to Cloudera Manager as a user with full system administrator privileges.
 2. From the Navigation panel, select Clusters and then WXM.
 3. From the Status Summary panel, select Admin API Server.
 4. Click the Configuration tab.
 5. Search for the Admin API Server Advanced Configuration Snippet (Safety Valve) for the wxm-conf/sigmaadminapi.properites option field.
 6. In the text field, replace the blobstore.purger.frequency value with none.
 7. Click Save Changes.
 8. From the Actions menu, select Restart this Admin API Server.
 9. In the Restart this Admin API Server message, confirm your changes by clicking Restart this Admin API Server.
 10. When the Restart API Server step window displays Completed, click Close.
- To stop a scheduled purge event whilst still running:
 1. In a supported web browser, log in to Cloudera Manager as a user with full system administrator privileges.
 2. From the Navigation panel, select Clusters and then WXM.
 3. From the Status Summary panel, select Admin API Server.
 4. From the Actions menu, select Stop this Admin API Server.
 5. Still in the Admin API Server page, click the Configuration tab.
 6. Search for the Admin API Server Advanced Configuration Snippet (Safety Valve) for the wxm-conf/sigmaadminapi.properites option field.
 7. Replace the blobstore.purger.frequency value with none.
 8. Click Save Changes.
 9. From the Actions menu, select Restart this Admin API Server.
 10. In the Restart this Admin API Server message, confirm your changes by clicking Restart this Admin API Server.
 11. When the Restart API Server step window displays Completed, click Close.

Troubleshooting

The Workload XM purge event does not delete directories and files that do not have the full wxm owner and file permissions. Files and directories may revert back to the hdfs owner when a restore is created from a snapshot. In this case and before creating an automatic or manual purge event you must verify the owner and file permissions of the required files to be purged.

To reset your HDFS files and directories as the wxm owner with full administrative permissions do the following:

1. In a terminal go to the /etc directory and open the hdfs password file by entering:


```
vim passwd
```
2. Search for the kafka parameter.
3. Replace /sbin/nologin with /bin/hash.
4. Save the file.
5. Grant full wxm access permissions to the hdfs password file by using the chown command.

Tracking your Purge Event from Log Entries

You can determine if the purge event was successful or identify potential problems from the Cloudera Manager Admin API Server log files.

The Admin API Server log file entries also list the names of the files and directories that were deleted and provide details about how many files and directories were deleted, the sum total size of the files and directories that were deleted, and the time they were deleted.