

Data Engineering II - 1TD076

Assignment 2 - A Practical Introduction to Apache Pulsar

Teame Hailu Gebru

April 24, 2025

Task 1: Setting up Apache pulsar

A. Start Pulsar in the bare VM

Pulsar has been started successfully in standalone mode. Part of screenshot showing output for successful setup is shown in figure 1.

```
scription(topic=persistent://public/functions/metadata, name=c-standalone-fw-localhost-8080-function-metadata-tailer-reader-860263e978), consumerId=4, consumerName=c-standalone-fw-localhost-8080-function-metadata-tailer, address=127.0.0.1:41948)
2025-04-18T18:45:52.717+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.persistent.PersistentSubscription - [persistent://public/functions/metadata][c-standalone-fw-localhost-8080-function-metadata-tailer-reader-860263e978] Successfully closed subscription [NonDurableCursorImpl{ledger=public/functions/persistent/metadata, ackPos=18-1, readPos=18-0}]
2025-04-18T18:45:52.717+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.persistent.PersistentSubscription - [persistent://public/functions/metadata][c-standalone-fw-localhost-8080-function-metadata-tailer-reader-860263e978] Successfully closed dispatcher for reader
2025-04-18T18:45:52.717+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.ServerCnx - [/127.0.0.1:41948] Closed consumer, consumerId=4
2025-04-18T18:45:52.718+0000 [pulsar-client-io-78-1] INFO org.apache.pulsar.client.impl.ConsumerImpl - [persistent://public/functions/metadata] [c-standalone-fw-localhost-8080-function-metadata-tailer-reader-860263e978] Closed consumer
2025-04-18T18:45:52.719+0000 [pulsar-external-listener-79-1] INFO org.apache.pulsar.functions.worker.FunctionMetadataTopicTailer - Stopped function metadata tailer
2025-04-18T18:45:52.719+0000 [pulsar-external-listener-79-1] INFO org.apache.pulsar.functions.worker.FunctionMetadataManager - FunctionMetadataManager done becoming leader
2025-04-18T18:45:52.876+0000 [pulsar-web-63-6] INFO org.eclipse.jetty.server.RequestLog - 127.0.0.1 - - [18/Apr/2025:18:45:52 +0000] "GET /admin/v2/persistent/public/functions/coordinate/status?getPreciseBacklog=false&subscriptionBacklogSize=false HTTP/1.1" 200 2080 "-" "Pulsar-Java-v2.9.4" 15
2025-04-18T18:45:52.890+0000 [worker-scheduler-0] INFO org.apache.pulsar.functions.worker.SchedulerManager - Schedule summary - execution time: 0.165683186 sec | total unassigned: 0 | stats: {"Added": 0, "Updated": 0, "Removed": 0}
{
  "c-standalone-fw-localhost-8080" : {
    "originalNumAssignments" : 0,
    "finalNumAssignments" : 0,
    "instancesAdded" : 0,
    "instancesRemoved" : 0,
    "instancesUpdated" : 0,
    "alive" : true
  }
}
2025-04-18T18:46:20.735+0000 [pulsar-web-63-7] INFO org.eclipse.jetty.server.RequestLog - 127.0.0.1 - - [18/Apr/2025:18:46:20 +0000] "GET /admin/v2/persistent/public/functions/coordinate/status?getPreciseBacklog=false&subscriptionBacklogSize=false HTTP/1.1" 200 2080 "-" "Pulsar-Java-v2.9.4" 15
2025-04-18T18:46:20.756+0000 [pulsar-web-63-1] INFO org.eclipse.jetty.server.RequestLog - 127.0.0.1 - - [18/Apr/2025:18:46:20 +0000] "GET /admin/v2/persistent/public/functions/coordinate/status?getPreciseBacklog=false&subscriptionBacklogSize=false HTTP/1.1" 200 2080 "-" "Pulsar-Java-v2.9.4" 8
```

