

Starting Hadoop

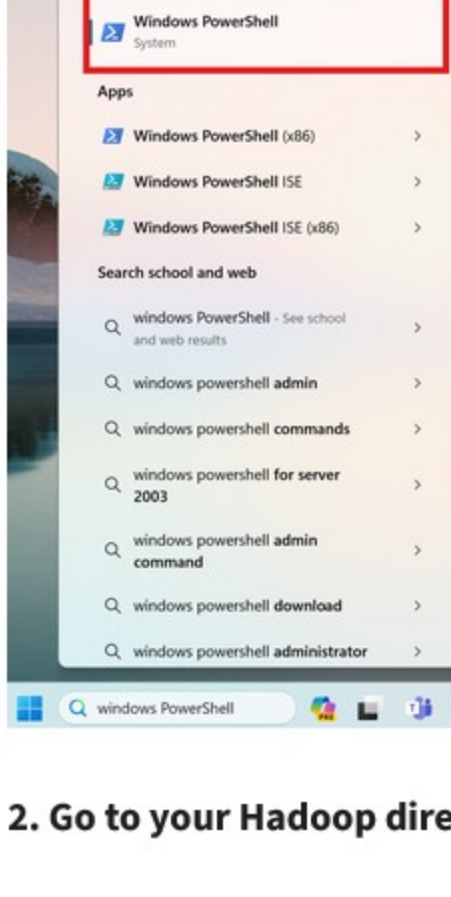
Learning Goals

By the end of this activity, you will be able to:

- Start Hadoop using Docker containers
- Copy files into and out of the Hadoop Distributed File System (HDFS).

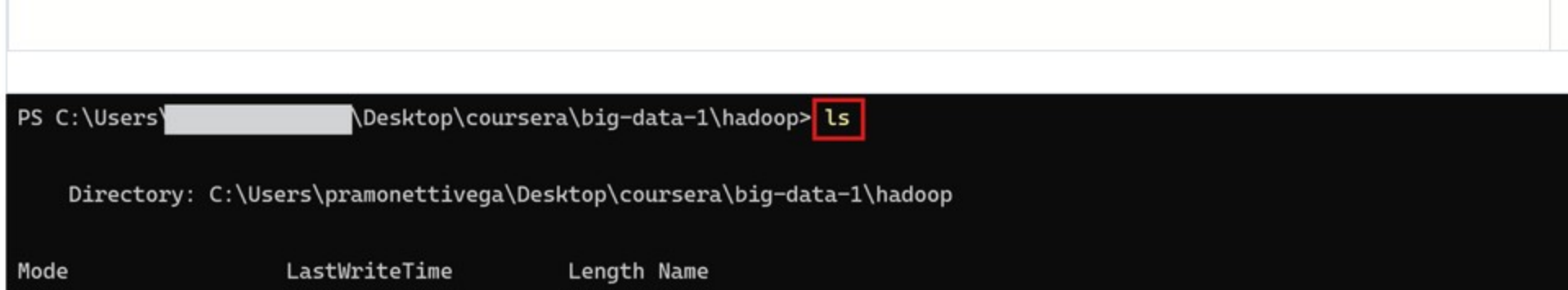
Instructions

1. **Open a terminal shell.** Open your local system terminal shell. In Windows, you can use Windows PowerShell, while in Mac, you can use Terminal.



2. **Go to your Hadoop directory.** In the terminal shell, change to your `big-data-1/hadoop` directory

