

NTI Graduation Project

At the beginning, I generated logs using a Python script and a pre-existing data file.

Data Content:

```
log-generator / defaultDataFile.txt
Code Blame 8153 lines (5861 loc) · 188 KB
21
22     A long time ago, in a galaxy far, far, away...
23
24     A vast sea of stars serves as the backdrop for the main title.
25     War drums echo through the heavens as a rollup slowly crawls
26     into infinity.
27
28     It is a period of civil war. Rebel spaceships,
29     striking from a hidden base, have won their first
30     victory against the evil Galactic Empire.
31
32     During the battle, Rebel spies managed to steal
33     secret plans to the Empire's ultimate weapon, the
34     Death Star, an armored space station with enough
35     power to destroy an entire planet.
36
37     Pursued by the Empire's sinister agents, Princess
38     Leia races home aboard her starship, custodian of
39     the stolen plans that can save her people and
40     restore freedom to the galaxy...
41
42     The awesome yellow planet of Tatooine emerges from a total
43     eclipse, her two moons glowing against the darkness. A tiny
44     silver spacecraft, a Rebel Blockade Runner firing lasers
45     from the back of the ship, races through space. It is pursued
46     by a giant Imperial Stardestroyer. Hundreds of deadly
47     laserbolts streak from the Imperial Stardestroyer, causing
48     the main solar fin of the Rebel craft to disintegrate.
49
```

Script:

```
README
log-generator

Simple python script to generate random time series data into a log file (GMT time) useful for generating sample
time-series based data to feed other downstream systems.

Usage

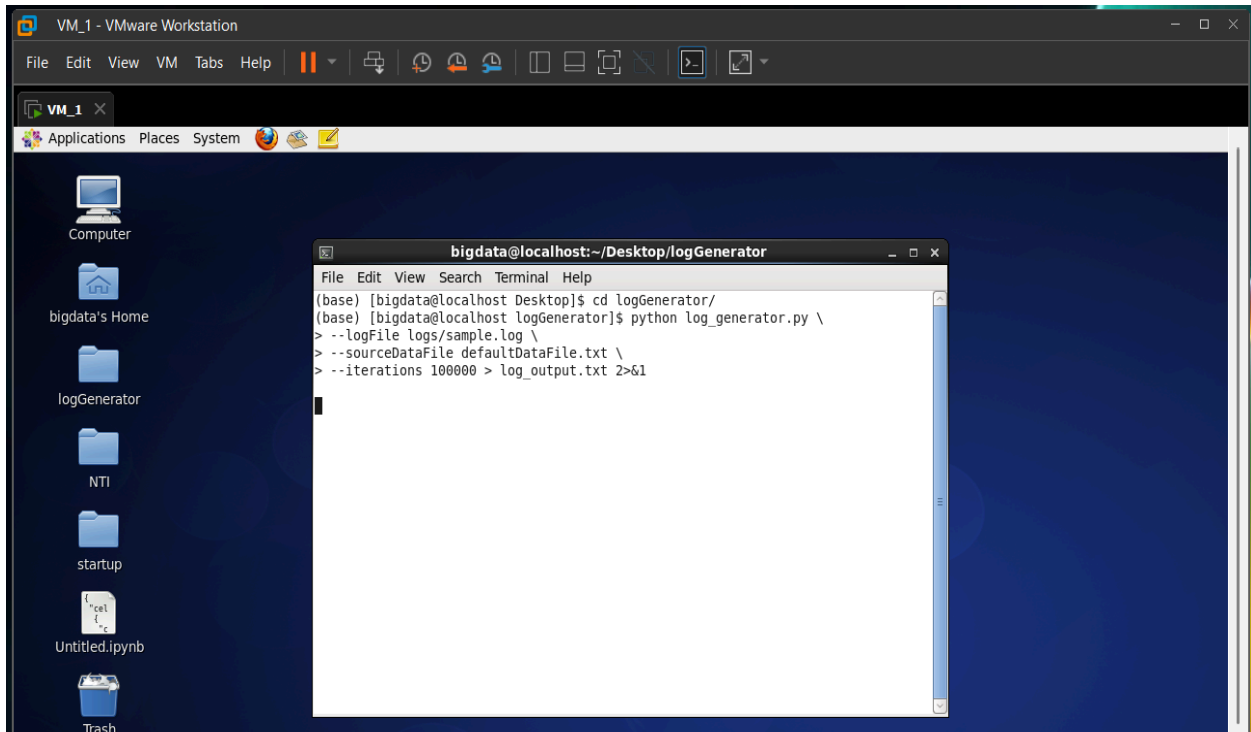
logGenerator.py --logFile <targetFile>
                [--minSleepMs <int>] [--maxSleepMs <int>]
                [--sourceDataFile <fileWithTextData>] [--iterations <long>]
                [--minLines <int>] [--maxLines <int>]
                [--logPattern <pattern>] [--datePattern <pattern>]

Defaults

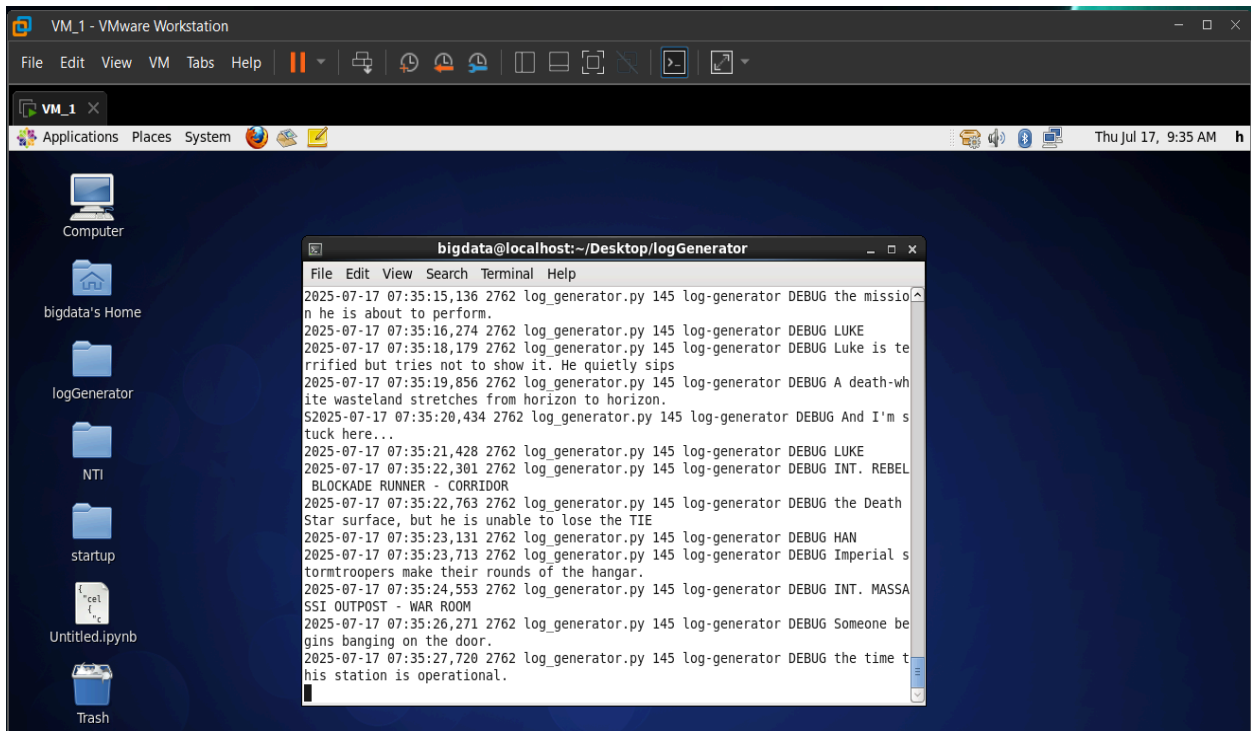
iterations = -1 # infinite
minSleep = 0.1
maxSleep = 1
minLines = 1
maxLines = 1
logFile = 'logGenerator.log'
sourceDataFile = 'defaultDataFile.txt'
logPattern = '%(asctime)s,%(msecs)d %(process)d %(filename)s %(lineno)d %(name)s %(levelname)s %(message'
datePattern = "%Y-%m-%d %H:%M:%S"
```

At the beginning, I ran the script using the following command:

Note: the log generator folder on the Desktop.



After that, I streamed the log data live using: [`tail -f logs/sample.log`]



After I started HDFS and created a directory called 'graduation', the logs were saved there using Flume.

```
hdfs dfs -mkdir -p /user/bigdata/graduation
```

VM_1 - VMware Workstation

File Edit View VM Tabs Help

VM_1

Applications Places System

Thu Jul 17, 9:37 AM

Computer

bigdata's Home

logGenerator

NTI

startup

Untitled.ipynb

Trash

flumeHdfs.conf

Browsing HDFS - Mozilla Firefox

File Edit View History Bookmarks Tools Help

Browsing HDFS

localhost:50070/explorer.html#/user/bigdata/graduation

Hadoop Overview Datanodes Snapshot Startup Progress Utilities

Browse Directory

/user/bigdata/graduation Go!

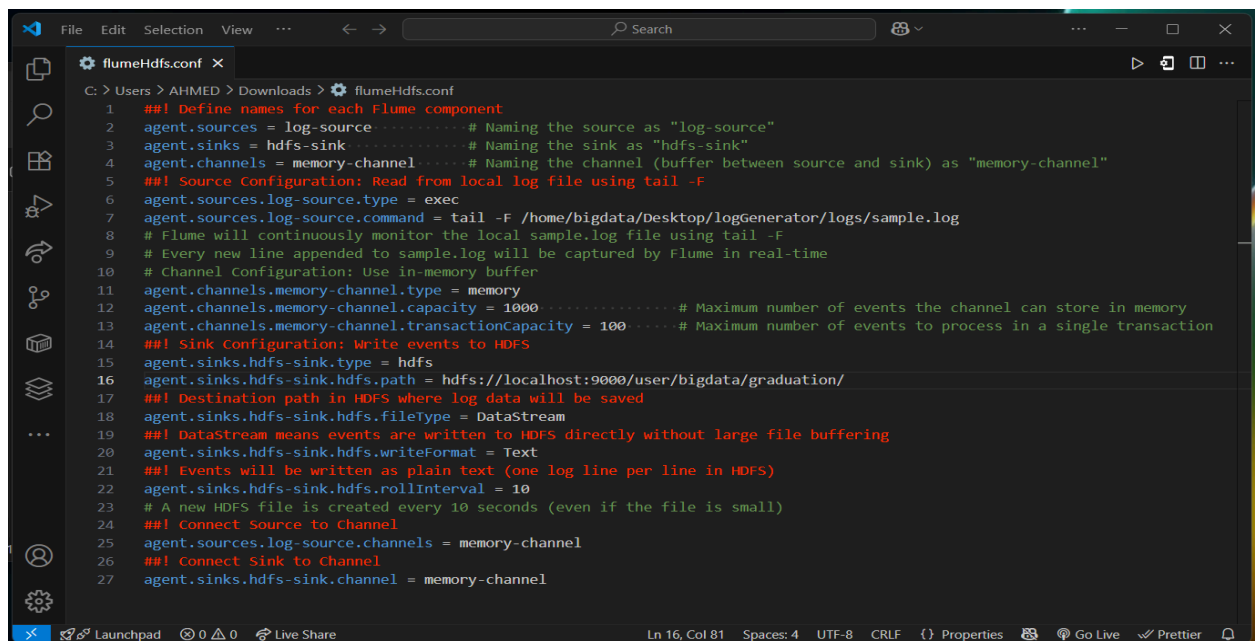
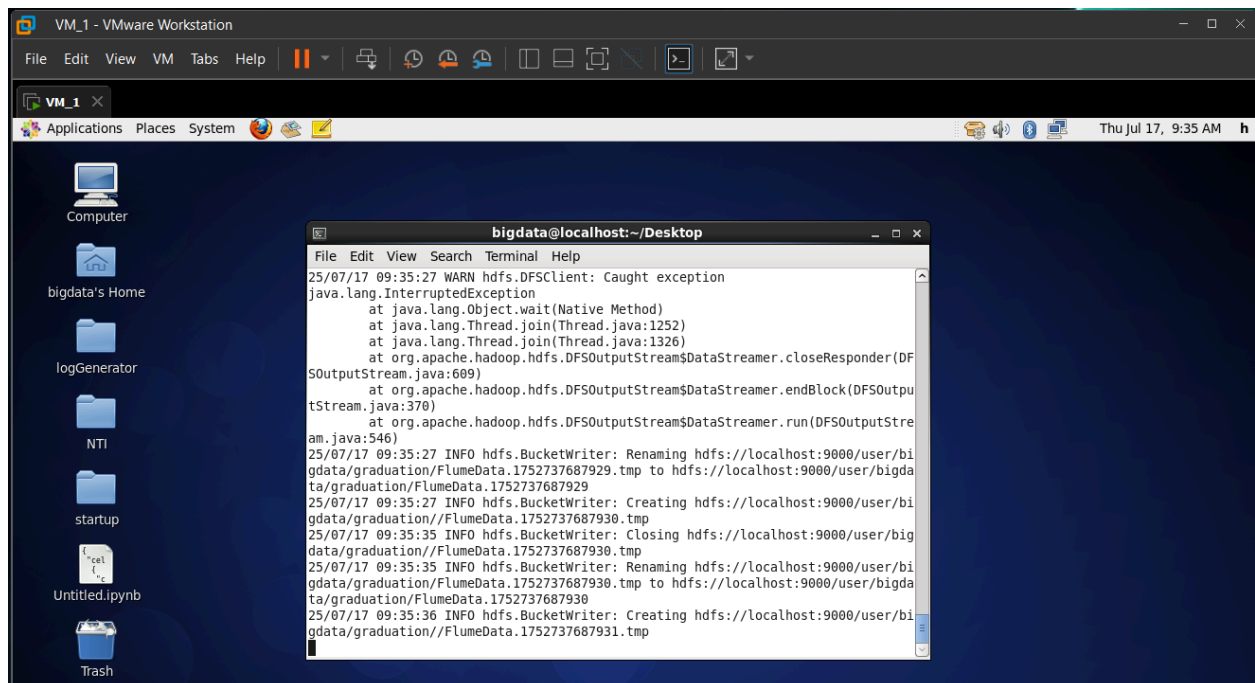
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	bigdata	supergroup	1.04 KB	Thu 17 Jul 2025 09:34:22 AM EET	1	128 MB	FlumeData.1752737655562
-rw-r--r--	bigdata	supergroup	982 B	Thu 17 Jul 2025 09:34:22 AM EET	1	128 MB	FlumeData.1752737655563
-rw-r--r--	bigdata	supergroup	960 B	Thu 17 Jul 2025 09:34:29 AM EET	1	128 MB	FlumeData.1752737655564

[bigdata@localhost:~/...]

To direct input to this VM, move the mouse pointer inside or press Ctrl+G.