Figure 1: Output for successful setup of Apache pulsar

B. Setting up pulsar in a container in a VM

The setup of Apache Pulsar in a Docker container within a virtual machine was successfully started, as shown in figure 2.

```

2025-04-18T18:48:43.636+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.persistent.PersistentSubscription - [persistent://public/functions/assignments][c-standalone-fw-localhost-8080-function-assignment-tailer-reader-af83ec9743] Successfully closed dispatcher for reader
2025-04-18T18:48:43.636+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.ServerCnx - [/127.0.0.1:57512] Closed consumer, consumerId=3
2025-04-18T18:48:43.637+0000 [pulsar-client-io-78-1] INFO org.apache.pulsar.client.impl.ConsumerImpl - [persistent://public/functions/assignments] [c-standalone-fw-localhost-8080-function-assignment-tailer-reader-af83ec9743] Closed consumer
2025-04-18T18:48:43.639+0000 [pulsar-external-listener-79-1] INFO org.apache.pulsar.functions.worker.FunctionMetadataManager - FunctionMetadataManager becoming leader by creating exclusive producer
2025-04-18T18:48:43.651+0000 [function-metadata-tailer-thread] INFO org.apache.pulsar.client.impl.ConsumerImpl - [persistent://public/functions/metadata][c-standalone-fw-localhost-8080-function-metadata-tailer-reader-af83ec9743] Get topic last message Id
2025-04-18T18:48:43.652+0000 [pulsar-client-io-78-1] INFO org.apache.pulsar.client.impl.ConsumerImpl - [persistent://public/functions/metadata][c-standalone-fw-localhost-8080-function-metadata-tailer-reader-af83ec9743] Successfully get last message Id -1
2025-04-18T18:48:43.653+0000 [function-metadata-tailer-thread] INFO org.apache.pulsar.functions.worker.FunctionMetadataTopicTailer - metadata tailer thread exiting
2025-04-18T18:48:43.653+0000 [pulsar-external-listener-79-1] INFO org.apache.pulsar.functions.worker.FunctionMetadataTopicTailer - Stopping function metadata tailer
2025-04-18T18:48:43.656+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.ServerCnx - [/127.0.0.1:57512] Closing consumer, consumerId=4
2025-04-18T18:48:43.657+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.AbstractDispatcherSingleActiveConsumer - Removing consumer Consumer[subscription=persistentSubscription:topic=persistent://public/functions/metadata, name=c-standalone-fw-localhost-8080-function-metadata-tailer-reader-af83ec9743], consumerId=4, consumerName=c-standalone-fw-localhost-8080-function-metadata-tailer, address=/127.0.0.1:57512]
2025-04-18T18:48:43.657+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.persistent.PersistentSubscription - [persistent://public/functions/metadata][c-standalone-fw-localhost-8080-function-metadata-tailer-reader-af83ec9743] Successfully closed subscription [NonDurableCursorImpl{ledger=public/functions/persistent/metadata, ackPos=85, readPos=85, 0}]
2025-04-18T18:48:43.658+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.persistent.PersistentSubscription - [persistent://public/functions/metadata][c-standalone-fw-localhost-8080-function-metadata-tailer-reader-af83ec9743] Successfully closed dispatcher for reader
2025-04-18T18:48:43.660+0000 [pulsar-io-29-4] INFO org.apache.pulsar.broker.service.ServerCnx - [/127.0.0.1:57512] Closed consumer, consumerId=4
2025-04-18T18:48:43.660+0000 [pulsar-client-io-78-1] INFO org.apache.pulsar.client.impl.ConsumerImpl - [persistent://public/functions/metadata] [c-standalone-fw-localhost-8080-function-metadata-tailer-reader-af83ec9743] Closed consumer
2025-04-18T18:48:43.663+0000 [pulsar-external-listener-79-1] INFO org.apache.pulsar.functions.worker.FunctionMetadataTopicTailer - Stopped function metadata tailer
2025-04-18T18:48:43.663+0000 [pulsar-external-listener-79-1] INFO org.apache.pulsar.functions.worker.FunctionMetadataManager - FunctionMetadataManager done becoming leader
2025-04-18T18:48:43.731+0000 [pulsar-web-63-6] INFO org.eclipse.jetty.server.RequestLog - 127.0.0.1 - - [18/Apr/2025:18:48:43 +0000] "GET /admin/v2/persistent/public/functions/coordinates/status?getPreciseBacklog=false&subscriptionBacklogSize=false HTTP/1.1" 200 2081 "-" "Pulsar-Java-v2.9.4" 12
2025-04-18T18:49:17.691+0000 [pulsar-web-63-4] INFO org.eclipse.jetty.server.RequestLog - 127.0.0.1 - - [18/Apr/2025:18:49:12 +0000] "GET /admin/v2/persistent/public/functions/co
0 | stats: {"Added": 0, "Updated": 0, "Removed": 0}
{
  "c-standalone-fw-localhost-8080": {
    "originalNumAssignments": 0,
    "finalNumAssignments": 0,
    "instancesAdded": 0,
    "instancesRemoved": 0,
    "instancesUpdated": 0,
    "alive": true
  }
}
2025-04-18T18:49:12.682+0000 [pulsar-web-63-8] INFO org.eclipse.jetty.server.RequestLog - 127.0.0.1 - - [18/Apr/2025:18:49:12 +0000] "GET /admin/v2/persistent/public/functions/coordinates/status?getPreciseBacklog=false&subscriptionBacklogSize=false HTTP/1.1" 200 2081 "-" "Pulsar-Java-v2.9.4" 12
2025-04-18T18:49:17.691+0000 [pulsar-web-63-4] INFO org.eclipse.jetty.server.RequestLog - 127.0.0.1 - - [18/Apr/2025:18:49:12 +0000] "GET /admin/v2/persistent/public/functions/co

