

End Lab 00:49:46

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more](#).

[Open Google Console](#)

Username

student-02-5a7887f0834



Password

6ho9kG8jtn9z



GCP Project ID

qwiklabs-gcp-00-af4419



Student Resources

[Flexible, Easy Data Pipelines on Google Cloud with Cloud Composer \(NEXT'18\)](#)

1 hour 30 minutes

Free



Overview

Workflows are a common theme in data analytics - they involve ingesting, transforming, and analyzing data to figure out the meaningful information within. In Google Cloud Platform (GCP), the tool for hosting workflows is Cloud Composer which is a hosted version of the popular open source workflow tool Apache Airflow.

In this lab, you use the GCP Console to set up a Cloud Composer environment. You then use Cloud Composer to go through a simple workflow that verifies the existence of a data file, creates a Cloud Dataproc cluster, runs an Apache Hadoop wordcount job on the Cloud Dataproc cluster, and deletes the Cloud Dataproc cluster afterwards.

What you'll do

- Use GCP Console to create the Cloud Composer environment
- View and run the DAG (Directed Acyclic Graph) in the Airflow web interface
- View the results of the wordcount job in storage.

Setup and requirements

Qwiklabs setup

For each lab, you get a new GCP project and set of resources for a fixed time at no cost.

1. Make sure you signed into Qwiklabs using an **incognito window**.
2. Note the lab's access time (for example, **02:00:00**) and make sure you can finish in that time block.

There is no pause feature. You can restart if needed, but you have to start at the beginning.

3. When ready, click **START LAB**.

4. Note your lab credentials. You will use them to sign in to Cloud Platform Console.

[Overview](#)

15/15

[Setup and requirements](#)

Ensure that the Kubernetes Engine API is successfully enabled

Ensure that the Cloud Composer API is successfully enabled

[Create Cloud Composer environment](#)[Airflow and core concepts](#)[Defining the workflow](#)[Viewing environment information](#)[Using the Airflow UI](#)[Setting Airflow variables](#)[Uploading the DAG to Cloud Storage](#)[Congratulations!](#)[Next steps](#)[End your lab](#)

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked.
[Learn more.](#)

[Open Google Console](#)

Username

Password

GCP Project ID

5. Click **Open Google Console**.

6. Click **Use another account** and copy/paste credentials for **this** lab into the prompts.

If you use other credentials, you'll get errors or incur charges.

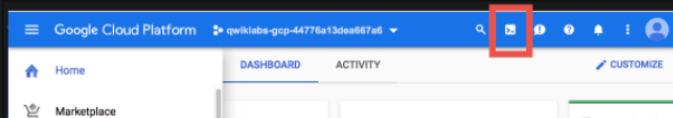
7. Accept the terms and skip the recovery resource page.

Do not click **End Lab** unless you are finished with the lab or want to restart it. This clears your work and removes the project.

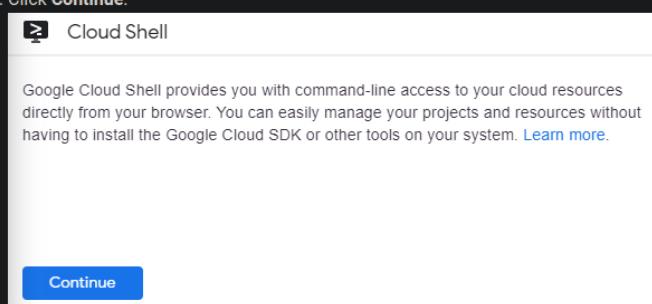
Activate Google Cloud Shell

Google Cloud Shell is a virtual machine that is loaded with development tools. It offers a persistent 5GB home directory and runs on the Google Cloud. Google Cloud Shell provides command-line access to your GCP resources.

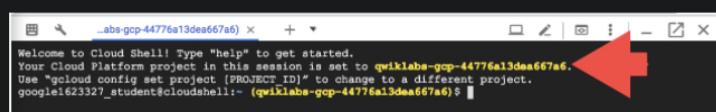
1. In GCP console, on the top right toolbar, click the Open Cloud Shell button.



2. Click **Continue**.



It takes a few moments to provision and connect to the environment. When you are connected, you are already authenticated, and the project is set to your **PROJECT_ID**. For example:



gcloud is the command-line tool for Google Cloud Platform. It comes pre-installed on Cloud Shell and supports tab-completion.

