



**TAP PROJECT:
Behavioral Analysis for E-Commerce**

Capobianco Francesco Massimo

Indice Relazione del Progetto:

1. Introduzione e Obiettivi del progetto

- Descrizione caso di studio: E-Commerce “Guardarobe”
- Il Problema: Analisi dei dati in tempo reale vs Batch Processing
- L’Obiettivo: Monitoraggio KPI e Profilazione Utenti

2. Architettura del sistema

- Panoramica della pipeline TAP (Transazione, Analisi, Presentazione)
- Tecnologie usate
- Pipeline flusso dati

3. Data ingestion

- Generazione dei dati: WordPress/WooCommerce
- Fluent-Bit: Configurazione e ruolo di Log Collector
- Dettagli Tecnico del JSON

4. Il Cuore del sistema: Apache Kafka

- Kafka: Broker
- Struttura dei Topic Creati
- Monitoraggio tramite kafka UI

5. Elaborazione Real-Time: Apache Spark

- Spark SQL
- Spark MLlib

6. Indicizzazione e Visualizzazione:

- Logstash
- Elasticsearch
- Kibana
- KPI Realizzati e Analizzati
- Riepilogo Finale e spiegazione Use-Case

Introduzione e Obiettivi del progetto

Descrizione caso di studio: E-Commerce “Garderobe”

Il progetto simula un ambiente di e-commerce reale denominato “Garderobe”. In uno scenario di mercato moderno, un negozio online genera una mole continua di dati eterogenei: log di sistema, transazioni finanziarie, interazioni degli utenti e aggiornamenti di inventario. **La gestione e la comprensione di questi dati** sono cruciali per la competitività aziendale.

Il Problema: Analisi dei dati in tempo reale vs Batch Processing

Tradizionalmente, l'analisi dei dati avviene tramite **processi Batch** (ad esempio, report generati ogni notte). Questo approccio presenta un limite critico: la latenza. Sapere che un prodotto è andato esaurito o che una promozione non funziona con 24 ore di ritardo comporta perdite economiche.

La soluzione da adottare è l'**approccio Real-Time Streaming**: i dati vengono elaborati istantaneamente non appena vengono generati, permettendo decisioni immediate.

L'Obiettivo: Monitoraggio KPI e Profilazione Utenti

L'obiettivo principale è costruire una pipeline completa che trasformi il dato grezzo in informazione di valore. Nello specifico:

1. **Monitoraggio KPI**: Visualizzare metriche finanziarie (fatturato, vendite per prodotto).
2. **Profilazione (Machine Learning)**: Utilizzare l'Intelligenza Artificiale per classificare il comportamento dell'utente (es. acquisto Impulsivo vs Ragionato) basandosi sulla durata della sessione.

Architettura del Sistema

Panoramica della pipeline TAP (Transazione, Analisi, Presentazione)

L'architettura è caratterizzata dalla pipeline TAP, che prevede:

- **Transazione (Sorgente)**: Dove il dato nasce (WordPress/WooCommerce).
- **Analisi (Elaborazione)**: Dove il dato viene trasformato e arricchito (Spark).
- **Presentazione (Visualizzazione)**: Dove il dato viene mostrato agli stakeholder (Kibana).