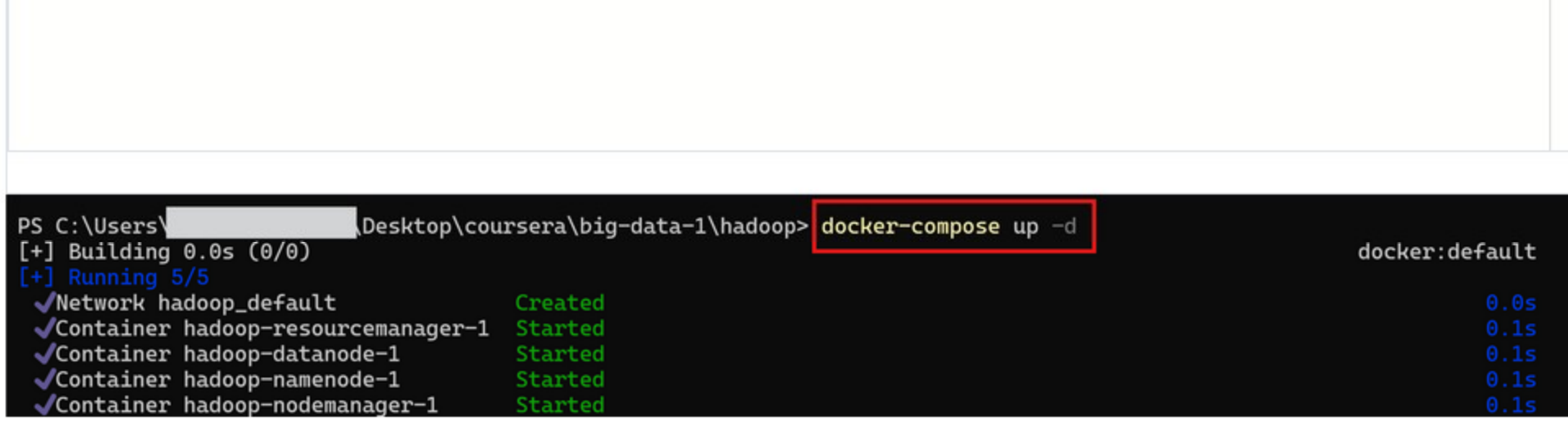
Type `ls` to see the files within the directory.



3. **Start Docker.** In order to start working, we need to start the Docker Engine. The easiest way to do this is to simply open Docker Desktop. Wait for it to initiate.

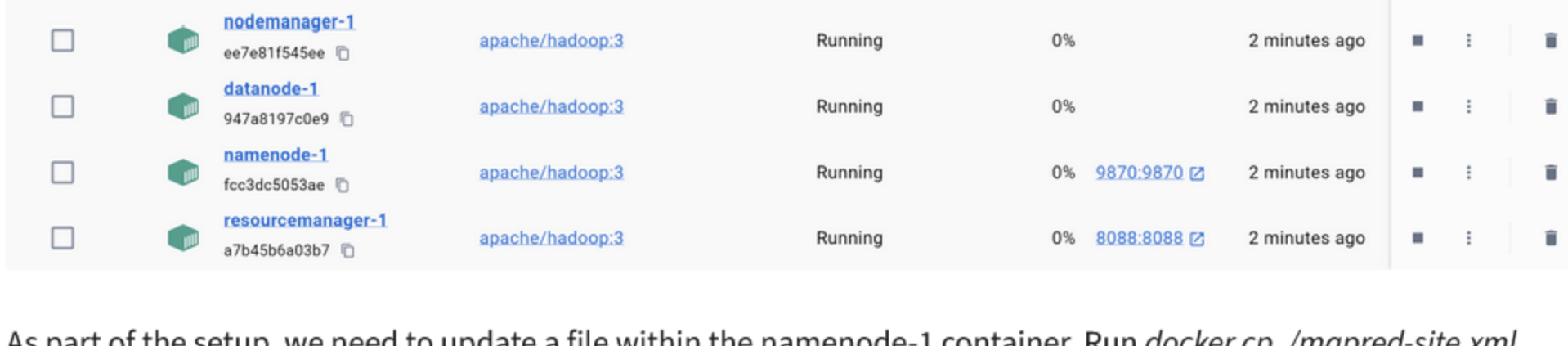
4. **Build your Hadoop cluster.** To run a Hadoop cluster locally, we are going to run multiple Docker containers, each of them running a different service. The `docker-compose.yaml` file within the directory is the recipe to achieve this.