You can list the active account name with this command:

```
gcloud auth list
```



Output:

```
Credentialed accounts:  
- <myaccount>@<mydomain>.com (active)
```

Example output:

```
Credentialed accounts:  
- google1623327_student@qwiklabs.net
```

You can list the project ID with this command:

```
gcloud config list project
```



Output:

```
[core]  
project = <project_ID>
```

Example output:

```
[core]  
project = qwiklabs-gcp-44776a13dea667a6
```

Full documentation of **gcloud** is available on [Google Cloud gcloud Overview](#).

Check project permissions

Before you begin your work on Google Cloud, you need to ensure that your project has the correct permissions within Identity and Access Management (IAM).

1. In the Google Cloud console, on the **Navigation menu** (≡), click **IAM & Admin** > **IAM**.
2. Confirm that the default compute Service Account `<project-number>-compute@developer.gserviceaccount.com` is present and has the `editor` role assigned. The account prefix is the project number, which you can find on **Navigation menu** > **Home**.

Name	Role
Compute Engine default service account	Editor
Cloud Build Service Account	
Google APIs Service Agent	Editor
Owner	
App Engine Admin	
BigQuery Admin	

If the account is not present in IAM or does not have the `editor` role, follow the steps below to assign the required role.

- In the Google Cloud console, on the **Navigation menu**, click **Home**.
- Copy the project number (e.g. 729328892908).
- On the **Navigation menu**, click **IAM & Admin > IAM**.
- At the top of the **IAM** page, click **Add**.
- For **New principals**, type:

```
{project-number}-compute@developer.gserviceaccount.com
```



Replace `{project-number}` with your project number.

- For **Role**, select **Project (or Basic) > Editor**. Click **Save**.

Ensure that the Kubernetes Engine API is successfully enabled

To ensure access to the necessary APIs, restart the connection to the Kubernetes Engine API.

1. In the Google Cloud Console, enter **Kubernetes Engine API** in the top search bar. Click on the result for **Kubernetes Engine API**.
2. Click **Manage**.
3. Click **Disable API**.

If asked to confirm, click **Disable**.

Again, when prompted Do you want to disable Kubernetes Engine API and its dependent APIs? , Click **Disable**.

4. Click **Enable**.

When the API has been enabled again, the page will show the option to disable.

The screenshot shows the Google Cloud Platform API Services interface. The navigation bar at the top includes the Google Cloud logo, the text "Google Cloud Platform", and a dropdown menu showing "qwiklabs-gcp-04-0f110d8924dd". Below the navigation bar, there is a search bar and a filter icon. The main content area displays a table with one row. The row contains a blue hexagonal icon representing the Kubernetes Engine API, the text "Kubernetes Engine API", the word "Overview" in bold, and a blue "DISABLE API" button. The entire row is highlighted with a light gray background.

Ensure that the Cloud Composer API is successfully enabled

Restart the connection to the Cloud Composer API. In the prior step, restarting the Kubernetes Engine API forced the Cloud Composer API to be disabled.

1. In the Google Cloud Console, enter **Cloud Composer API** in the top search bar. Click on the result for **Cloud Composer API**.
2. Click **Enable**.

The screenshot shows the Google Cloud Platform API Services interface. The navigation bar at the top includes the Google Cloud logo, the text "Google Cloud Platform", and a dropdown menu showing "qwiklabs-gcp-04-0f110d8924dd". Below the navigation bar, there is a search bar and a filter icon. The main content area displays a table with one row. The row contains a blue hexagonal icon representing the Cloud Composer API, the text "Cloud Composer API", the word "Overview" in bold, and a blue "DISABLE API" button. The entire row is highlighted with a light gray background.

Google

Manages Apache Airflow environments on Google Cloud Platform.

ENABLE TRY THIS API

When the API has been enabled again, the page will show the option to disable.

Google Cloud Platform

APIs & Services

Cloud Composer API

Overview

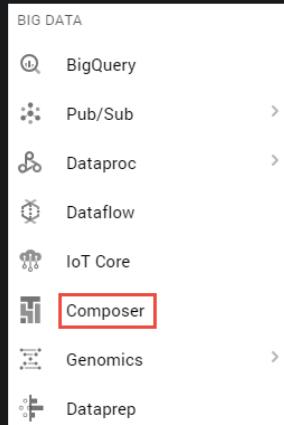
DISABLE API

Create Cloud Composer environment

In this section, you create a Cloud Composer environment.