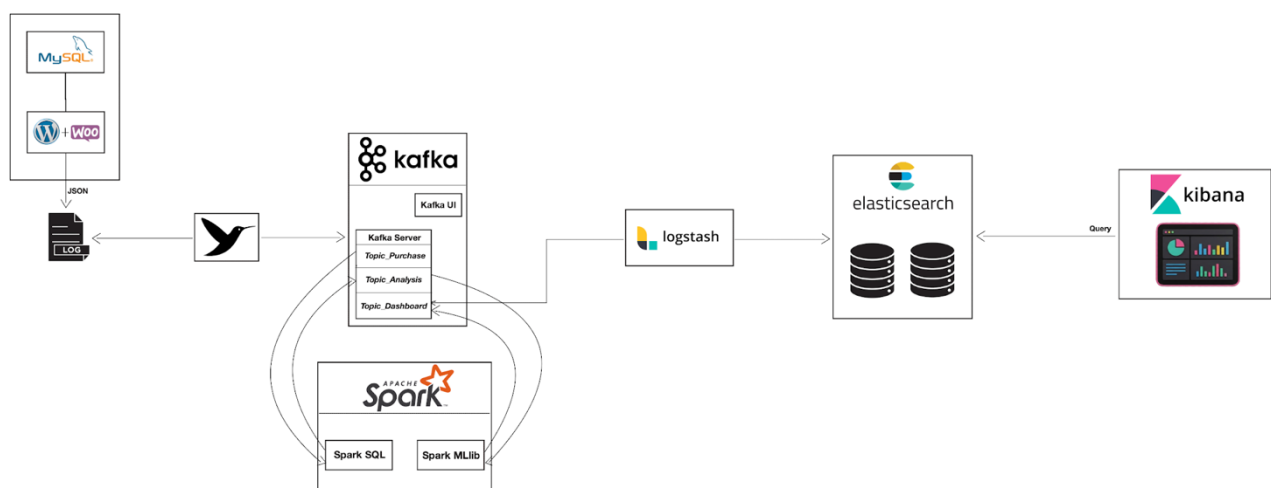
Tecnologie usate

Per garantire portabilità, isolamento e facilità di deploy, l'intera infrastruttura è basata su microservizi gestiti tramite Docker e Docker Compose. Questo permette di simulare un intero data center su una singola macchina.

A causa dei limiti del calcolatore utilizzato per realizzare il progetto, la **demo** verrà visualizzata e presentata accendendo e spegnendo i **servizi** in modo orchestrato.

Pipeline flusso dati

Il flusso segue una direzione lineare: dalla generazione del log alla dashboard finale.



Data ingestion

Generazione dei dati

Il lato **front-end** del sistema è un **sito WordPress** con **plugin WooCommerce**. Questa scelta simula fedelmente un ambiente di produzione reale. Ogni volta che un utente completa un ordine, il sistema genera un evento di log che contiene i dettagli della transazione (chi ha comprato, cosa, a che prezzo).

Fluent-Bit: Configurazione e ruolo di Log Collector

Per estrarre i dati da WordPress senza appesantire il database, utilizziamo Fluent-Bit.

- **Ruolo:** Agisce come un Log Collector che legge i file di log generati dal web server.
- **Funzionamento:** Appena una nuova riga viene scritta nel log, Fluent-Bit la cattura, la formatta e la invia al Message Broker (Kafka).

<pre>1 [SERVICE] 2 Flush 1 3 Log_Level info 4 Parsers_File /fluent-bit/etc/parsers.conf 5 6 [INPUT] 7 Name tail 8 Path /input_logs/fluent-bit-orders.log 9 Parser json_parser 10 Tag e-commerce.purchase 11 12 [OUTPUT] 13 Name stdout 14 Match e-commerce.purchase 15 16 [OUTPUT] 17 Name kafka 18 Match e-commerce.purchase 19 Brokers e-commerce_kafkaserver:9092 20 Topics Topic_Purchase 21 Format json 22 Message_Key_Field order_id</pre>	<pre>1 [PARSER] 2 Name json_parser 3 Format json 4 Time_Key timestamp 5 Time_Format %Y-%m-%d %H:%M:%S</pre>
---	---

Dettagli Tecnico del JSON

I dati viaggiano in **formato JSON** (**JavaScript Object Notation**).

È il **formato standard** per lo scambio di dati perché è **leggero** e **leggibile** sia dalle macchine che dagli umani.

Ecco un esempio di payload che include campi come: customer_email, items, total, session_duration.

```
$order_data = [
  'order_id'      => $order->get_id(),
  'timestamp'     => current_time('mysql'), // Formato YYYY-MM-DD HH:MM:SS
  'items'         => [],
  'customer'     => [
    'first_name' => $order->get_billing_first_name(),
    'email'      => $order->get_billing_email(),
  ],
  'currency'     => $order->get_currency(),
  'total'        => (float)$order->get_total(),
  'duration_session' => $duration_session,
];
```

WordPress:

Ecco inoltre il codice di quello che succede subito dopo che un utente effettua un acquisto sull'E-Commerce, ciò corrisponde all'inizio del viaggio del dato.

```
1 <?php
2 /**
3  * Blocksy functions and definitions
4  *
5  * @link https://developer.wordpress.org/themes/basics/theme-functions/
6  *
7  * @package Blocksy
8  */
9
10 if (version_compare(PHP_VERSION, '5.7.0', '<')) {
11     require get_template_directory() . '/inc/php-fallback.php';
12     return;
13 }
14
15 require get_template_directory() . '/inc/init.php';
16
17 // Codice Inserito:
18
19 add_action('init', 'start_session_timer');
20
21 function start_session_timer() {
22     if (!isset($_COOKIE['visit_start_time'])) {
23         setcookie('visit_start_time', time(), time() + 86400, "/");
24     }
25 }
26
27 add_action('wp_login', 'reset_session_on_login');
28
29 function reset_session_on_login() {
30     setcookie('visit_start_time', time(), time() + 86400, "/");
31 }
32
33 add_action('woocommerce_thankyou', 'generate_tap_order_json', 10, 1);
34
35 function generate_tap_order_json($order_id) {
36     if (!$order_id) return;
37
38     $order = wc_get_order($order_id);
39
40     // Duration Session:
41     $now = time();
42     $start_time = isset($_COOKIE['visit_start_time']) ? intval($_COOKIE['visit_start_time']) : $now;
43     $duration_session = $now - $start_time;
44
45     $order_data = [
46         'order_id' => $order->get_id(),
47         'timestamp' => current_time('mysql'), // Formato YYYY-MM-DD HH:MM:SS
48         'items' => [],
49         'customer' => [
50             'first_name' => $order->get_billing_first_name(),
51             'email' => $order->get_billing_email(),
52         ],
53         'currency' => $order->get_currency(),
54         'total' => (float)$order->get_total(),
55         'duration_session' => $duration_session,
56     ];
57
58     foreach ($order->get_items() as $item_id => $item) {
59         $order_data['items'][] = [
60             'product_name' => $item->get_name(),
61             'quantity' => $item->get_quantity(),
62             'total' => (float)$item->get_total(),
63         ];
64     }
65
66     $json_payload = json_encode($order_data);
67
68     $file_path = '/var/www/html/wp-content/uploads/fluent-bit-orders.log';
69     file_put_contents($file_path, $json_payload . PHP_EOL, FILE_APPEND);
70 }
71
72 }
```

Il Cuore del sistema: Apache Kafka

Kafka: Broker

Apache Kafka è il componente centrale che garantisce la resilienza del sistema.

Funge da broker tra chi produce i dati (Fluent-Bit) e chi li consuma (Spark).

Se Spark dovesse spegnersi per manutenzione o crash, i dati non andrebbero persi, ma rimarrebbero salvati in Kafka in attesa di essere letti. Questo disaccoppia i sistemi ed evita colli di bottiglia.

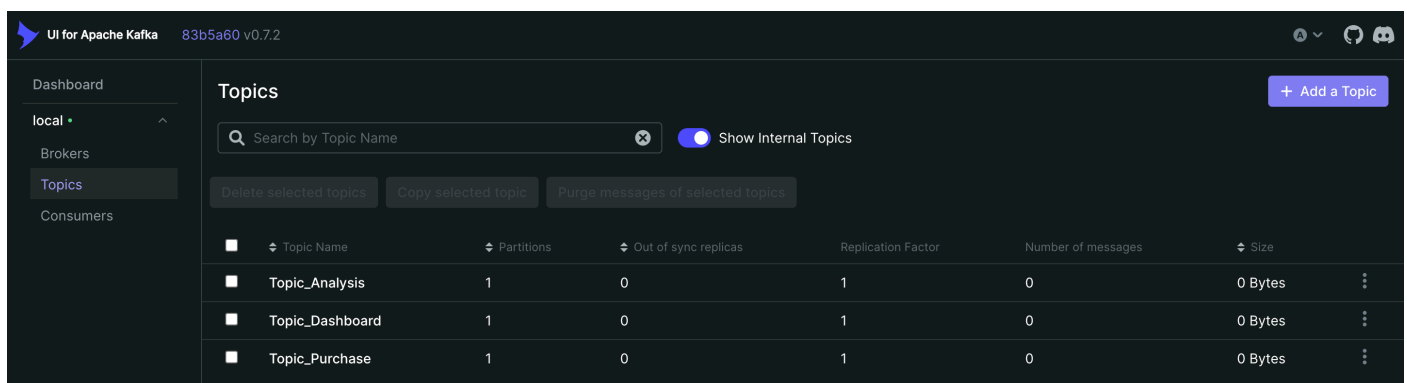
Struttura dei Topic Creati

I dati vengono organizzati in "canali" tematici chiamati Topic:

- **Topic_Purchase:** Contiene i dati grezzi appena arrivati da Fluent-Bit.
- **Topic_Analysis:** Contiene i dati elaborati da Spark SQL.
- **Topic_Dashboard:** Contiene i dati finali, arricchiti con il *modello K-Means*, pronti per Kibana.

