

Group Project - 1

Group Number: 2

Brian Honea (responsible for part B ONLY SHOWN in this version)

Harshita Kala

Dorsey Kirpatrick

Basil Muhammad

Project Part B: Flume Ingestion

Screen shot of python code: I used Twelve Data API to pull data for AAPL, IBM, GOOGL, MSFT, and TSLA. I took time to figure out how to get the API to communicate appropriately and finally got it to work. I understand this submission is late.



```
import requests
import json
import socket
import time
from datetime import datetime

# Replace with your Twelve Data API key
API_KEY = "3953749cbf684e8306a2176ac1c3d478"

# Function to fetch stock data from Twelve Data API
def fetch_stock_data(symbol):
    url = "https://api.twelvedata.com/time_series?symbol={}&interval=1min&apikey={}".format(symbol, API_KEY)
    response = requests.get(url)

    # Check if the response is successful
    if response.status_code == 200:
        try:
            data = response.json()
            if 'values' in data:
                return data
            else:
                print("Error in response:", data)
                return {}
        except ValueError as e:
            print("JSON Decode Error:", e)
            return {}
    else:
        print("Error fetching data for {}: Status Code: {}, Response Content: {}".format(symbol, response.status_code, response.text))
        return {}

# Define the stock symbols to monitor
symbols = ["AAPL", "IBM", "GOOGL", "MSFT", "TSLA"]

# Define the localhost server settings
HOST = 'localhost'
PORT = 12345 # You can choose any available port

# Function to send data to localhost server
def send_data_to_localhost(data):
    s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
    s.connect((HOST, PORT))
    s.sendall(data + "\n") # Send data with a newline character
    s.close()

# Function to determine if the market is open (9:30 am to 4:00 pm ET)
def is_market_open():
    return True

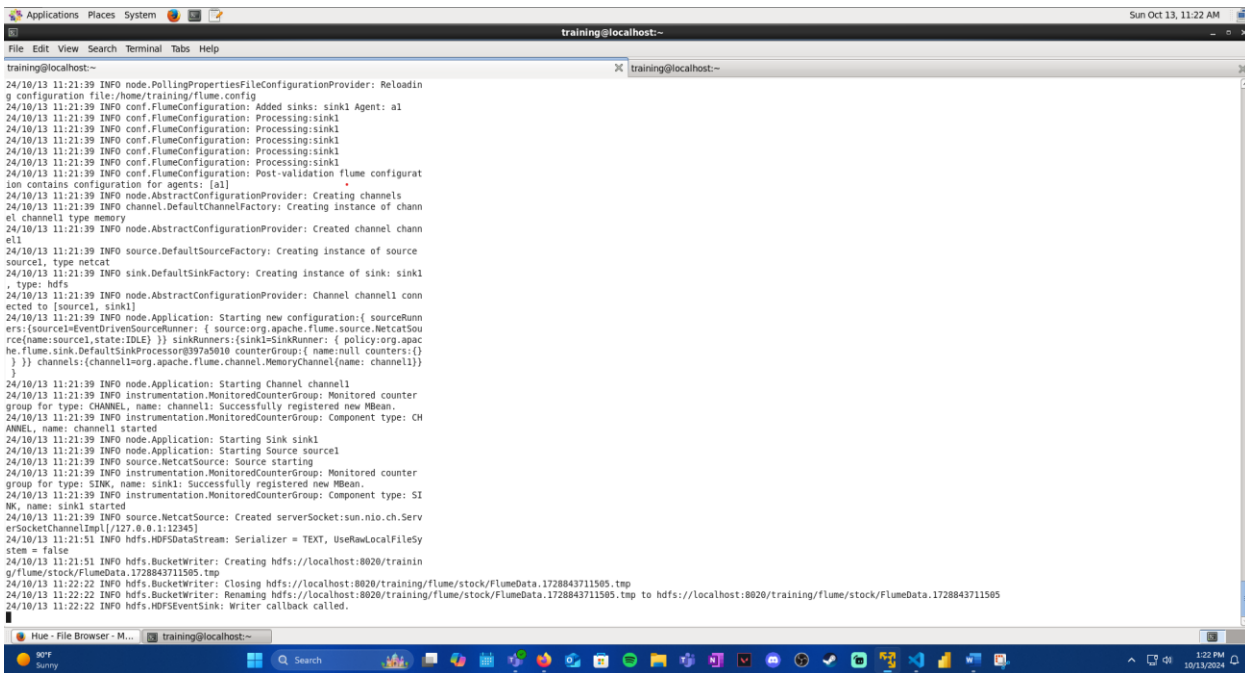
    #now = datetime.now()
    #market_open = now.replace(hour=9, minute=30, second=0, microsecond=0)
    #market_close = now.replace(hour=16, minute=0, second=0, microsecond=0)

    # Checking if it's a weekday (Mon-Fri)
    #if now.weekday() < 5:
    #    return market_open <= now <= market_close
    #return False

# Main loop to fetch stock data and send to localhost
```