Before proceeding further, make sure that you have performed earlier tasks to ensure that the required APIs are successfully enabled. If not, then please perform those tasks otherwise Cloud Composer environment creation will fail.

1. Go to **Navigation menu > Composer**:



2. Click **CREATE ENVIRONMENT** and select **Composer 1**. Set the following for your environment:

Property	Value
Name	highcpu
Location	us-central1
Zone	us-central1-a
Machine type	n1-highcpu-4

Leave all other settings as default.

3. Click **Create**.

The environment creation process is completed when the green checkmark displays to the left of the environment name on the Environments page in the GCP Console.

It can take 10-20 minutes for the environment to complete the setup process. Continue with the lab while the environment spins up.

Click **Check my progress** to verify the objective.



Create Cloud Composer environment.

[Check my progress](#)

Assessment Completed!

Create a Cloud Storage bucket

Create a Cloud Storage bucket in your project. This buckets will be used as output for the Hadoop job from Dataproc.

1. Go to **Navigation menu > Cloud Storage > Browser** and then click **Create bucket**.

2. Give your bucket a universally unique name, then click **Create**.

Remember the Cloud Storage bucket name as you'll use it as an Airflow variable later in the lab.

Click **Check my progress** to verify the objective.



Create a Cloud Storage bucket.

[Check my progress](#)

Assessment Completed!

Airflow and core concepts

While waiting for your Composer environment to get created, review some terms that are used with Airflow.

[Airflow](#) is a platform to programmatically author, schedule and monitor workflows.

Use Airflow to author workflows as directed acyclic graphs (DAGs) of tasks. The airflow scheduler executes your tasks on an array of workers while following the specified dependencies.

Core concepts

DAG

A Directed Acyclic Graph is a collection of all the tasks you want to run, organized in a way that reflects their relationships and dependencies.

Operator

The description of a single task, it is usually atomic. For example, the *BashOperator* is used to execute bash command.

Task

A parameterised instance of an Operator; a node in the DAG.

Task Instance

A specific run of a task; characterized as: a DAG, a Task, and a point in time. It has an indicative state: *running, success, failed, skipped, ...*

You can read more about the concepts [here](#).

Defining the workflow

Now let's discuss the workflow you'll be using. Cloud Composer workflows are comprised of [DAGs \(Directed Acyclic Graphs\)](#). DAGs are defined in standard Python files that are placed in Airflow's `DAG_FOLDER`. Airflow will execute the code in each file to dynamically build the DAG objects. You can have as many DAGs as you want, each describing an arbitrary number of tasks. In general, each one should correspond to a single logical workflow.