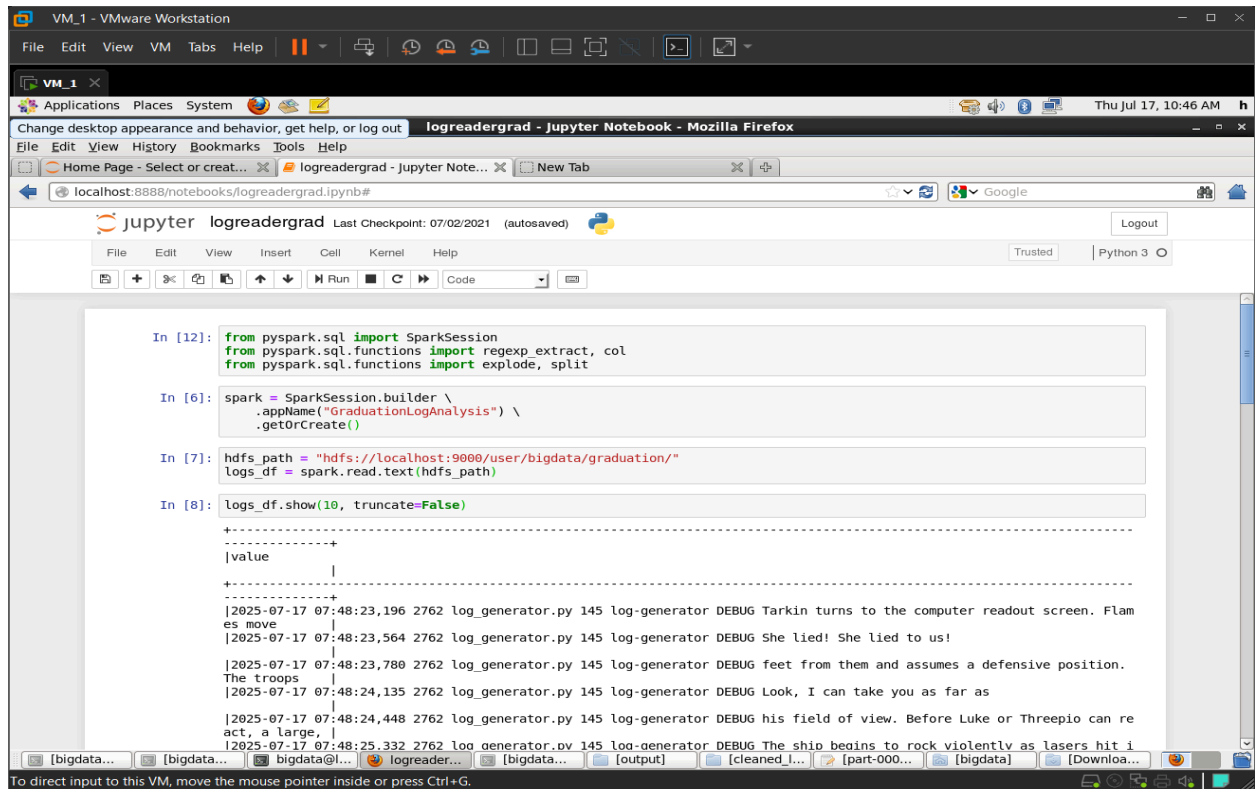
Then, I ran Flume with the following command after creating the configuration file:

```
flume-ng agent \
--conf conf \
--conf-file /home/bigdata/Desktop/flumeHdfs.conf \
--name agent \
-Dflume.root.logger=INFO,console
```



After that, I launched the Jupyter-Notebook and initialized PySpark.

[jupyter notebook]



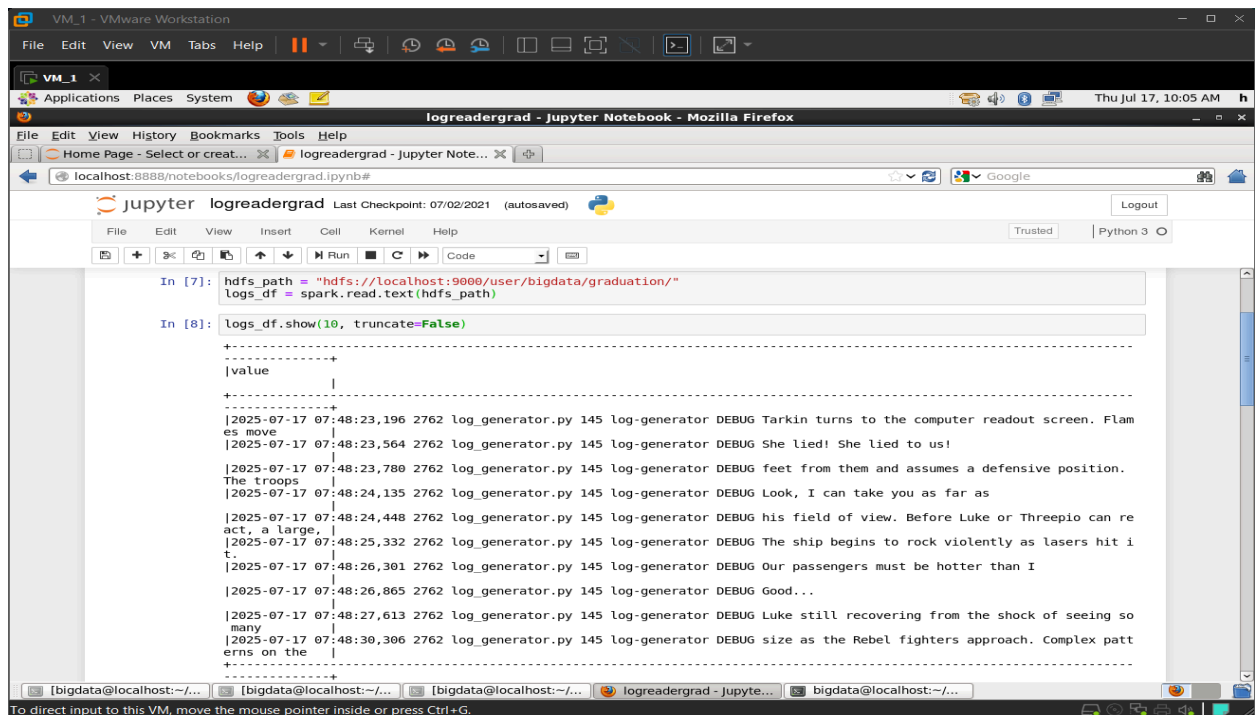
```
In [12]: from pyspark.sql import SparkSession
from pyspark.sql.functions import regexp_extract, col
from pyspark.sql.functions import explode, split

In [6]: spark = SparkSession.builder \
        .appName("GraduationLogAnalysis") \
        .getOrCreate()

In [7]: hdfs_path = "hdfs://localhost:9000/user/bigdata/graduation/"
logs_df = spark.read.text(hdfs_path)

In [8]: logs_df.show(10, truncate=False)

+-----+
|value|
+-----+
|2025-07-17 07:48:23,196 2762 log_generator.py 145 log-generator DEBUG Tarkin turns to the computer readout screen. Flamm  
es move|
|2025-07-17 07:48:23,564 2762 log_generator.py 145 log-generator DEBUG She lied! She lied to us!|
|2025-07-17 07:48:23,780 2762 log_generator.py 145 log-generator DEBUG feet from them and assumes a defensive position.  
The troops|
|2025-07-17 07:48:24,135 2762 log_generator.py 145 log-generator DEBUG Look, I can take you as far as|
|2025-07-17 07:48:24,448 2762 log_generator.py 145 log-generator DEBUG his field of view. Before Luke or Threepio can re  
act, a large,|
|2025-07-17 07:48:25,332 2762 log_generator.py 145 log-generator DEBUG The ship begins to rock violently as lasers hit i
```