Python output in console:

As you can see the python file is running as it should in the VM server. It is outputting the data gathered from the API for AAPL, IBM, GOOGL, MSFT, and TSLA.

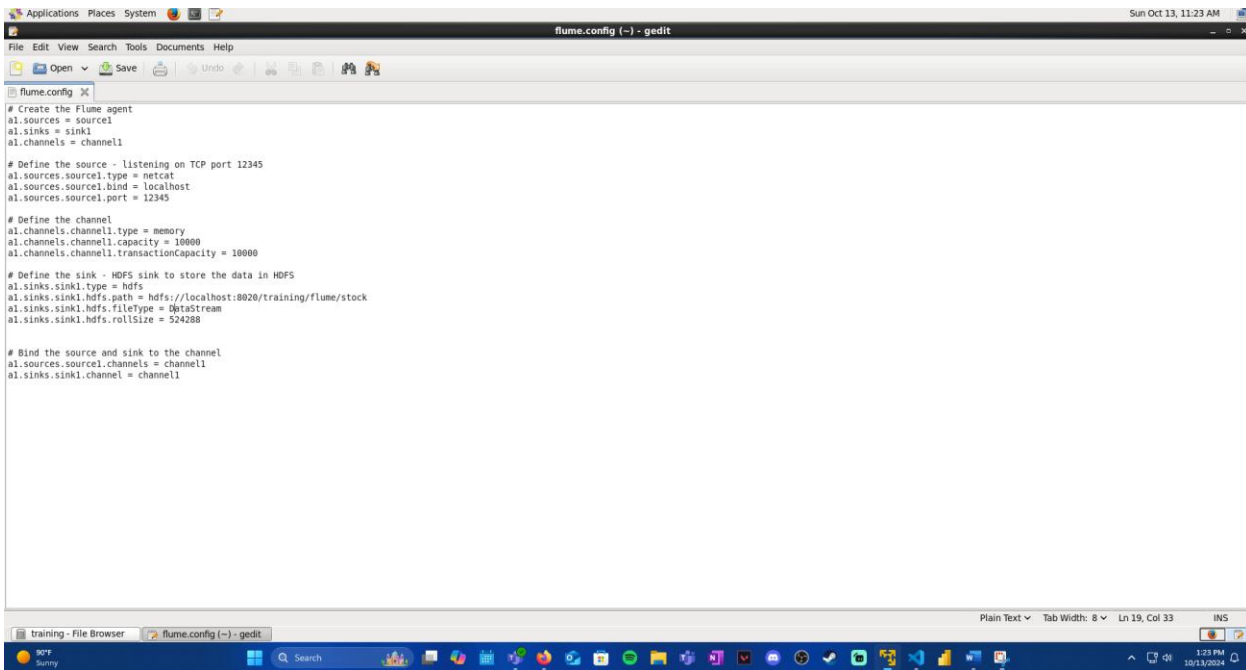


```
24/10/13 11:21:39 INFO node.PollingPropertiesFileConfigurationProvider: Reloading configuration file:/home/training/flume.conf
24/10/13 11:21:39 INFO conf.FlumeConfiguration: Added sinks: sink1 Agent: a1
24/10/13 11:21:39 INFO conf.FlumeConfiguration: Processing:sink1
24/10/13 11:21:39 INFO conf.FlumeConfiguration: Processing:sink1
24/10/13 11:21:39 INFO conf.FlumeConfiguration: Processing:sink1
24/10/13 11:21:39 INFO conf.FlumeConfiguration: Processing:sink1
24/10/13 11:21:39 INFO conf.FlumeConfiguration: Post-validation flume configuration contains configuration for agents: [a1]
24/10/13 11:21:39 INFO node.AbstractConfigurationProvider: Creating channels
24/10/13 11:21:39 INFO channel.DefaultChannelFactory: Creating instance of channel: channel1 type: memory
24/10/13 11:21:39 INFO node.AbstractConfigurationProvider: Created channel: channel1
24/10/13 11:21:39 INFO source.DefaultSourceFactory: Creating instance of source: source1, type: netcat
24/10/13 11:21:39 INFO sink.DefaultSinkFactory: Creating instance of sink: sink1, type: hdfs
24/10/13 11:21:39 INFO node.AbstractConfigurationProvider: Channel: channel1 connected to [source1, sink1]
24/10/13 11:21:39 INFO node.Application: Starting new configuration: { sourceRunners: {source1:org.apache.flume.source.NetcatSource, state:IDLE } }, sinkRunners: {sink1:org.apache.flume.sink.DefaultSinkProcessor@897a5010, counterGroup: { name:null, counters:{} } }, channels: {channel1:org.apache.flume.channel.MemoryChannel(name: channel1)} }
24/10/13 11:21:39 INFO node.Application: Starting Channel: channel1
24/10/13 11:21:39 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: CHANNEL, name: channel1: Successfully registered new MBean.
24/10/13 11:21:39 INFO instrumentation.MonitoredCounterGroup: Component type: CHANNEL, name: channel1 started
24/10/13 11:21:39 INFO node.Application: Starting Sink: sink1
24/10/13 11:21:39 INFO source.NetcatSource: Source starting
24/10/13 11:21:39 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: SINK, name: sink1: Successfully registered new MBean.
24/10/13 11:21:39 INFO instrumentation.MonitoredCounterGroup: Component type: SINK, name: sink1 started
24/10/13 11:21:39 INFO source.NetcatSource: Created serverSocket: sun.nio.ch.ServerSocketChannelImpl[/127.0.0.1:12345]
24/10/13 11:21:51 INFO hdfs.HDFSDataStream: Serializer = TEXT, UseRawLocalFileSystem = false
24/10/13 11:21:51 INFO hdfs.BucketWriter: Creating hdfs://localhost:8020/training/flume/stock/FlumeData.1728843711505.tmp
24/10/13 11:22:22 INFO hdfs.BucketWriter: Closing hdfs://localhost:8020/training/flume/stock/FlumeData.1728843711505.tmp
24/10/13 11:22:22 INFO hdfs.BucketWriter: Renaming hdfs://localhost:8020/training/flume/stock/FlumeData.1728843711505.tmp to hdfs://localhost:8020/training/flume/stock/FlumeData.1728843711505
24/10/13 11:22:22 INFO hdfs.HDFSEventSink: Writer callback called.
```