Below is the `hadoop_tutorial.py` workflow code, also referred to as the DAG:

```
"""Example Airflow DAG that creates a Cloud Dataproc cluster, runs the
Hadoop
wordcount example, and deletes the cluster.
This DAG relies on three Airflow variables
https://airflow.apache.org/concepts.html#variables
* gcp_project - Google Cloud Project to use for the Cloud Dataproc
cluster.
* gce_zone - Google Compute Engine zone where Cloud Dataproc cluster
should be
    created.
* gcs_bucket - Google Cloud Storage bucket to used as output for the
Hadoop jobs from Dataproc.
    See https://cloud.google.com/storage/docs/creating-buckets for
creating a
    bucket.
"""

import datetime
import os
from airflow import models
from airflow.contrib.operators import dataproc_operator
from airflow.utils import trigger_rule
# Output file for Cloud Dataproc job.
output_file = os.path.join(
    models.Variable.get('gcs_bucket'), 'wordcount',
    datetime.datetime.now().strftime('%Y%m%d-%H%M%S')) + os.sep
# Path to Hadoop wordcount example available on every Dataproc cluster.
WORDCOUNT_JAR = (
    'file:///usr/lib/hadoop-mapreduce/hadoop-mapreduce-examples.jar'
)
# Arguments to pass to Cloud Dataproc job.
wordcount_args = ['wordcount', 'gs://pub/shakespeare/rose.txt',
output_file]
yesterday = datetime.datetime.combine(
    datetime.datetime.today() - datetime.timedelta(1),
    datetime.datetime.min.time())
default_dag_args = {
    # Setting start date as yesterday starts the DAG immediately when it
is
    # detected in the Cloud Storage bucket.
    'start_date': yesterday,
    # To email on failure or retry set 'email' arg to your email and
enable
    # emailing here.
    'email_on_failure': False,
    'email_on_retry': False,
    # If a task fails, retry it once after waiting at least 5 minutes
    'retries': 1,
    'retry_delay': datetime.timedelta(minutes=5),
    'project_id': models.Variable.get('gcp_project')
}
with models.DAG(
    'composer_sample_quickstart',
    # Continue to run DAG once per day
    schedule_interval=datetime.timedelta(days=1),
    default_args=default_dag_args) as dag:
    # Create a Cloud Dataproc cluster.
    create_dataproc_cluster =
```

```

create_dataproc_cluster =
    dataproc_operator.DataproClusterCreateOperator(
        task_id='create_dataproc_cluster',
        # Give the cluster a unique name by appending the date
        scheduled.
        # See https://airflow.apache.org/code.html#default-variables
        cluster_name='composer-hadoop-tutorial-cluster-{{ ds_nodash }}',
        num_workers=2,
        region='us-central1',
        zone=models.Variable.get('gce_zone'),
        image_version='2.0',
        master_machine_type='n1-standard-2',
        worker_machine_type='n1-standard-2')
    # Run the Hadoop wordcount example installed on the Cloud Dataproc
    cluster
    # master node.
    run_dataproc_hadoop = dataproc_operator.DataProcHadoopOperator(
        task_id='run_dataproc_hadoop',
        region='us-central1',
        main_jar=WORDCOUNT_JAR,
        cluster_name='composer-hadoop-tutorial-cluster-{{ ds_nodash }}',
        arguments=wordcount_args)
    # Delete Cloud Dataproc cluster.
    delete_dataproc_cluster =
        dataproc_operator.DataproClusterDeleteOperator(
            task_id='delete_dataproc_cluster',
            region='us-central1',
            cluster_name='composer-hadoop-tutorial-cluster-{{ ds_nodash }}',
            # Setting trigger_rule to ALL_DONE causes the cluster to be
            deleted
            # even if the Dataproc job fails.
            trigger_rule=trigger_rule.TriggerRule.ALL_DONE)
    # Define DAG dependencies.
    create_dataproc_cluster >> run_dataproc_hadoop >>
    delete_dataproc_cluster

```

To orchestrate the three workflow tasks, the DAG imports the following operators:

1. `DataproClusterCreateOperator`: Creates a Cloud Dataproc cluster.
2. `DataProcHadoopOperator`: Submits a Hadoop wordcount job and writes results to a Cloud Storage bucket.
3. `DataproClusterDeleteOperator`: Deletes the cluster to avoid incurring ongoing Compute Engine charges.

The tasks run sequentially, which you can see in this section of the file:

```

# Define DAG dependencies.
create_dataproc_cluster >> run_dataproc_hadoop >>
delete_dataproc_cluster

```

The name of the DAG is `quickstart`, and the DAG runs once each day.

```

with models.DAG(
    'composer_sample_quickstart',
    # Continue to run DAG once per day
    schedule_interval=datetime.timedelta(days=1),
    default_args=default_dag_args) as dag:

```

Because the `start_date` that is passed into `default_dag_args` is set to `yesterday`, Cloud Composer schedules the workflow to start immediately after the DAG uploads.

Viewing environment information

1. Go back to **Composer** to check the status of your environment.
2. Once your environment has been created, click the name of the environment (highcpu) to see its details.

On the **Environment details** you'll see information such as the Airflow web interface URL, Kubernetes Engine cluster ID, and a link to the DAGs folder, which is stored in your bucket.

Note: Cloud Composer only schedules the workflows in the `/dags` folder.

Using the Airflow UI

To access the Airflow web interface using the GCP Console:

1. Go back to the **Environments** page.
2. In the **Airflow webserver** column for the environment, click **Airflow**.
3. Click on your lab credentials.
4. The Airflow web interface opens in a new browser window.

Setting Airflow variables

Airflow variables are an Airflow-specific concept that is distinct from [environment variables](#).