Go back to your terminal shell and run the following command:

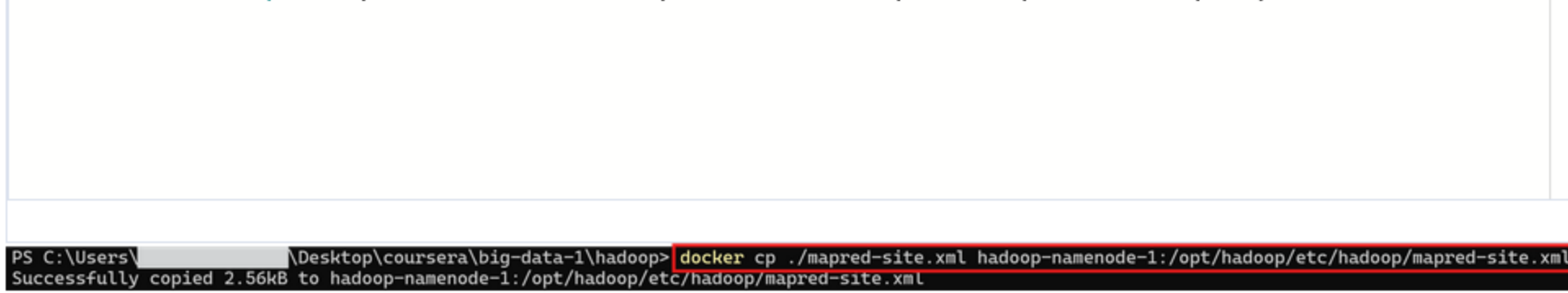


This will start the creation of the multiple containers needed for this activity. This process might take a few minutes.

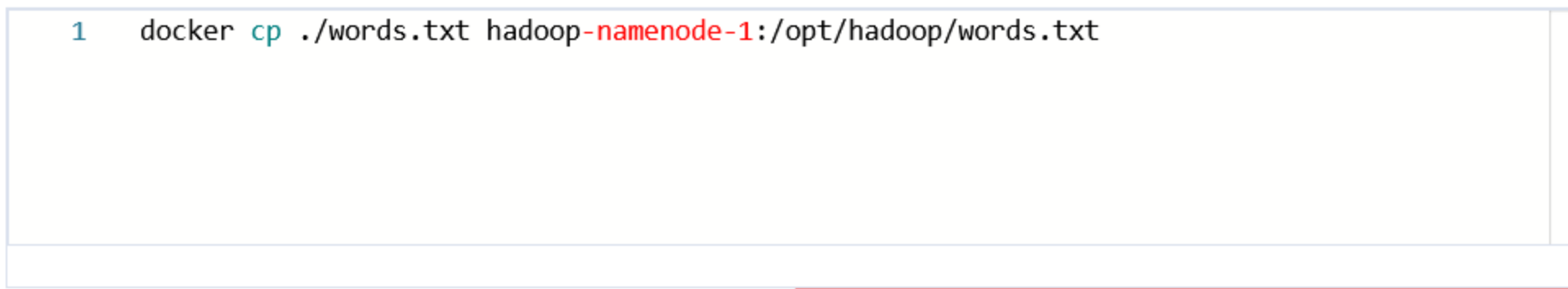
Go to Docker Desktop, click on *Containers* in the top left, and make sure the containers are running.



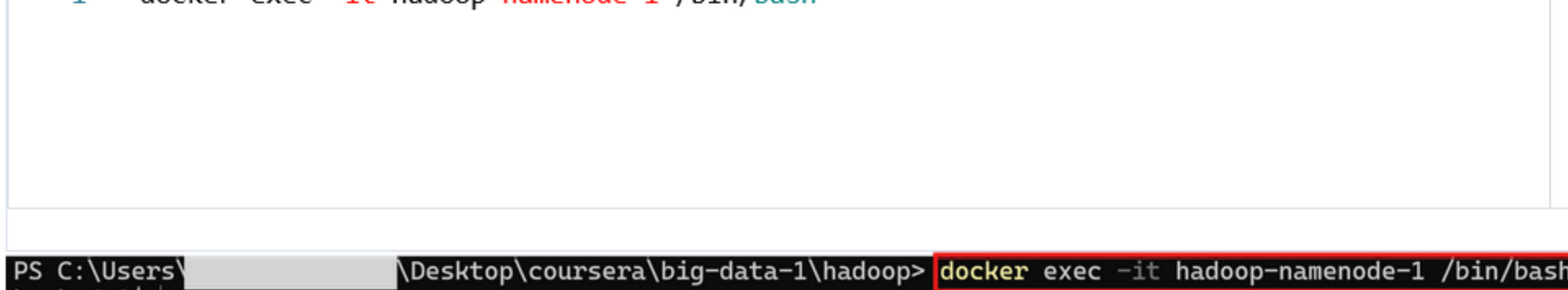
As part of the setup, we need to update a file within the namenode-1 container. Run `docker cp ./mapred-site.xml hadoop-namenode-1:/opt/hadoop/etc/hadoop/mapred-site.xml` to replace the `mapred-site.xml` file within the container.



We also need to transfer the `words.txt` file (which is a collection of Williams Shakespeare's works) to the container. Run `docker cp ./words.txt hadoop-namenode-1:/opt/hadoop/words.txt` to copy the file from your local directory to the container.



5. **Access containers terminal shell.** Run `docker exec -it hadoop-namenode-1 /bin/bash` to access namenode-1 shell.

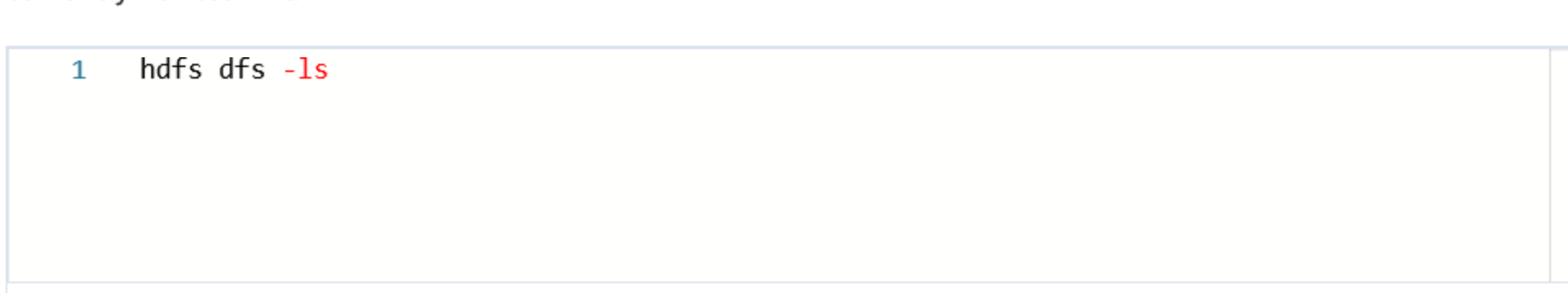


Now we are inside our local Hadoop cluster.

