

Scalable Cloud-Based Sentiment Analysis System for Amazon Electronics Reviews

Shwe Moe Thant
National College of Ireland
Email: x24170399@student.ncirl.ie

I. PHASE 1: DESIGN AND SETUP

A. Problem Analysis and Use Case

The Amazon Electronics Reviews dataset, obtained from <https://www.kaggle.com/datasets/zhipegluo89/electronics> [1], contains millions of consumer reviews in JSON format. Each review includes attributes such as overall (rating, 1–5), reviewText, reviewTime, reviewerID, asin, verified, and others.

```
1 {
2   "overall": 5,
3   "verified": true,
4   "reviewTime": "07 17, 2002",
5   "reviewerID": "A1N070NS9CJQ2I",
6   "asin": "0060009810",
7   "style": {"Format": " Hardcover"},
8   "reviewerName": "Teri Adams",
9   "reviewText": "This was the first time I read Garcia-
10    Aguilera... I closed the book with a feeling I had
11    grown emotionally as well.",
12   "summary": "Hit The Spot!",
13   "unixReviewTime": 1026864000
14 }
```

Listing 1. Sample JSON review from the dataset

Given the dataset's 4 GB size and unstructured nature, tasks such as extracting sentiment polarity or identifying recurring terms (e.g., “battery” or “defective”) require substantial computational resources. Sequential processing is inadequate for real-time e-commerce needs due to latency. Hence, a scalable cloud-based architecture is proposed to support parallel processing for sentiment analysis and keyword frequency detection.

Motivations for scalability:

- **Volume:** Efficiently handle millions of records using distributed computing.
- **Velocity:** Enable near real-time feedback analysis with low latency.
- **Scalability:** Dynamically scale compute resources using AWS services.
- **Resilience:** Ensure robustness and fault tolerance via distributed architecture.

This solution aligns with the core principles of cloud-native design, enabling rapid, reliable decision-making [2].

B. Data Ingestion Pipeline

To simplify processing, a 4 GB subset of the original 11 GB dataset was created using a Python script:

```
1 import os
2
3 input_file = os.path.expanduser("~/Downloads/Electronics.
4   json")
5 output_file = os.path.expanduser("~/Downloads/
6   Electronics_4GB.json")
7 target_size = 4 * 1024 * 1024 * 1024 # 4 GB in bytes
8
9 with open(input_file, 'r', encoding='utf-8') as infile, \
10    open(output_file, 'w', encoding='utf-8') as outfile:
11     total_bytes = 0
12     for line in infile:
13         line_bytes = len(line.encode('utf-8'))
14         if total_bytes + line_bytes > target_size:
15             break
16         outfile.write(line)
17         total_bytes += line_bytes
18
19 print(f"Output file size: {total_bytes / (1024**3):.2f} GB")
```

Listing 2. Python script to extract 4 GB subset

The resulting subset was uploaded to Amazon S3 (s3://electronics-reviews-bucket/Electronics_4GB.json) for processing.

C. Distributed Ingestion Using PySpark on EMR

An EMR cluster (release emr-7.9.0, Spark 3.5) was provisioned via AWS Learner Lab with one master, two core, and one task node (all m5.xlarge, 4 vCPUs, 16 GB RAM). A bootstrap action ensured TextBlob was available for later sentiment tasks:

```
1 #!/bin/bash
2 sudo yum install -y python3 python3-pip
3 pip3 install textblob
4 python3 -m textblob.download_corpora
```

Listing 3. Bootstrap script for TextBlob installation

After SSH access and dependency setup (e.g., python3, psutil), the PySpark ingestion script was executed:

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, lower
3 from pyspark.sql.types import StructType, StructField,
4   StringType, DoubleType
5 import time, psutil, os
6 from pyspark.sql.utils import AnalysisException
7
8 def get_memory_usage():
9     process = psutil.Process(os.getpid())
10    return process.memory_info().rss / (1024 ** 2) # MB
11
12 spark = SparkSession.builder \
13     .appName("ElectronicsDataIngestion") \
14     .config("spark.hadoop.fs.s3a.endpoint", "s3.us-east-1.
15     amazonaws.com") \
```