1. Select **Admin > Variables** from the Airflow menu bar, then **Create**.

	Key	Val	Is Encrypted
There are no items in the table.			

2. Create the following Airflow variables, `gcp_project`, `gcs_bucket`, and `gce_zone`:

KEY	VALUE	Details
<code>gcp_project</code>	<your project-id>	The Google Cloud Platform project you're using for this quickstart.
<code>gcs_bucket</code>	<code>gs://<my-bucket></code>	Replace <code><my-bucket></code> with the name of the Cloud Storage bucket you made earlier. This bucket stores the output from the Hadoop jobs from Dataproc.
<code>gce_zone</code>	<code>us-central1-a</code>	This is the Compute Engine zone where your Cloud Dataproc cluster will be created. To chose a different zone, see Available regions & zones .

Click **Save and Add Another** after adding first two variable and click **Save** for the third variable. Your Variables table should look like this when you're finished:

Variables

Choose File		No file chosen	Import Variables		
List (3)		Create	Add Filter	With selected	Search
		Key	Val		
		gcp_project	project-id		
		gcs_bucket	gs://bucket-name		
		gce_zone	zone		

Uploading the DAG to Cloud Storage

To upload the DAG:

1. In Cloud Shell, upload a copy of the `hadoopTutorial.py` file to the Cloud Storage bucket that was automatically created when you created the environment.

Replace `<dag_folder_path>` in the following command:

```
gsutil cp gs://cloud-
training/datawarehousing/lab_assets/hadoopTutorial.py
<dag_folder_path>
```

with the path to the DAGs folder. You can get the path by going to **Composer**. Click on the environment you created earlier and then click on the **Environment Configuration** tab to see the details of the environment. Find `DAGs folder` and copy the path.

The screenshot shows the 'Environment Configuration' section of the Composer interface. It includes fields for Network access control (All IP addresses have access), Machine type (composer-n1-webserver-2 (2 vCPU, 1.6 GB memory)), and DAGs folder (highlighted with a red box) which contains the value `gs://us-central1-highcpu-0682d8c0-bucket/dags`. Below it is the Airflow web UI URL: `https://00ehf683ea37daflytp.appspot.com`.

The revised command to upload the file will look similar to the one below:

```
gsutil cp gs://cloud-
training/datawarehousing/lab_assets/hadoopTutorial.py gs://us-central1-
highcpu-0682d8c0-bucket/dags
```

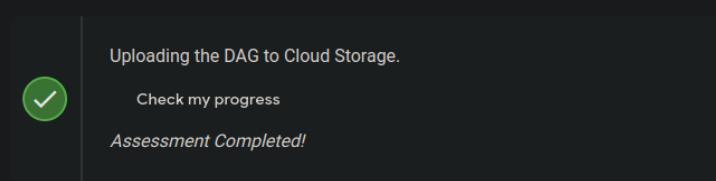
Once the file has been successfully uploaded to the DAGs directory, open `dags` folder in the bucket and you will see the file in the **Objects** tab of the Bucket details.

The screenshot shows the 'Bucket details' page for the bucket `us-central1-highcpu-0682d8c0-bucket`. The left sidebar has 'Storage' selected, with 'Browser' highlighted. The main area shows the 'OBJECTS' tab with two files listed: `airflowMonitoring.py` and `hadoopTutorial.py`, both circled with a red box. There are also tabs for 'CONFIGURATION' and 'PERMISSIONS'.

When a DAG file is added to the DAGs folder, Cloud Composer adds the DAG to Airflow and schedules it automatically. DAG changes occur within 3-5 minutes.

You can see the task status of the `composer_hadoop_tutorial` DAG in the Airflow web interface.

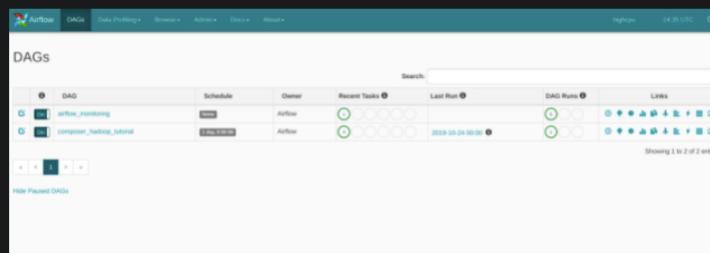
Click **Check my progress** to verify the objective.