```

Figure 2: Output for successful setup of Apache pulsar in Docker container

Task 2: Producing and Consuming messages

A. Running Apache pulsar, consumer and producer on a single machine

After creating consumer.py and producer.py files in two terminal in the same VM and successfully started as shown in the figure below.

```

ubuntu@teame-hailu-gebru-11:~/Data_EngineeringII_Project/Project_II$ python3 consumer.py
2025-04-18 09:04:57.683 INFO [140608974942208] Client:88 | Subscribing on Topic :DEtopic
2025-04-18 09:04:57.684 INFO [140608954955328] ExecutorService:41 | Run io_service in a single thread
2025-04-18 09:04:57.684 INFO [140608974942208] ClientConnection:189 | [None] -> pulsar://localhost:6650] Create Client
Connection, timeout=10000
2025-04-18 09:04:57.684 INFO [140608974942208] ConnectionPool:96 | Created connection for pulsar://localhost:6650
2025-04-18 09:04:57.686 INFO [140608954955328] ClientConnection:375 | [127.0.0.1:43836 -> 127.0.0.1:6650] Connected to
broker
2025-04-18 09:04:57.716 INFO [140608954955328] HandlerBase:61 | [persistent://public/default/DEtopic, DE-sub, 0] Gettin
g connection from pool
2025-04-18 09:04:57.716 INFO [140608938169920] ExecutorService:41 | Run io_service in a single thread
2025-04-18 09:04:58.363 INFO [140608954955328] ConsumerImpl:218 | [persistent://public/default/DEtopic, DE-sub, 0] Crea
ted consumer on broker [127.0.0.1:43836 -> 127.0.0.1:6650]
Received message : 'b>Welcome to Data Engineering Course!'
2025-04-18 09:07:10.352 INFO [140608974942208] ClientImpl:500 | Closing Pulsar client with 0 producers and 1 consumers
2025-04-18 09:07:10.352 INFO [140608974942208] ConsumerImpl:885 | [persistent://public/default/DEtopic, DE-sub, 0] Clos
ing consumer for topic persistent://public/default/DEtopic
2025-04-18 09:07:10.380 INFO [140608954955328] ConsumerImpl:935 | [persistent://public/default/DEtopic, DE-sub, 0] Clos
ed consumer 0
2025-04-18 09:07:10.380 INFO [140608929777216] ClientConnection:1561 | [127.0.0.1:43836 -> 127.0.0.1:6650] Connection c
losed
2025-04-18 09:07:10.381 INFO [140608929777216] ClientConnection:263 | [127.0.0.1:43836 -> 127.0.0.1:6650] Destroyed con
nection
2025-04-18 09:07:10.381 INFO [140608954955328] ExecutorService:47 | Event loop of ExecutorService exits successfully
2025-04-18 09:07:10.381 INFO [140608938169920] ExecutorService:47 | Event loop of ExecutorService exits successfully
ubuntu@teame-hailu-gebru-11:~/Data_EngineeringII_Project/Project_II$