```

14 .config("spark.default.parallelism", "32") \
15 .config("spark.executor.memory", "6g") \
16 .config("spark.executor.cores", "4") \
17 .config("spark.executor.instances", "4") \
18 .config("spark.sql.files.maxPartitionBytes", "67108864") \
19 .config("spark.dynamicAllocation.enabled", "true") \
20 .config("spark.dynamicAllocation.minExecutors", "1") \
21 .config("spark.dynamicAllocation.maxExecutors", "8") \
22 .getOrCreate()
23
24 input_path = "s3://electronics-reviews-bucket/
   Electronics_4GB.json"
25 output_path = "s3://electronics-reviews-bucket/
   processed_data/"
26
27 start_time = time.time()
28 schema = StructType([
29     StructField("asin", StringType(), True),
30     StructField("overall", DoubleType(), True),
31     StructField("reviewText", StringType(), True),
32     StructField("reviewTime", StringType(), True)
33 ])
34
35 try:
36     df = spark.read.schema(schema).option("mode", "
   PERMISSIVE").json(input_path)
37     df = df.repartition(32)
38     record_count = df.count()
39     df = df.select(
40         col("asin"),
41         col("overall"),
42         col("reviewText"),
43         col("reviewTime"),
44         lower(col("reviewText")).alias("cleaned_reviewText")
45     )
46     df.write.mode("overwrite").parquet(output_path)
47
48     print(f"Records processed: {record_count}")
49     print(f"Execution time: {time.time() - start_time:.2f}
   seconds")
50     print(f"Memory usage: {get_memory_usage():.2f} MB")
51     print(f"Output saved to: {output_path}")
52
53 except AnalysisException as e:
54     print(f"AnalysisException: {str(e)}")
55 except Exception as e:
56     print(f"Error: {str(e)}")
57 finally:
58     spark.stop()

```

Listing 4. PySpark script for scalable data ingestion

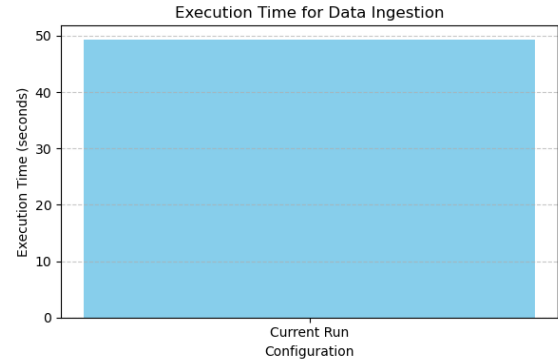


Fig. 1. Execution Time for Data Ingestion

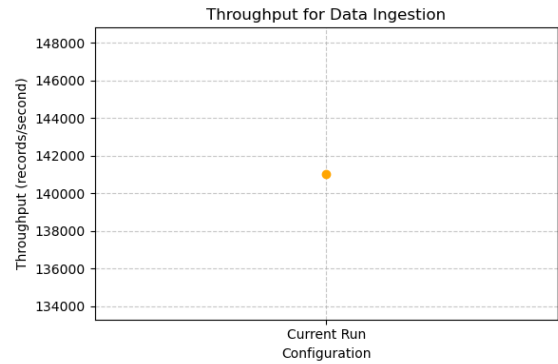


Fig. 2. Throughput during Ingestion

D. Performance Evaluation

Key performance metrics from the ingestion pipeline:

- **Records Processed:** 6,956,621
- **Execution Time:** 49.32 seconds
- **Throughput:** 141,037.83 records/second
- **Partitions:** 32 (parallel execution)
- **Memory Usage:** 51.28 MB (driver), 3.0 GB per executor

REFERENCES

- [1] J. Ni, J. Li, and J. McAuley, "Justifying Recommendations using Distantly-Labeled Reviews and Fine-Grained Aspects," in *Proc. 2019 Conf. Empir. Methods Natural Lang. Process. (EMNLP-IJCNLP)*, Hong Kong, China, 2019, pp. 188–197.
- [2] I. Foster and D. B. Gannon, *Cloud Computing for Science and Engineering*. MIT Press, 2017.
- [3] H. Karau, A. Konwinski, P. Wendell, and M. Zaharia, *Learning Spark: Lightning-Fast Big Data Analysis*. O'Reilly Media, 2015.

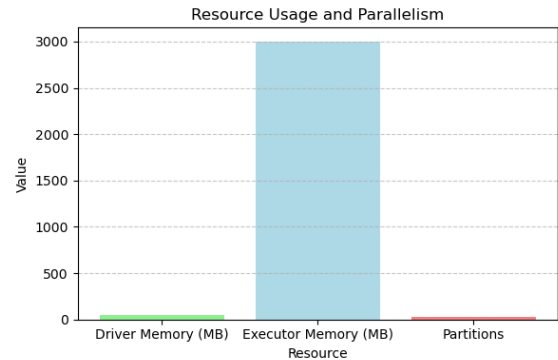


Fig. 3. Executor Resource Usage and Parallelism