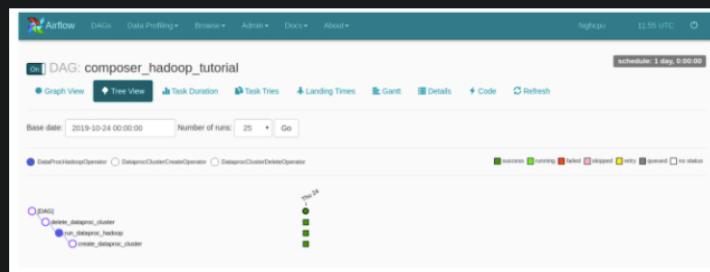
Exploring DAG runs

When you upload your DAG file to the `dags` folder in Cloud Storage, Cloud Composer parses the file. If no errors are found, the name of the workflow appears in the DAG listing, and the workflow is queued to run immediately.

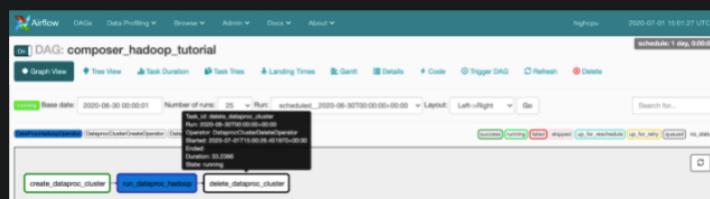
Make sure that you're on the DAGs tab in the Airflow web interface. It takes several minutes for this process to complete. Refresh your browser to make sure you're looking at the latest information.



1. In Airflow, click **composer_hadoop_tutorial** to open the DAG details page. This page includes several representations of the workflow tasks and dependencies.



2. In the toolbar, click **Graph View**. Mouseover the graphic for each task to see its status. Note that the border around each task also indicates the status (green border = running; red = failed, etc.).



3. Click the "Refresh" link to make sure you're looking at the most recent information. The borders of the processes change colors as the state of the process changes

Note: If your Dataproc cluster already exists, you can run the workflow again to reach the success state by clicking `create_dataproc_cluster` graphic and then click **Clear** to reset the three tasks and click **OK** to confirm.

4. Once the status for `create_dataproc_cluster` has changed to "running", go to **Navigation menu > Dataproc**, then click on:

- **Clusters** to monitor cluster creation and deletion. The cluster created by the workflow is ephemeral; it only exists for the duration of the workflow and is deleted as part of

is ephemeral. It only exists for the duration of the workflow and is deleted as part of the last workflow task.

- **Jobs** to monitor the Apache Hadoop wordcount job. Click the Job ID to see job log output.
5. Once Dataproc gets to a state of "Running", return to Airflow and click **Refresh** to see that the cluster is complete.

When the `run_dataproc_hadoop` process is complete, go to **Navigation menu > Cloud Storage > Browser** and click on the name of your bucket to see the results of the wordcount in the `wordcount` folder.

6. Once all the steps are complete in the DAG, each step has a dark green border.
Additionally the Dataproc cluster that was created is now deleted.

Congratulations!

You've successfully run a Cloud Composer workflow!

Next steps

- Check out when Cloud Composer was presented at NEXT 18 in San Francisco:
<https://www.youtube.com/watch?v=GeNFEtt-D4k>
- To see the value of a variable, run the Airflow CLI sub-command `variables` with the `get` argument or use the [Airflow web interface](#).
- For information about the Airflow web interface, see [Accessing the web interface](#).

End your lab

When you have completed your lab, click **End Lab**. Qwiklabs removes the resources you've used and cleans the account for you.

You will be given an opportunity to rate the lab experience. Select the applicable number of stars, type a comment, and then click **Submit**.

The number of stars indicates the following:

- 1 star = Very dissatisfied
- 2 stars = Dissatisfied
- 3 stars = Neutral
- 4 stars = Satisfied
- 5 stars = Very satisfied

You can close the dialog box if you don't want to provide feedback.

For feedback, suggestions, or corrections, please use the **Support** tab.

Copyright 2021 Google LLC All rights reserved. Google and the Google logo are trademarks of Google LLC. All other company and product names may be trademarks of the respective companies with which they are associated.