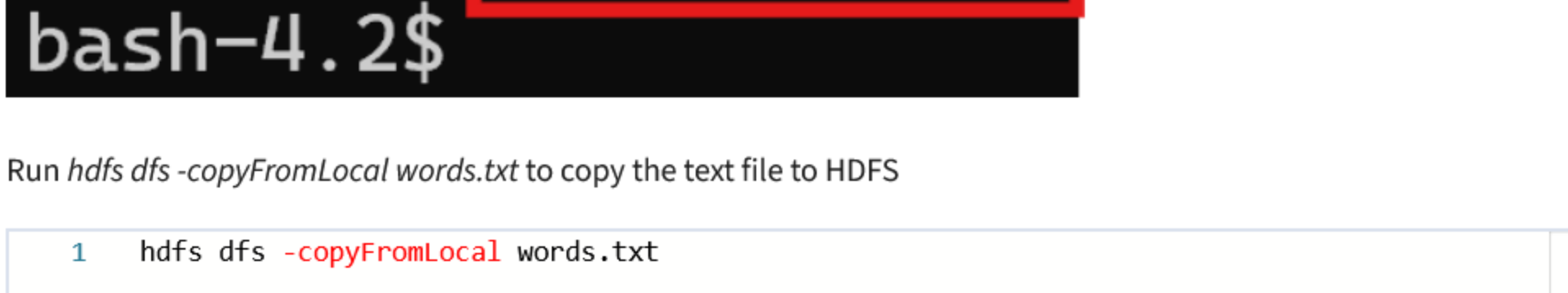
6. **Creating the HDFS root folder.** To start interacting with the HDFS, we first need to create a user folder to be able to place our files in it. Run `hdfs dfs -mkdir -p /user/hadoop`



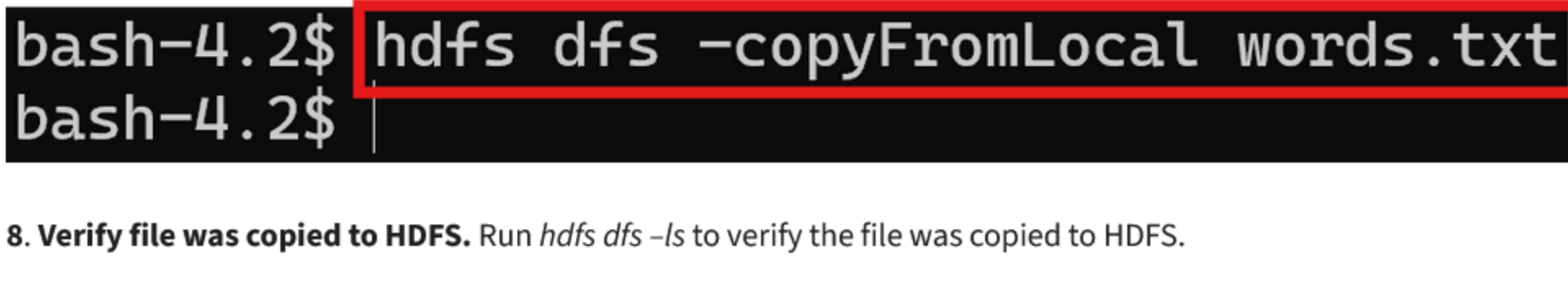
7. **Copy file to HDFS.** Run `hdfs dfs -ls` to see the files within the file system. Nothing will come out, as there are currently no files in it.



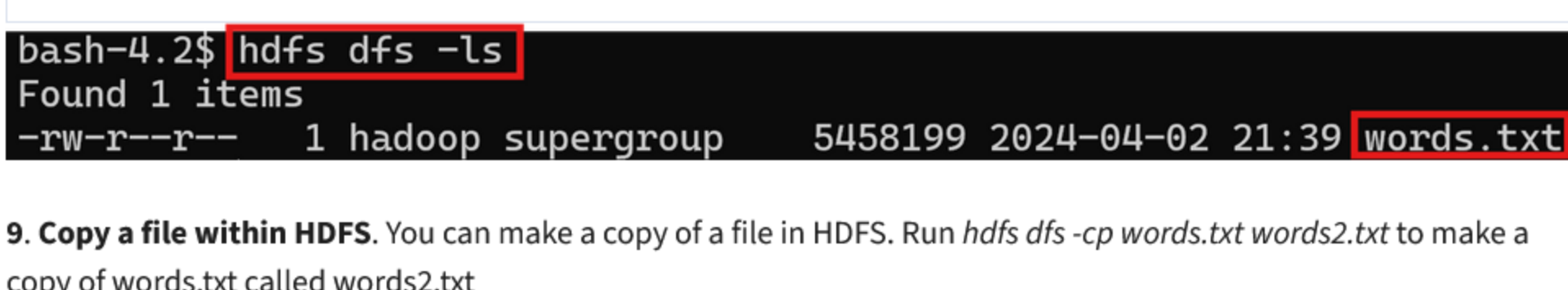
Run `hdfs dfs -copyFromLocal words.txt` to copy the text file to HDFS