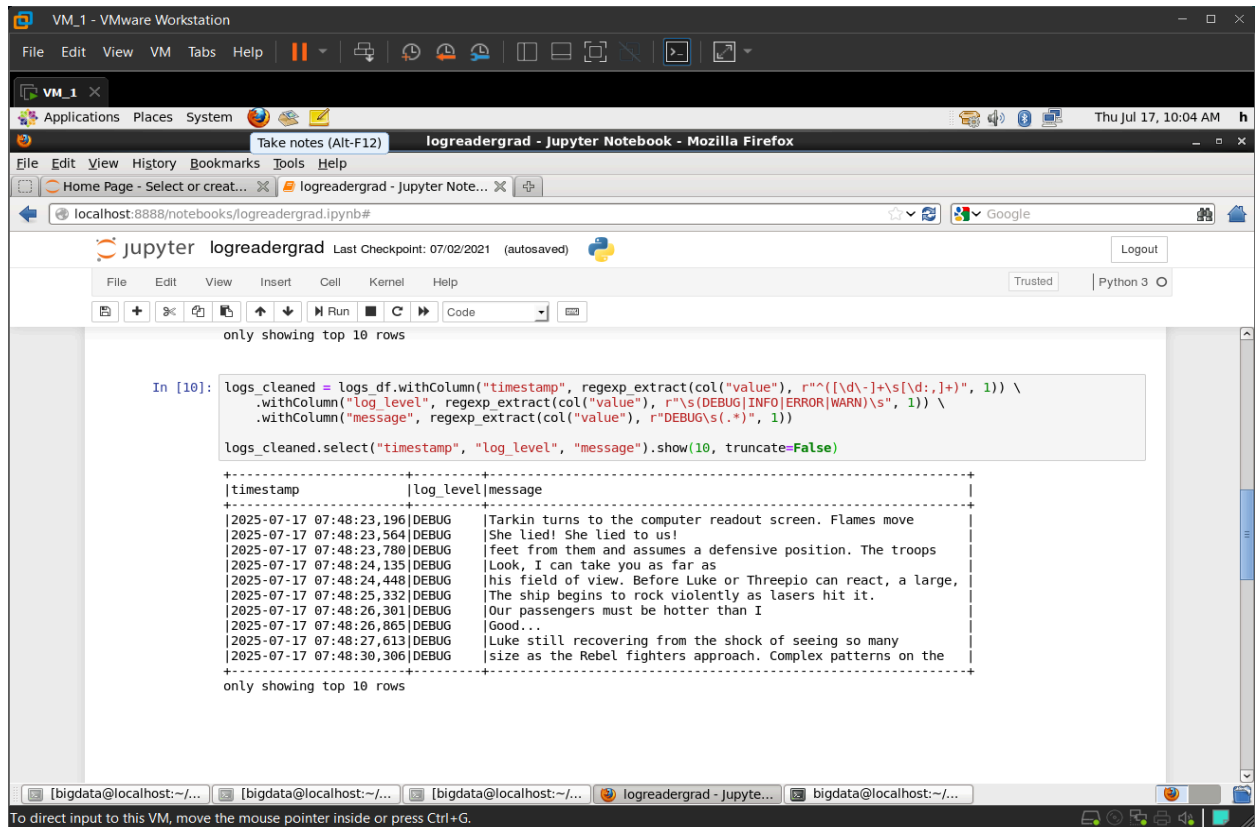


```
In [7]: hdfs_path = "hdfs://localhost:9000/user/bigdata/graduation/"
logs_df = spark.read.text(hdfs_path)

In [8]: logs_df.show(10, truncate=False)

+-----+
|value|
+-----+
|2025-07-17 07:48:23,196 2762 log_generator.py 145 log-generator DEBUG Tarkin turns to the computer readout screen. Flamm  
es move|
|2025-07-17 07:48:23,564 2762 log_generator.py 145 log-generator DEBUG She lied! She lied to us!|
|2025-07-17 07:48:23,780 2762 log_generator.py 145 log-generator DEBUG feet from them and assumes a defensive position.  
The troops|
|2025-07-17 07:48:24,135 2762 log_generator.py 145 log-generator DEBUG Look, I can take you as far as|
|2025-07-17 07:48:24,448 2762 log_generator.py 145 log-generator DEBUG his field of view. Before Luke or Threepio can re  
act, a large,|
|2025-07-17 07:48:25,332 2762 log_generator.py 145 log-generator DEBUG The ship begins to rock violently as lasers hit i  
t.|
|2025-07-17 07:48:26,301 2762 log_generator.py 145 log-generator DEBUG Our passengers must be hotter than I|
|2025-07-17 07:48:26,865 2762 log_generator.py 145 log-generator DEBUG Good...|
|2025-07-17 07:48:27,613 2762 log_generator.py 145 log-generator DEBUG Luke still recovering from the shock of seeing so  
many|
|2025-07-17 07:48:30,306 2762 log_generator.py 145 log-generator DEBUG size as the Rebel fighters approach. Complex patt  
erns on the|
```