```

Figure 3: Output for successful setup of consumer

```

ubuntu@teame-hailu-gebru-11:~/Data_EngineeringII_Project/Project_II$ nano producer.py
ubuntu@teame-hailu-gebru-11:~/Data_EngineeringII_Project/Project_II$ python3 producer.py
2025-04-18 09:07:10.268 INFO [139919566173760] ExecutorService:41 | Run io_service in a single thread
2025-04-18 09:07:10.268 INFO [139919586160640] ClientConnection:189 | [<none> -> pulsar://localhost:6650] Create Client
Connection, timeout=10000
2025-04-18 09:07:10.268 INFO [139919586160640] ConnectionPool:96 | Created connection for pulsar://localhost:6650
2025-04-18 09:07:10.270 INFO [139919566173760] ClientConnection:375 | [127.0.0.1:40658 -> 127.0.0.1:6650] Connected to
broker
2025-04-18 09:07:10.283 INFO [139919566173760] HandlerBase:61 | [persistent://public/default/DEtopic, ] Getting connect
ion from pool
2025-04-18 09:07:10.295 INFO [139919566173760] ProducerImpl:173 | [persistent://public/default/DEtopic, ] Created produ
cer on broker [127.0.0.1:40658 -> 127.0.0.1:6650]
2025-04-18 09:07:10.322 INFO [139919586160640] ClientImpl:500 | Closing Pulsar client with 1 producers and 0 consumers
2025-04-18 09:07:10.322 INFO [139919586160640] ProducerImpl:590 | [persistent://public/default/DEtopic, standalone-1-0]
Closing producer for topic persistent://public/default/DEtopic
2025-04-18 09:07:10.341 INFO [139919566173760] ProducerImpl:630 | [persistent://public/default/DEtopic, standalone-1-0]
Closed producer
2025-04-18 09:07:10.342 INFO [139919549388352] ClientConnection:1561 | [127.0.0.1:40658 -> 127.0.0.1:6650] Connection c
losed
2025-04-18 09:07:10.342 INFO [139919549388352] ClientConnection:263 | [127.0.0.1:40658 -> 127.0.0.1:6650] Destroyed con
nection
2025-04-18 09:07:10.342 INFO [139919566173760] ExecutorService:47 | Event loop of ExecutorService exits successfully
2025-04-18 09:07:10.345 INFO [139919586160640] ProducerImpl:559 | Producer - [persistent://public/default/DEtopic, stan
dalone-1-0], [batching = off]
ubuntu@teame-hailu-gebru-11:~/Data_EngineeringII_Project/Project_II$

```

Figure 4: Output for successful setup producer

Task 3: Conceptual questions

1. What features does Apache Pulsar support have which the previous distributed data stream framework (e.g., Kafka) does not support?

Apache Pulsar support several features that enhance flexibility and scalability compared to the previous data frame works:

- **Multi-tenancy:** Unlike Kafka, Pulsar supports multi-tenancy, allowing multiple users to share the same Pulsar cluster securely.
- **Built-in geo-replication:** Pulsar ensures data availability across multiple geographical locations through its built-in support for geo-replication. This feature enables replication of data across data centers without the need for additional tools, making it highly effective when producers and consumers are distributed across different locations.
- **Segmented architecture (Broker + BookKeeper):** In contrast to Kafka's monolithic broker+storage approach, Pulsar separates serving (brokers) and storage (BookKeeper), which improves scalability and resilience.

2. What is the issue with using a batch processing approach on data at scale? How do modern data stream processing systems such as Apache pulsar can overcome the issue of batch processing?

Batch processing introduces latency and is inefficient for real-time analytics, especially as data volume grows. It waits for all data to arrive before processing, causing delays and potential system bottlenecks. Modern data processing systems overcome this by:

- trying to combine batch and stream processing capabilities into a single or multiple parallel data processing pipelines.
- Enabling real-time message streaming, where data is processed immediately upon arrival and meets the challenges of incremental processing, scalability, and fault tolerance.

- Supporting low-latency communication between producers and consumers, evident from the immediate message transmission in the consumer terminal output.
- Allowing fine-grained scaling, so system resources are used more efficiently over time, unlike static batch jobs.

By using streaming, Pulsar facilitates quicker insight generation and supports time-sensitive applications like fraud detection or live monitoring.

3. What is the underlying messaging pattern used by Apache Pulsar? What is the advantage of such a messaging pattern?

Apache Pulsar uses the publish-subscribe (pub-sub) messaging pattern. In this pattern, producers publish messages to logical channels called topics. Consumers can subscribe to these topics and receive messages in real-time.

The main advantage of the pub-sub pattern is decoupling—producers and consumers do not need to know about each other. This enables better scalability, parallel processing, and system flexibility.

4. What are different modes of subscription? When are each modes of subscription used?

- Exclusive: Only one consumer can subscribe to a topic per subscription. Ideal for guaranteed single processing.
- Failover: One active consumer, others standby. Ensures high availability.
- Shared (Round-robin): Messages are distributed across multiple consumers. Used for parallel processing.
- Key_Shared: Messages with the same key go to the same consumer. Best for ordered processing based on keys.

5. What is the role of the ZooKeeper? Is it possible to ensure reliability in streaming frameworks such as Pulsar or Kafka without ZooKeeper?

ZooKeeper is used for:

- Cluster coordination: It keeps track of which Pulsar brokers are active and manages leader election for components like the Pulsar broker service
- Metadata storage: Zookeeper stores important metadata such as topic ownership, namespace information, and broker assignments.
- configuration Management: It helps in managing configuration changes across the cluster consistently.

While Pulsar originally used ZooKeeper, newer versions aim to eliminate this dependency via the Pulsar broker-based metadata management. This makes deployment lighter and improves fault tolerance.

For small-scale deployments ZooKeeper may still be bundled or optional. For large, production-grade clusters, decoupling ZooKeeper improves maintainability.

6. Enlist different components of Pulsar?

1. Broker – Handles message dispatch between clients and BookKeeper.
2. BookKeeper – Manages persistent message storage.
3. ZooKeeper – Manages cluster metadata (optional in modern setups).
4. Topic – Logical channel for messages.
5. Namespace – A grouping mechanism within a tenant.
6. Tenant – Multi-tenant boundary (not explicitly used in Task 1 or 2).

7. Mention what is the meaning of broker in Pulsar?

A broker is a core component that acts as a message router between producers and consumers. The key responsibilities of pulsar broker are:

- Receives messages from producers.
- Dispatches messages to consumers.
- Handles subscription and delivery logic.

8. What is the role of a tenant in pulsar?

A tenant represents an isolated unit in a multi-tenant Pulsar deployment and supports isolation, authentication, authorization and quotas.

- Has its own namespaces, topics, and roles.
- Useful for multi-user environments.

9. What will happen if there is no log compaction in Apache Pulsar?

Pulsar retains all historical messages.

- Storage may fill up rapidly, increasing latency and cost.
- Consumers re-reading messages might encounter redundant or outdated data.

Log compaction helps maintain a compacted view of the latest key-value pairs—critical in data consistency for large-scale systems.

Task 4: Splitting and merging operations with Apache Pulsar

Question 1: Issue in Current Implementation

The current implementation of `conversion.py` processes each word of a string sequentially using a for-loop. This approach works for small-scale data but is not scalable for large datasets, especially when the number of words reaches millions.

The main issues include:

- **Lack of Parallelism:** The conversion operation is applied sequentially to each word. This significantly slows down processing for large input sizes.
- **Memory Constraints:** The entire input and output data must be stored in memory, which is inefficient for large-scale text data.
- **No Distribution:** The implementation is limited to a single machine, making it unable to leverage distributed computing environments.
- **No Fault Tolerance:** If the process fails midway, all data must be reprocessed, as there is no message-level reliability.

These limitations make the current approach unsuitable for handling big data in real-time or large-scale scenarios.