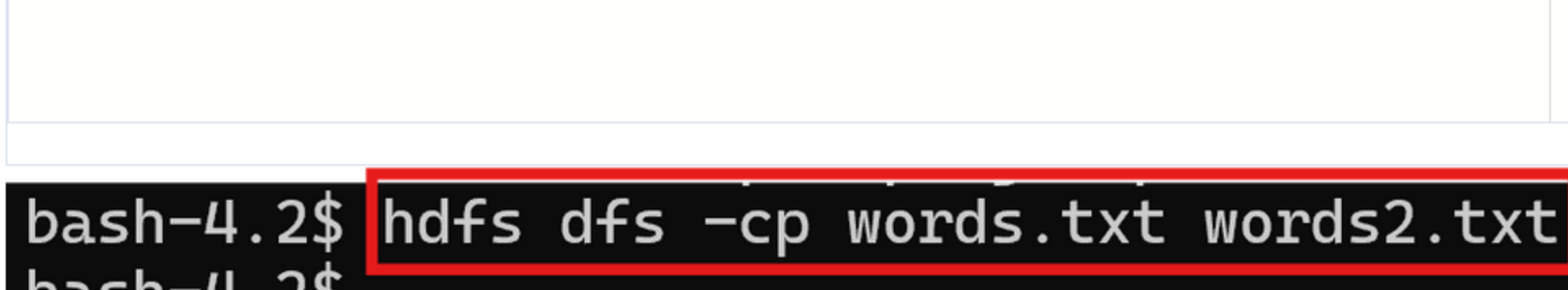
8. **Verify file was copied to HDFS.** Run `hdfs dfs -ls` to verify the file was copied to HDFS.



9. **Copy a file within HDFS.** You can make a copy of a file in HDFS. Run `hdfs dfs -cp words.txt words2.txt` to make a copy of `words.txt` called `words2.txt`



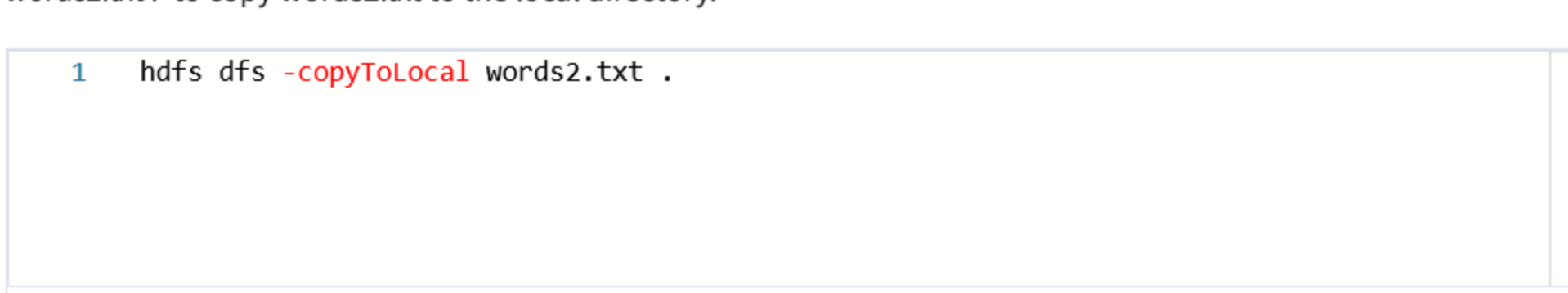
We can see the both files by running `hdfs dfs -ls`



10. **Copy a file from HDFS.** We can also copy a file from HDFS to the local file system. Run `hdfs dfs -copyToLocal words2.txt .` to copy `words2.txt` to the local directory.



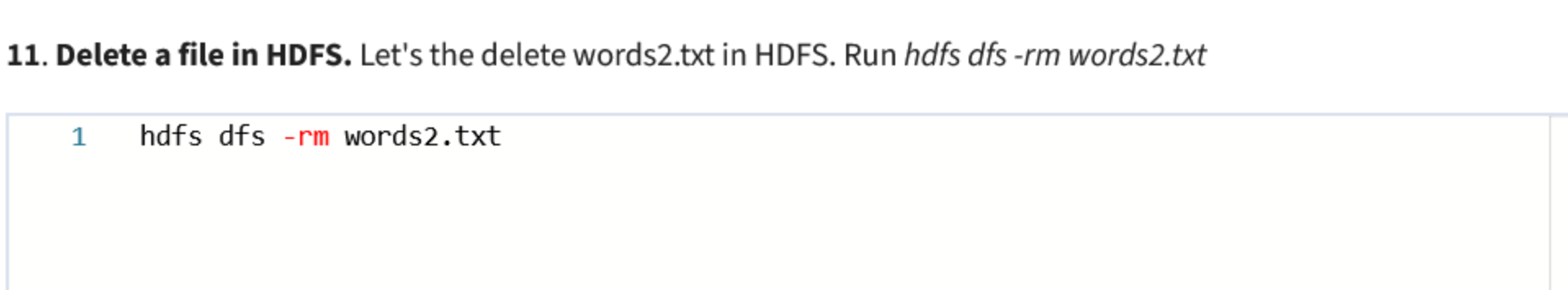
Let's run `ls` to see that `words2.txt` is there.



11. **Delete a file in HDFS.** Let's delete `words2.txt` in HDFS. Run `hdfs dfs -rm words2.txt`

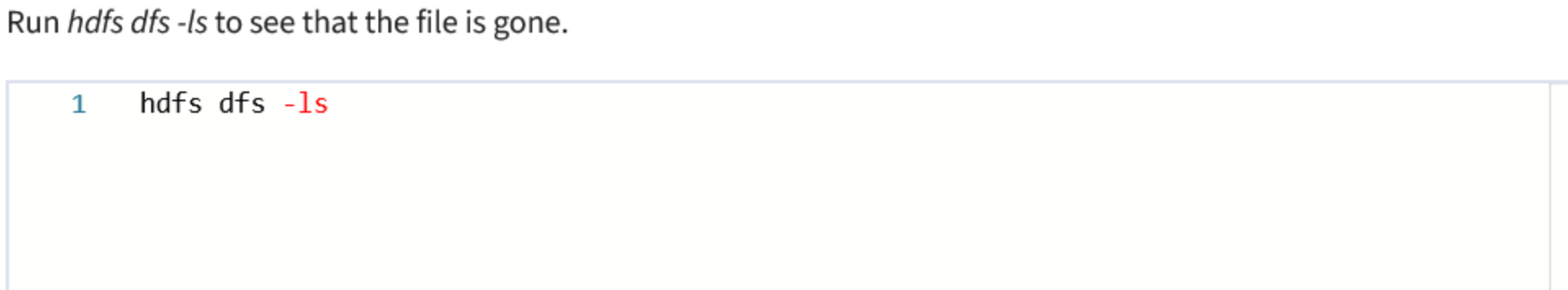


Run `hdfs dfs -ls` to see that the file is gone.



12. **Deleting cluster.** From here, you can continue to the next activity. If you're planning to stop working and come back to the course later, make sure to delete the cluster first and repeat steps 3-5 next time you come back. Given the current configuration, if you simply stop and restart the cluster from Docker Desktop, the cluster won't work properly.

Run `exit` to exit the container's shell. You will be sent back to your local terminal shell.



To delete the containers, run `docker compose down`.