Monitoraggio tramite Kafka UI

Per gestire visivamente i flussi e verificare che i messaggi arrivino correttamente, utilizziamo Kafka UI, un'interfaccia grafica che permette di ispezionare il contenuto dei topic e lo stato dei consumer.



The screenshot shows the Kafka UI interface. On the left is a sidebar with navigation links: Dashboard, local (selected), Brokers, Topics, and Consumers. The main area is titled 'Topics' and contains a search bar, a 'Show Internal Topics' toggle, and buttons for 'Delete selected topics', 'Copy selected topic', and 'Purge messages of selected topics'. Below this is a table listing topics:

Topic Name	Partitions	Out of sync replicas	Replication Factor	Number of messages	Size
Topic_Analysis	1	0	1	0	0 Bytes
Topic_Dashboard	1	0	1	0	0 Bytes
Topic_Purchase	1	0	1	0	0 Bytes

Elaborazione Real-Time: Apache Spark

Spark SQL (process_orders.py)

Utilizziamo Spark Streaming per elaborare i dati in movimento.

Lo script process_orders.py esegue tre operazioni fondamentali:

- **Estrazione dati:** Legge lo stream continuo da Topic_Purchase.
- **Trasformazione:** Seleziona solo le colonne utili, calcola i totali e applica regole di business (es. definire lo spender_type come "High" o "Standard" in base all'importo).
- **Caricamento su Topic:** Scrive il risultato pulito su Topic_Analysis.

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import from_json, col, to_json, struct, when
3 from pyspark.sql.types import StructType, StructField, StringType, IntegerType, FloatType, ArrayType
4
5 def main():
6     # 1. conf. iniziale
7     spark = SparkSession.builder \
8         .appName("E-Commerce-BehaviorAnalysis") \
9         .getOrCreate()
10
11     spark.sparkContext.setLogLevel("WARN")
12
13     # 2. schema json: recupero dati
14     json_schema = StructType([
15         StructField("order_id", IntegerType()),
16         StructField("timestamp", StringType()),
17         StructField("items", ArrayType(StructType([
18             StructField("product_name", StringType()),
19             StructField("quantity", IntegerType()),
20             StructField("total", FloatType())
21         ]))),
22         StructField("customer", StructType([
23             StructField("first_name", StringType()),
24             StructField("email", StringType())
25         ])),
26         StructField("currency", StringType()),
27         StructField("total", FloatType()),
28         StructField("duration_session", IntegerType())
29     ])
30
31     # 3. Lettura da Topic_Purchase
32     raw_stream = spark.readStream \
33         .format("kafka") \
34         .option("kafka.bootstrap.servers", "e-commerce-kafkaserver:9092") \
35         .option("subscribe", "Topic_Purchase") \
36         .option("startingOffsets", "earliest") \
37         .option("failOnDataLoss", "false") \
38         .load() # 'earliest' per prendere tutto lo storico riletto da FluentBit, altrimenti 'latest'
39
40     # Processo col_Value -> String -> JSON
41     json_stream = raw_stream.select(
42         from_json(col("value").cast("string"), json_schema).alias("data")
43     ).select("data.*")
44
45     # 4. Data Enrichment:
46
47     processed_stream = json_stream.withColumn(
48         "spender_type",
49         when(col("total") > 150, "High Spender").otherwise("Standard Spender")
50     ).withColumn(
51         "stock_alert",
52         when(col("total") > 500, "High Demand").otherwise("Normal Demand")
53     )
54
55     # 5. Scrittura su Topic_Analysis
56     query = processed_stream.select(
57         to_json(struct(col("data")).alias("value")) \
58         .writeStream \
59         .format("kafka") \
60         .option("kafka.bootstrap.servers", "e-commerce-kafkaserver:9092") \
61         .option("topic", "Topic_Analysis") \
62         .option("checkpointLocation", "/opt/spark-apps/checkpoint_sql") \
63         .trigger(processingTime='20 seconds') \
64         .start()
65
66     print("---- Streaming Behaviour_Analysis started ----")
67     query.awaitTermination()
68
69 if __name__ == "__main__":
70     main()
```

Spark MLlib (Machine Learning)

La componente di **intelligenza artificiale** è gestita dalla **libreria MLlib**.

- **Modello:** Viene utilizzato il *K-Means*, un modello ML di clustering.
- **Logica:** Il modello analizza una feature: la durata della sessione.

In base a questo dato, classifica l'utente in due cluster:

1. **Impulsive:** *Bassa durata*, pochi click e acquisto rapido.
2. **Reasoned:** *Alta durata* (tempo di riflessione), acquisto ponderato.

Il risultato della categoria grazie a K-Means viene aggiunto al dato e inviato a Topic_Dashboard.

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import from_json, col, to_json, struct, when
3 from pyspark.sql.types import StructType, StructField, StringType, IntegerType, FloatType, ArrayType
4
5 from pyspark.ml.feature import VectorAssembler
6 from pyspark.ml.clustering import KMeans
7 from pyspark.ml import Pipeline
8
9 def train_kmeans_model(spark):
10     csv_path = "/opt/spark-apps/training_data.csv"
11     try:
12         df_train = spark.read.option("header", "true") \
13             .option("inferSchema", "true") \
14             .csv(csv_path)
15
16         df_train = df_train.withColumn("duration_session", col("duration_session").cast("integer"))
17
18         assembler = VectorAssembler(inputCols=["duration_session"], outputCol="features")
19         kmeans = KMeans(k=2, seed=1)
20         pipeline = Pipeline(stages=[assembler, kmeans])
21
22         # Addestramento
23         model = pipeline.fit(df_train)
24
25         # Cluster
26         centers = model.stages[-1].clusterCenters()
27         impulsive_cluster = 0 if centers[0][0] < centers[1][0] else 1
28         return model, impulsive_cluster
29
30 except Exception as e:
31     print(f"Error Loading Training Data: {e}")
32     return None, 0
33
34 def main():
35     spark = SparkSession \
36         .builder \
37         .appName("E-Commerce_Machine_Learning") \
38         .getOrCreate()
39
40     spark.sparkContext.setLogLevel("WARN")
41
42     ml_model, impulsive_idx = train_kmeans_model(spark)
43
44     if ml_model is None:
45         print("Stopping due to training error")
46         return
47
48     # 1. Schema Input dal Topic_Analysis
49     json_schema = StructType([
50         StructField("order_id", IntegerType()),
51         StructField("timestamp", StringType()),
52         StructField("items", ArrayType(StructType([
53             StructField("product_name", StringType()),
54             StructField("quantity", IntegerType()),
55             StructField("total", FloatType())
56         ]))),
57         StructField("customer", StructType([
58             StructField("first_name", StringType()),
59             StructField("email", StringType())
60         ])),
61         StructField("currency", StringType()),
62         StructField("duration_session", IntegerType()),
63         StructField("spender_type", StringType()), # Extra field
64         StructField("stock_alert", StringType()) # Extra field
65     ])
66
67     # 2. Lettura da Topic_Analysis
68     raw_stream = spark.readStream \
69         .format("kafka") \
70         .option("kafka.bootstrap.servers", "e-commerce-kafkaserver:9092") \
71         .option("subscribe", "Topic_Analysis") \
72         .option("startingOffsets", "earliest") \
73         .option("failOnDataLoss", "false") \
74         .load()
75
76     json_stream = raw_stream.select(
77         from_json(col("value").cast("string"), json_schema).alias("data")
78     ).select("data.*")
79
80     # 3. ML Prediction
81     prediction_stream = ml_model.transform(json_stream)
82
83     final_stream = prediction_stream.withColumn(
84         "purchase_behaviour",
85         when(col("prediction") == impulsive_idx, "Impulsive")
86         .otherwise("Reasoned")
87     ).drop("features", "prediction")
88
89     # 4. Scrittura su Topic_Dashboard
90     query = final_stream.select(to_json(struct(col("data")).alias("value")) \
91         .writeStream \
92         .format("kafka") \
93         .option("kafka.bootstrap.servers", "e-commerce-kafkaserver:9092") \
94         .option("topic", "Topic_Dashboard") \
95         .option("checkpointLocation", "/opt/spark-apps/checkpoint_ml") \
96         .trigger(processingTime='20 seconds') \
97         .start()
98
99     print("---- Machine Learning Started ----")
100     query.awaitTermination()
101
102 if __name__ == "__main__":
103     main()
```

Indicizzazione e Visualizzazione:

Logstash:

Logstash preleva i dati finali da Kafka per inviarli al database.