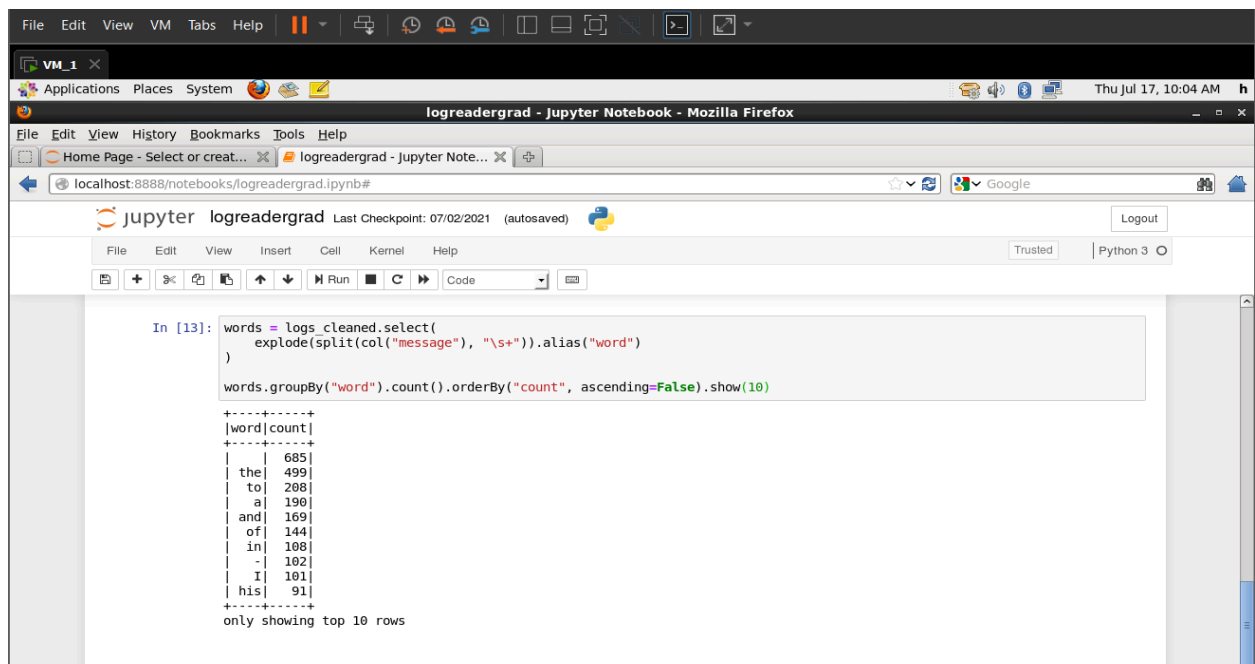
I performed some operations to improve the data presentation.



The screenshot shows a Jupyter Notebook interface within a VMware Workstation VM. The notebook is titled 'logreadergrad' and is running on a Python 3 kernel. The code in the cell performs the following operations:

- Creates a new column 'timestamp' by extracting the date and time from the 'value' column using a regular expression.
- Creates a new column 'log_level' by extracting the log level (DEBUG, INFO, ERROR, WARN) from the 'value' column using a regular expression.
- Creates a new column 'message' by extracting the message content from the 'value' column using a regular expression.
- Displays the first 10 rows of the cleaned data using `show(10, truncate=False)`.

timestamp	log_level	message
2025-07-17 07:48:23,196	DEBUG	Tarkin turns to the computer readout screen. Flames move
2025-07-17 07:48:23,564	DEBUG	She lied! She lied to us!
2025-07-17 07:48:23,780	DEBUG	feet from them and assumes a defensive position. The troops
2025-07-17 07:48:24,135	DEBUG	Look, I can take you as far as
2025-07-17 07:48:24,448	DEBUG	this field of view. Before Luke or Threepio can react, a large,
2025-07-17 07:48:25,332	DEBUG	The ship begins to rock violently as lasers hit it.
2025-07-17 07:48:26,301	DEBUG	Our passengers must be hotter than I
2025-07-17 07:48:26,865	DEBUG	Good...
2025-07-17 07:48:27,613	DEBUG	Luke still recovering from the shock of seeing so many
2025-07-17 07:48:30,306	DEBUG	size as the Rebel fighters approach. Complex patterns on the