Question 2: Redesign using Apache Pulsar

To address the scalability and performance limitations in the current implementation, Apache Pulsar can be used to split and merge data operations across multiple components.

a. Splitting and Merging using Pulsar

The system is redesigned as follows:

- **Producer:** Splits the input string into individual words and sends each word as a separate message to a Pulsar topic (`conversion-topic`).
- **Consumers:** Multiple consumers subscribe to the topic and apply a conversion operation (e.g. uppercase) on each word.
- **Merger:** A final component (or consumer) subscribes to the processed word topic (`result-topic`) and reconstructs the full converted string.

This architecture enables parallel processing and efficient message handling.

b. Architecture Diagram

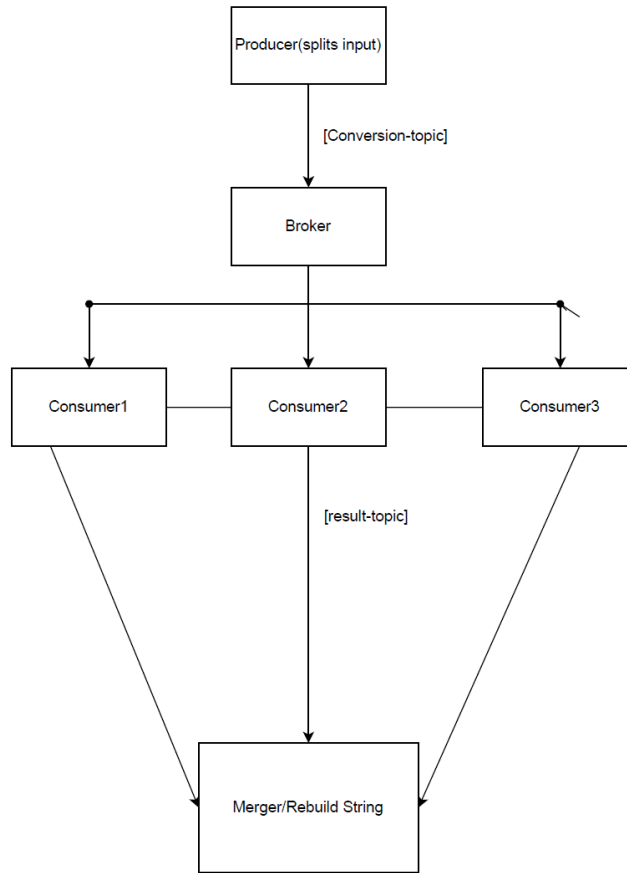


Figure 5: Architecture diagram

c. Labels

- **Broker:** Apache Pulsar broker that routes messages between producers and consumers.
- **Producer:** Sends individual words to the topic.
- **Consumers:** Perform conversion logic and send results to another topic.
- **Merger:** Reassembles the final output string from converted messages.

This redesigned system provides scalability, parallelism, and better handling of large-scale data streams.

1 Question 3. Implementation Overview (Python)

All components were implemented using the Apache Pulsar Python client (pulsar-client==2.9.4) on a single VM using multiple terminal sessions.

producer.py: Splits the string and sends each word to conversion-topic. For python script refer appendix A.1. Output showing successful execution of producer script is shown in figure 6 below.

```

ubuntu@teamehailu:~/Data_EngineeringII_Project/Project_II/Task4$ python3 producer.py
2025-04-24 19:18:12.563 INFO [140616605832960] ExecutorService:41 | Run io_service in a single thread
2025-04-24 19:18:12.563 INFO [140616624105280] ClientConnection:189 | [<none> -> pulsar://localhost:6650] Create Client
Connection, timeout=10000
2025-04-24 19:18:12.563 INFO [140616624105280] ConnectionPool:96 | Created connection for pulsar://localhost:6650
2025-04-24 19:18:12.568 INFO [140616605832960] ClientConnection:375 | [127.0.0.1:45802 -> 127.0.0.1:6650] Connected to
broker
2025-04-24 19:18:12.574 INFO [140616605832960] HandlerBase:61 | [persistent://public/default/conversion-topic, ] Gettin
g connection from pool
2025-04-24 19:18:12.583 INFO [140616605832960] ProducerImpl:173 | [persistent://public/default/conversion-topic, ] Crea
ted producer on broker [127.0.0.1:45802 -> 127.0.0.1:6650]
2025-04-24 19:18:12.683 INFO [140616624105280] ClientImpl:501 | Closing Pulsar client with 1 producers and 0 consumers
2025-04-24 19:18:12.683 INFO [140616624105280] ProducerImpl:590 | [persistent://public/default/conversion-topic, standa
alone-1-1] Closing producer for topic persistent://public/default/conversion-topic
2025-04-24 19:18:12.720 INFO [140616605832960] ProducerImpl:630 | [persistent://public/default/conversion-topic, standa
alone-1-1] Closed producer
2025-04-24 19:18:12.721 INFO [140616588965632] ClientConnection:1561 | [127.0.0.1:45802 -> 127.0.0.1:6650] Connection c
losed
2025-04-24 19:18:12.721 INFO [140616588965632] ClientConnection:263 | [127.0.0.1:45802 -> 127.0.0.1:6650] Destroyed con
nection
2025-04-24 19:18:12.721 INFO [140616605832960] ExecutorService:47 | Event loop of ExecutorService exits successfully
2025-04-24 19:18:12.723 INFO [140616624105280] ProducerImpl:559 | Producer - [persistent://public/default/conversion-to
pic, standalone-1-1], [batching = off]
ubuntu@teamehailu:~/Data_EngineeringII_Project/Project_II/Task4$

```

Figure 6: Output showing the successful execution of the producer script in the Pulsar-based system.

conversion_consumer.py: Listens to conversion-topic, converts each word to uppercase, and sends to result-topic. For python script refer appendix A.2. Output showing successful execution of consumer script is shown in figure 7 below.

```

ubuntu@teamehailu:~/Data_EngineeringII_Project/Project_II/Task4$ python3 consumer.py
2025-04-24 19:17:41.284 INFO [140189646575424] Client:88 | Subscribing on Topic :conversion-topic
2025-04-24 19:17:41.285 INFO [140189628303104] ExecutorService:41 | Run io_service in a single thread
2025-04-24 19:17:41.285 INFO [140189646575424] ClientConnection:189 | [<none> -> pulsar://localhost:6650] Create Client
Connection, timeout=10000
2025-04-24 19:17:41.285 INFO [140189646575424] ConnectionPool:96 | Created connection for pulsar://localhost:6650
2025-04-24 19:17:41.287 INFO [140189628303104] ClientConnection:375 | [127.0.0.1:60032 -> 127.0.0.1:6650] Connected to
broker
2025-04-24 19:17:41.303 INFO [140189628303104] HandlerBase:61 | [persistent://public/default/conversion-topic, conversi
on-sub, 0] Getting connection from pool
2025-04-24 19:17:41.303 INFO [140189536085760] ExecutorService:41 | Run io_service in a single thread
2025-04-24 19:17:41.363 INFO [140189628303104] ConsumerImpl:218 | [persistent://public/default/conversion-topic, conver
sion-sub, 0] Created consumer on broker [127.0.0.1:60032 -> 127.0.0.1:6650]
2025-04-24 19:17:41.366 INFO [140189628303104] HandlerBase:61 | [persistent://public/default/result-topic, ] Getting co
nnection from pool
2025-04-24 19:17:41.381 INFO [140189628303104] ProducerImpl:173 | [persistent://public/default/result-topic, ] Created
producer on broker [127.0.0.1:60032 -> 127.0.0.1:6650]

```

Figure 7: Output showing the successful execution of the consumer script in the Pulsar-based system.

merger.py: Listens to result-topic and reconstructs the final sentence. For python script refer appendix A.3. Output showing successful execution of merger script is shown in figure 8 below.

```

ubuntu@teamehailu:~/Data_EngineeringII_Project/Project_II/Task4$ python3 merger.py
2025-04-24 19:17:02.158 INFO [139794585003840] Client:88 | Subscribing on Topic :result-topic
2025-04-24 19:17:02.158 INFO [139794566731520] ExecutorService:41 | Run io_service in a single thread
2025-04-24 19:17:02.159 INFO [139794585003840] ClientConnection:189 | [<none> -> pulsar://localhost:6650] Create Client
Connection, timeout=10000
2025-04-24 19:17:02.160 INFO [139794585003840] ConnectionPool:96 | Created connection for pulsar://localhost:6650
2025-04-24 19:17:02.162 INFO [139794566731520] ClientConnection:375 | [127.0.0.1:38878 -> 127.0.0.1:6650] Connected to
broker
2025-04-24 19:17:02.179 INFO [139794566731520] HandlerBase:61 | [persistent://public/default/result-topic, result-sub,
0] Getting connection from pool
2025-04-24 19:17:02.179 INFO [139794474596096] ExecutorService:41 | Run io_service in a single thread
2025-04-24 19:17:02.467 INFO [139794566731520] ConsumerImpl:218 | [persistent://public/default/result-topic, result-sub
, 0] Created consumer on broker [127.0.0.1:38878 -> 127.0.0.1:6650]
Resultant String: I WANT TO BE CAPITALIZED
2025-04-24 19:18:12.756 INFO [139794585003840] ClientImpl:501 | Closing Pulsar client with 0 producers and 1 consumers
2025-04-24 19:18:12.756 INFO [139794585003840] ConsumerImpl:885 | [persistent://public/default/result-topic, result-sub
, 0] Closing consumer for topic persistent://public/default/result-topic
2025-04-24 19:18:12.793 INFO [139794566731520] ConsumerImpl:935 | [persistent://public/default/result-topic, result-sub
, 0] Closed consumer
2025-04-24 19:18:12.794 INFO [139794466203392] ClientConnection:1561 | [127.0.0.1:38878 -> 127.0.0.1:6650] Connection c
losed
2025-04-24 19:18:12.794 INFO [139794466203392] ClientConnection:263 | [127.0.0.1:38878 -> 127.0.0.1:6650] Destroyed con
nection
2025-04-24 19:18:12.794 INFO [139794566731520] ExecutorService:47 | Event loop of ExecutorService exits successfully
2025-04-24 19:18:12.794 INFO [139794474596096] ExecutorService:47 | Event loop of ExecutorService exits successfully
ubuntu@teamehailu:~/Data_EngineeringII_Project/Project_II/Task4$

```

Figure 8: Output showing the successful execution of the merger script in the Pulsar-based system.

Each conversion consumer runs independently, and multiple instances can be launched in parallel

to support horizontal scaling using Pulsar's Shared subscription mode.

A Appendix

A.1 Producer

```
1 import pulsar
2
3 INPUT_STRING = "I want to be capitalized"
4 client = pulsar.Client('pulsar://localhost:6650')
5 producer = client.create_producer('conversion-topic')
6
7 for word in INPUT_STRING.split():
8     producer.send(word.encode('utf-8'))
9
10 client.close()
```

A.2 Conversion_consumer

```
1 import pulsar
2
3 def convert(word):
4     return word.upper()
5
6 client = pulsar.Client('pulsar://localhost:6650')
7 consumer = client.subscribe('conversion-topic', subscription_name='conversion-sub')
8 producer = client.create_producer('result-topic')
9
10 while True:
11     msg = consumer.receive()
12     try:
13         converted = convert(msg.data().decode('utf-8'))
14         producer.send(converted.encode('utf-8'))
15         consumer.acknowledge(msg)
16     except:
17         consumer.negative_acknowledge(msg)
18
19 client.close()
```

A.3 Merger

```
1 import pulsar
2
3 client = pulsar.Client('pulsar://localhost:6650')
4 consumer = client.subscribe('result-topic', subscription_name='result-sub')
5
6 result = []
7 # For demo: assuming we know the number of words
8 expected_words = 5
9
10 for _ in range(expected_words):
```

```
11     msg = consumer.receive()
12     try:
13         result.append(msg.data().decode('utf-8'))
14         consumer.acknowledge(msg)
15     except:
16         consumer.negative_acknowledge(msg)
17
18 print("Resultant String:", ' '.join(result))
19 client.close()
```