- **Gestione Offset:** Un punto critico affrontato è stato la configurazione dell'offset (earliest vs latest). Ho configurato Logstash in modo da poter leggere dall'inizio (earliest) e utilizzare un group_id specifico, garantendo che nessun dato storico venisse perso durante i riavvii del sistema.

```
1 input {
2   kafka {
3     bootstrap_servers => "e-commerce_kafkaserver:9092"
4     topics => ["Topic_Dashboard"]
5     codec => "json"
6     auto_offset_reset => "earliest"
7     group_id => "logstash_v1"
8   }
9 }
10
11 filter {
12   mutate {
13     add_field => { "source_system" => "spark_mllib" }
14   }
15 }
16
17 output {
18   elasticsearch {
19     hosts => ["elasticsearch:9200"]
20     index => "ecommerce-predictions"
21     ssl => false
22     ssl_certificate_verification => false
23   }
24   stdout { codec => rubydebug }
25 }
```

Logstash preleva i dati completi e arricchiti dal Machine Learning e li inserisce ordinatamente nel database di Elasticsearch per poterli visualizzare.

Elasticsearch:

Elasticsearch è il motore di ricerca e analisi. A differenza dei database tradizionali (SQL), è un database NoSQL orientato ai documenti. Questo lo rende velocissimo nell'indicizzare i JSON provenienti da Spark e nel permettere interrogazioni complesse in tempo reale.

Kibana:

Kibana è l'interfaccia visuale che interroga Elasticsearch.

Abbiamo realizzato una Dashboard interattiva divisa in più “widget”:

- **CEO View:** Mostra cosa è successo (Fatturato totale, Prodotti più venduti, Profile Analysis etc.)



KPI Realizzati e Analizzati

I Key Performance Indicators implementati includono:

- **Total Revenue:** Somma in tempo reale degli incassi.
- **Top Selling Products:** Classifica dei prodotti più popolari.
- **Behavioral Segmentation:** Percentuale di acquisti impulsivi vs ragionati.

Riepilogo Finale e Spiegazione Use-Case

Il progetto dimostra come **un'architettura con tecnologie moderne** possa trasformare un *semplice log di acquisto* ("T-Shirt venduta a 20€") *in una visione strategica* ("L'utente Lorenzo è un cliente impulsivo che spende sopra la media"). **Questo permette al management di un qualsiasi E-Commerce di reagire istantaneamente alle tendenze di mercato e personalizzare l'offerta per massimizzare i guadagni e raggiungere i budget mensili prefissati.**

Comandi utilizzati spesso:

- Accensione Demo orchestrato:

```
0. rm -rf checkpoint_sql && rm -rf checkpoint_mllib  
1. docker-compose up db WordPress fluent-bit kafka kafka-ui topics sparksql sparkmllib  
2. docker-compose down db WordPress fluent-bit sparksql sparkmllib  
3. docker-compose up logstash elasticsearch kibana
```

- Accesso al container WordPress:

```
docker exec -it e-commerce_wordpress bash
```

- Installazione dell'editor di testo (Nano):

```
apt-get update && apt-get install nano -y
```

- Installazione di WP-CLI (Command Line Interface per WordPress):

Serve per creare utenti e gestire il sito da terminale.

```
curl -O https://raw.githubusercontent.com/wp-cli/builds/gh-pages/phar/wp-cli.phar  
chmod +x wp-cli.phar  
mv wp-cli.phar /usr/local/bin/wp
```

- Modifica del file `functions.php`:

```
cd /var/www/html/wp-content/themes/blocksy  
nano functions.php
```

- Creazione Utenti (Data Generation): *Comando per creare un singolo utente (es. Mario):*

```
wp user create mario mario@test.com --role=customer --user_pass="password123" --allow-root
```

- Verifica Ricezione Dati (Consumer da Terminale oppure direttamente kafka-ui):

Comando per "ascoltare" il topic `Topic_Purchase` direttamente dal container Kafka:

```
docker exec -it e-commerce_kafkaserver /opt/kafka/bin/kafka-console-consumer.sh -  
-bootstrap-server localhost:9092 --topic Topic_Purchase --from-beginning
```