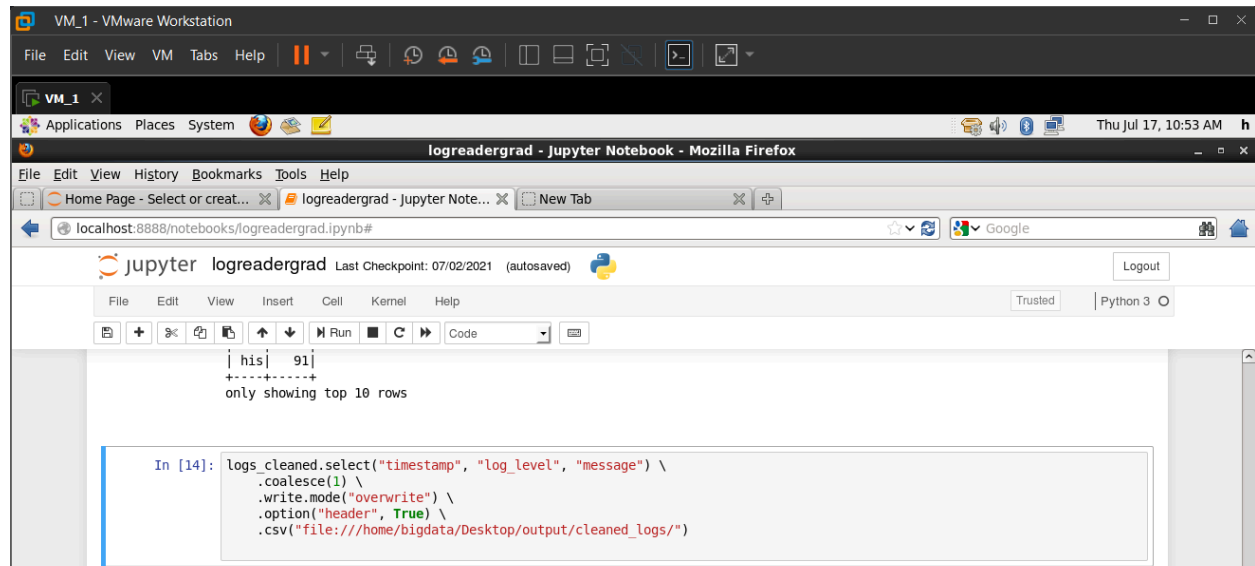


The screenshot shows the same Jupyter Notebook interface. The code in the cell performs the following operations:

- Creates a new column 'word' by splitting the 'message' column into individual words.
- Groups the data by 'word' and counts the frequency of each word.
- Orders the results by 'count' in descending order and displays the top 10 words.

word	count
I	685
the	499
to	208
a	190
and	169
of	144
in	100
-	102
I	101
his	91

Finally, I saved the processed results using the following method:



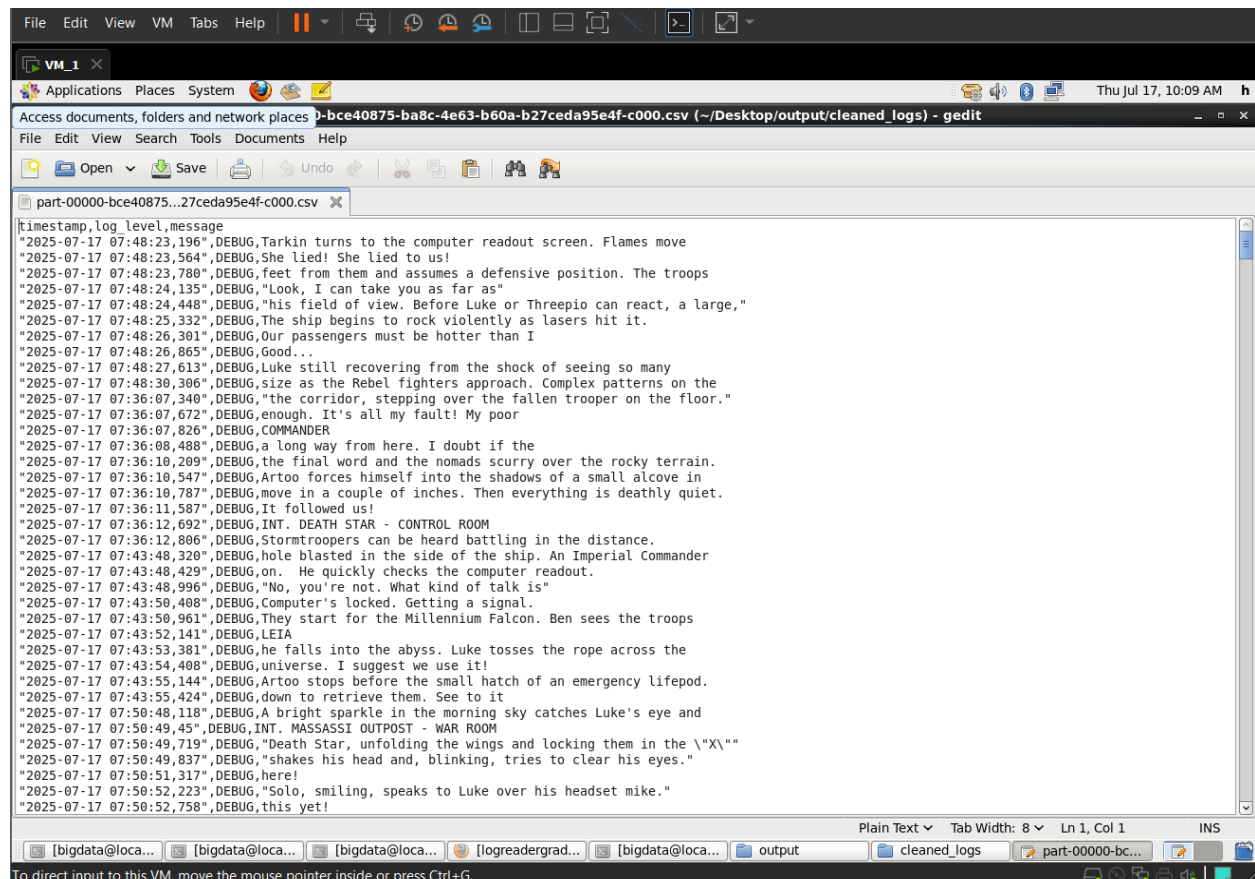
The screenshot shows a VMware Workstation window with a Jupyter Notebook titled 'logreadergrad - Jupyter Notebook - Mozilla Firefox'. The notebook is running on localhost:8888. The code in the cell is as follows:

```
In [14]: logs_cleaned.select("timestamp", "log_level", "message") \
        .coalesce(1) \
        .write.mode("overwrite") \
        .option("header", True) \
        .csv("file:///home/bigdata/Desktop/output/cleaned_logs/")
```

Below the code, there is a text area with the following text:

```
his| 91|
-----+
only showing top 10 rows
```

The .csv:



The screenshot shows a Gedit window titled 'part-00000-bce40875...27ceda95e4f-c000.csv (~/.Desktop/output/cleaned_logs) - gedit'. The window displays a CSV file with the following content:

```
timestamp,log_level,message
"2025-07-17 07:48:23,196",DEBUG,Tarkin turns to the computer readout screen. Flames move
"2025-07-17 07:48:23,564",DEBUG,She lied! She lied to us!
"2025-07-17 07:48:23,780",DEBUG,feet from them and assumes a defensive position. The troops
"2025-07-17 07:48:24,135",DEBUG,"Look, I can take you as far as"
"2025-07-17 07:48:24,448",DEBUG,"his field of view. Before Luke or Threepio can react, a large,"
"2025-07-17 07:48:25,332",DEBUG,The ship begins to rock violently as lasers hit it.
"2025-07-17 07:48:26,301",DEBUG,Our passengers must be hotter than I
"2025-07-17 07:48:26,865",DEBUG,Good...
"2025-07-17 07:48:27,613",DEBUG,Luke still recovering from the shock of seeing so many
"2025-07-17 07:48:30,306",DEBUG,size as the Rebel fighters approach. Complex patterns on the
"2025-07-17 07:36:07,340",DEBUG,"the corridor, stepping over the fallen trooper on the floor."
"2025-07-17 07:36:07,672",DEBUG,enough. It's all my fault! My poor
"2025-07-17 07:36:07,926",DEBUG,COMMANDER
"2025-07-17 07:36:08,488",DEBUG,a long way from here. I doubt if the
"2025-07-17 07:36:10,209",DEBUG,the final word and the nomads scurry over the rocky terrain.
"2025-07-17 07:36:10,547",DEBUG,Artoo forces himself into the shadows of a small alcove in
"2025-07-17 07:36:10,787",DEBUG,move in a couple of inches. Then everything is deathly quiet.
"2025-07-17 07:36:11,587",DEBUG,It followed us!
"2025-07-17 07:36:12,692",DEBUG,INT. DEATH STAR - CONTROL ROOM
"2025-07-17 07:36:12,806",DEBUG,Stormtroopers can be heard battling in the distance.
"2025-07-17 07:43:48,320",DEBUG,hole blasted in the side of the ship. An Imperial Commander
"2025-07-17 07:43:48,429",DEBUG,on. He quickly checks the computer readout.
"2025-07-17 07:43:48,996",DEBUG,"No, you're not. What kind of talk is"
"2025-07-17 07:43:50,408",DEBUG,Computer's locked. Getting a signal.
"2025-07-17 07:43:50,961",DEBUG,They start for the Millennium Falcon. Ben sees the troops
"2025-07-17 07:43:52,141",DEBUG,LEIA
"2025-07-17 07:43:53,381",DEBUG,he falls into the abyss. Luke tosses the rope across the
"2025-07-17 07:43:54,408",DEBUG,universe. I suggest we use it!
"2025-07-17 07:43:55,144",DEBUG,Artoo stops before the small hatch of an emergency lifepod.
"2025-07-17 07:43:55,424",DEBUG,down to retrieve them. See to it
"2025-07-17 07:50:48,118",DEBUG,A bright sparkle in the morning sky catches Luke's eye and
"2025-07-17 07:50:49,45",DEBUG,INT. MASSASSI OUTPOST - WAR ROOM
"2025-07-17 07:50:49,719",DEBUG,"Death Star, unfolding the wings and locking them in the "\X""
"2025-07-17 07:50:49,837",DEBUG,"shakes his head and, blinking, tries to clear his eyes."
"2025-07-17 07:50:51,317",DEBUG,here!
"2025-07-17 07:50:52,223",DEBUG,"Solo, smiling, speaks to Luke over his headset mike."
"2025-07-17 07:50:52,758",DEBUG,this yet!
```

Notes:

- All directories used during the project, including log files and output folders, were created and managed locally on my **Desktop** for easier access and testing.
- All **scripts**, **Python code**, **configuration** files (including Flume configurations), and notebooks used in this project are fully available on my personal GitHub repository for reference and future use

GitHub Profile: <https://github.com/khattabx>

Repo Name: [NTI-Graduation-Project](#)