Flume Configuration File:

I realized I had an error in my last file, I specified the HDFS path as “/user/training/~” where I needed to remove “User” from the path. Now it correctly outputs to HDFS using “/training/flume/stock”. I tried many variations to the sink setup, and no matter what I try I cannot seem to get the files to only roll when the data reaches 512 KB as the instructions specify. I have tried adding other options and even a buffer but it would still roll before the file reached 512KB, as you will see in my HDFS screenshot below. I believe this may be due to the data not coming over fast enough/ at a larger size. I tried with this code and still had it rolling into HDFS before reaching 512KB in size:

```
# Define the sink - HDFS sink to store the data in HDFS
a1.sinks.sink1.type = hdfs
a1.sinks.sink1.hdfs.path = hdfs://localhost:8020/training/flume/stock
a1.sinks.sink1.hdfs.fileType = DataStream
a1.sinks.sink1.hdfs.rollSize = 524288 # Set to 512 KB
a1.sinks.sink1.hdfs.rollInterval = 0 # Do not roll based on time
a1.sinks.sink1.hdfs.rollCount = 0 # Do not roll based on event count
a1.sinks.sink1.hdfs.batchSize = 10000 # Larger batch size to ensure more data is buffered
a1.sinks.sink1.hdfs.bufferSize = 131072 # 128 KB buffer size for better performance
```



Here is the file I ended up executing the code with

Here is my HDFS file output:

