# Spark Day 3 Assignment

_____

## 1. Pipeline 1 — RDBMS → Spark → Amazon  Keyspaces (Cassandra)

**Purpose:** Migrate and denormalize relational data from RDS MySQL to Amazon Keyspaces for analytics and high-performance queries.
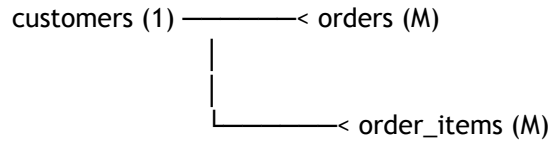
**Architecture:**

- Source: AWS RDS MySQL (3 normalized tables: customers, orders, order_items)
- Processing: Apache Spark (Scala) with JDBC and Cassandra Connector
- Target: Amazon Keyspaces table sales_data in keyspace spark_keyspace

**Database Schema:**

```
CREATE TABLE customers (
    customer_id INT UNSIGNED AUTO_INCREMENT PRIMARY KEY,
    name       VARCHAR(100) NOT NULL,
    email      VARCHAR(150) UNIQUE NOT NULL,
    city       VARCHAR(80)  NOT NULL
);

CREATE TABLE orders (
    order_id    INT UNSIGNED AUTO_INCREMENT PRIMARY KEY,
    customer_id INT UNSIGNED NOT NULL,
    order_date  DATE NOT NULL,
    amount      DECIMAL(12,2) NOT NULL,
    CONSTRAINT fk_orders_customer
       FOREIGN KEY (customer_id) REFERENCES customers(customer_id)
       ON DELETE RESTRICT ON UPDATE CASCADE
);

CREATE TABLE order_items (
    item_id       INT UNSIGNED AUTO_INCREMENT PRIMARY KEY,
    order_id      INT UNSIGNED NOT NULL,
    product_name  VARCHAR(200) NOT NULL,
    quantity      INT UNSIGNED NOT NULL DEFAULT 1,
    CONSTRAINT fk_order_items_order
       FOREIGN KEY (order_id) REFERENCES orders(order_id)
       ON DELETE CASCADE ON UPDATE CASCADE
);
```

**Table Relationship:**

```
customers (1) ───────< orders (M)
                │
                │
                └─────────< order_items (M)
```

**Data Flow:**

1. Extract: Read 3 tables from RDS MySQL via JDBC connector
2. Transform: Perform inner joins on customer_id and order_id to create denormalized view
3. Load: Write to Keyspaces using Spark Cassandra Connector

**Key Configuration:**

- SSL enabled for Keyspaces connection (port 9142)
- IAM service-specific credentials for authentication
- Truststore configured for SSL certificate validation

**Output Schema:**

- Partition Key: customer_id (int)
- Clustering Key: order_id (int)
- Additional columns: name, email, city, order_date, amount, item_id, product_name, quantity

## 2.Read Keyspaces → Write  Parquet to S3

**Purpose**: Extract denormalized sales data from Amazon Keyspaces and export to S3 as partitioned Parquet files for analytics and data lake storage.

**Architecture**:

- **Source**: Amazon Keyspaces table sales_data in keyspace spark_keyspace
- **Processing**: Apache Spark (Scala) with Cassandra Connector
- **Target**: AWS S3 bucket as partitioned Parquet files (s3://sparkdemo-bucket-1
- /sales/parquet/)

**Data Flow**:

1. **Extract**: Read sales_data from Keyspaces using Spark Cassandra Connector
2. **Transform**: Select subset of columns (customer_id, order_id, amount, product_name, quantity)
3. **Load**: Write to S3 as Parquet format, partitioned by customer_id

**Key Configuration**:

- SSL enabled for Keyspaces connection (port 9142)
- S3A file system with SimpleAWSCredentialsProvider for S3 access
- Extended timeouts for Keyspaces latency handling

## Output Structure:

- Format: Parquet (columnar, compressed)
- Partitioning: By customer_id for optimized query performance
- Location: s3a://sparkdemo-bucket-1/sales/parquet/customer_id=X/

## Directory Structure:

```
s3://sparkdemo-bucket-1/
│
├── sales/
│   └── parquet/
│       ├── customer_id=1/
│       │   ├── part-00000-xxx.snappy.parquet
│       │   └── part-00001-xxx.snappy.parquet
│       ├── customer_id=2/
│       │   └── part-00000-xxx.snappy.parquet
│       └── customer_id=3/
│           └── part-00000-xxx.snappy.parquet
```

## 3. Read Parquet → Aggregate → Write JSON

**Purpose**: Read partitioned Parquet sales data from S3, compute product-level aggregations, and export results as JSON for downstream reporting and API consumption.

**Architecture**:

- **Source**: S3 Parquet files from Pipeline 2 (s3://sparkdemo-bucket-1/sales/parquet/)
- **Processing**: Apache Spark (Scala) with DataFrame aggregations
- **Target**: S3 JSON output (s3://sparkdemo-bucket-1/aggregates/products.json)

**Data Flow**:

1. **Extract**: Read partitioned Parquet files from S3 (customer_id partitions)
2. **Transform**: Aggregate sales data by product_name
   - Sum quantity → total_quantity
   - Sum amount → total_revenue
3. **Load**: Write aggregated results as single JSON file to S3

**Aggregation Logic**:

- **Group By**: product_name
- **Metrics**:
  - total_quantity = SUM(quantity)
  - total_revenue = SUM(amount)
- **Sort**: Descending by total_revenue (top revenue products first)

**Output Format**:

json
{"product_name":"Widget A","total_quantity":150,"total_revenue":45000.0}

{"product_name":"Widget B","total_quantity":120,"total_revenue":36000.0}

**Key Configuration**:

- S3A file system with SimpleAWSCredentialsProvider
- Configuration values stored as variables at top of code for easy management
- .coalesce(1) to produce single output file

**Directory Structure:**

```
s3://sparkdemo-bucket-1/
│
├── aggregates/
│   └── products.json/
│       └── part-00000-xxx.json
```



## 4.Spark Structured Streaming → Check RDBMS for New Records → Write to Kafka (Avro)

**Purpose:** Stream newly inserted order records from MySQL into Kafka as Avro messages for real-time processing and event-driven workflows.

**Architecture:**

- **Source:** MySQL table new_orders
- **Processing:** Spark Structured Streaming (Scala), JDBC incremental scan
- **Target:** Kafka topic orders_avro_topic (Avro-encoded)

**Data Flow:**

- **Extract:** Poll MySQL every 5 seconds using incrementColumn = order_id

- **Transform:** Convert each new row into Avro (using orders.avsc)
- **Load:** Publish Avro bytes to Kafka with order_id as the message key

**Key Configuration:**

- MySQL JDBC with incremental reading
- Spark trigger interval: 5 seconds
- Avro GenericRecord serialization
- Kafka producer using ByteArraySerializer

**Avro Structure: orders.avsc**

```
{
  "type": "record",
  "name": "OrderRecord",
  "namespace": "com.retail",
  "fields": [
    { "name": "order_id",    "type": "int" },
    { "name": "customer_id", "type": "int" },
    { "name": "amount",      "type": "double" },
    { "name": "created_at",  "type": "string" }
  ]
}
```

**Output Structure:**

- **Format:** Avro binary
- **Topic:** orders_avro_topic
- **Key:** order_id
- **Value:** Avro payload containing order_id, customer_id, amount, created_at

```
Last login: Thu Dec  4 02:37:18 on ttys007
[racit@192 ~ % kafka-topics.sh --create --topic orders_avro_topic --bootstrap-server localhost:9092 --partitions 3 --replication-factor 1
WARNING: Due to limitations in metric names, topics with a period ('.') or underscore ('_') could collide. To avoid issues it is best to use either, but not both.
Created topic orders_avro_topic.
[racit@192 ~ % kafka-topics.sh --list --bootstrap-server localhost:9092
__consumer_offsets
allocation.events
maintenance.events
notifications
orders_avro_topic
overdue.events
people-topic
reservation-notifications
[racit@192 ~ % kafka-console-consumer.sh --bootstrap-server localhost:9092 --topic orders_avro_topic --from-beginning
?b@&2025-12-04 02:53:13
???(\V@&2025-12-04 02:53:13
@z@&2025-12-04 02:53:13
?b@&2025-12-04 02:57:00

???(\V@&2025-12-04 02:57:00

@z@&2025-12-04 02:57:00
```

# 5. Kafka Consumer Stream → Write JSON to S3

**Purpose**: Consume Avro messages from Kafka in real-time, decode them, and stream as JSON to S3.

**Architecture:**

- **Source**: Kafka topic orders_avro_topic (Avro binary from Pipeline 4)
- **Processing**: Spark Structured Streaming with Avro deserialization
- **Target**: S3 JSON files (s3://sparkdemo-bucket-1/stream/json/)
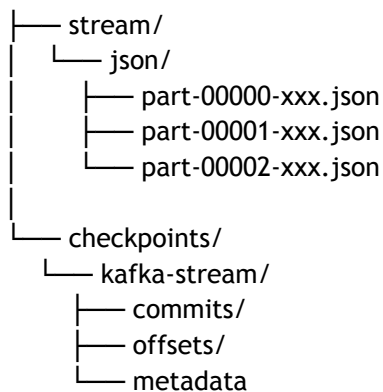
**Data Flow:**

1. **Consume**: Read from Kafka topic in real-time
2. **Decode**: Deserialize Avro binary using schema (order_id, customer_id, amount, created_at)
3. **Load**: Write JSON to S3 every 10 seconds

**Key Configuration:**

- Starting offsets: earliest (reads all messages)
- Trigger: 10-second micro-batches
- Checkpoint: /tmp/kafka-s3-checkpoint (fault tolerance)
- Output mode: Append

**Directory Structure:**

```
s3://sparkdemo-bucket-1/

├── stream/
│   └── json/
│       ├── part-00000-xxx.json
│       ├── part-00001-xxx.json
│       └── part-00002-xxx.json
│
└── checkpoints/
    └── kafka-stream/
        ├── commits/
        ├── offsets/
        └── metadata
```

# Output Format:

```
{} part-00000-029cd60b-b66d-4d89-af9f-4f829aa71966-c000.json 1 ×

Users > racit > Downloads > {} part-00000-029cd60b-b66d-4d89-af9f-4f829aa71966-c000.json > ...
1  {"order_id":13,"customer_id":1,"amount":150.75,"created_at":"2025-12-04 11:03:48","kafka_timestamp":"2025-12-04T11:07:09.915+05:30"}
2  {"order_id":14,"customer_id":2,"amount":89.99,"created_at":"2025-12-04 11:03:48","kafka_timestamp":"2025-12-04T11:07:09.922+05:30"}
3  {"order_id":15,"customer_id":3,"amount":420.0,"created_at":"2025-12-04 11:03:48","kafka_timestamp":"2025-12-04T11:07:09.922+05:30"}
4
```

Amazon S3 > Buckets > sparkdemo-bucket-1 > stream/ > json/

## Amazon S3

- ▼ Buckets
  - General purpose buckets
  - Directory buckets
  - Table buckets
  - Vector buckets  New
- ▼ Access management and security
  - Access Points
  - Access Points for FSx
  - Access Grants
  - IAM Access Analyzer
- ▼ Storage management and insights
  - Storage Lens
  - Batch Operations
- Account and organization settings
- ▶ AWS Marketplace for S3

## json/

Objects | Properties

### Objects (4)

Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more

| | Name | Type | Last modified | Size | Storage class |
|---|---|---|---|---|---|
| | _spark_metadata/ | Folder | - | - | - |
| | part-00000-029cd60b-b66d-4d89-af9f-4f829aa71966-c000.json | json | December 4, 2025, 11:07:20 (UTC+05:30) | 397.0 B | Standard |
| | part-00000-7d6f3009-b500-48c3-bf15-c72db67249a1-c000.json | json | December 4, 2025, 11:08:04 (UTC+05:30) | 397.0 B | Standard |
| | part-00000-eae4c6eb-a9f5-400d-b3f9-f9877dd2050d-c000.json | json | December 4, 2025, 11:07:09 (UTC+05:30) | 397.0 B | Standard |