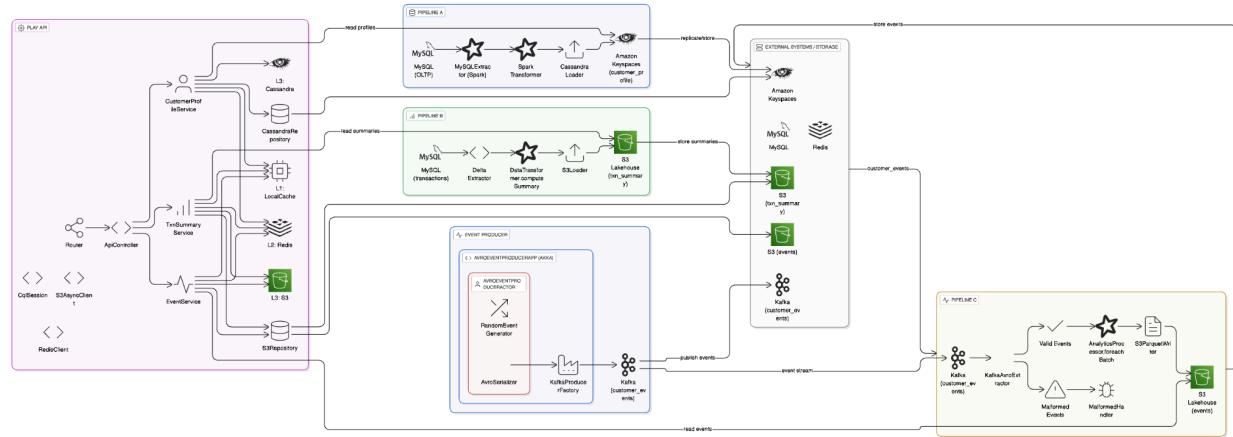


# Customer Data Engineering

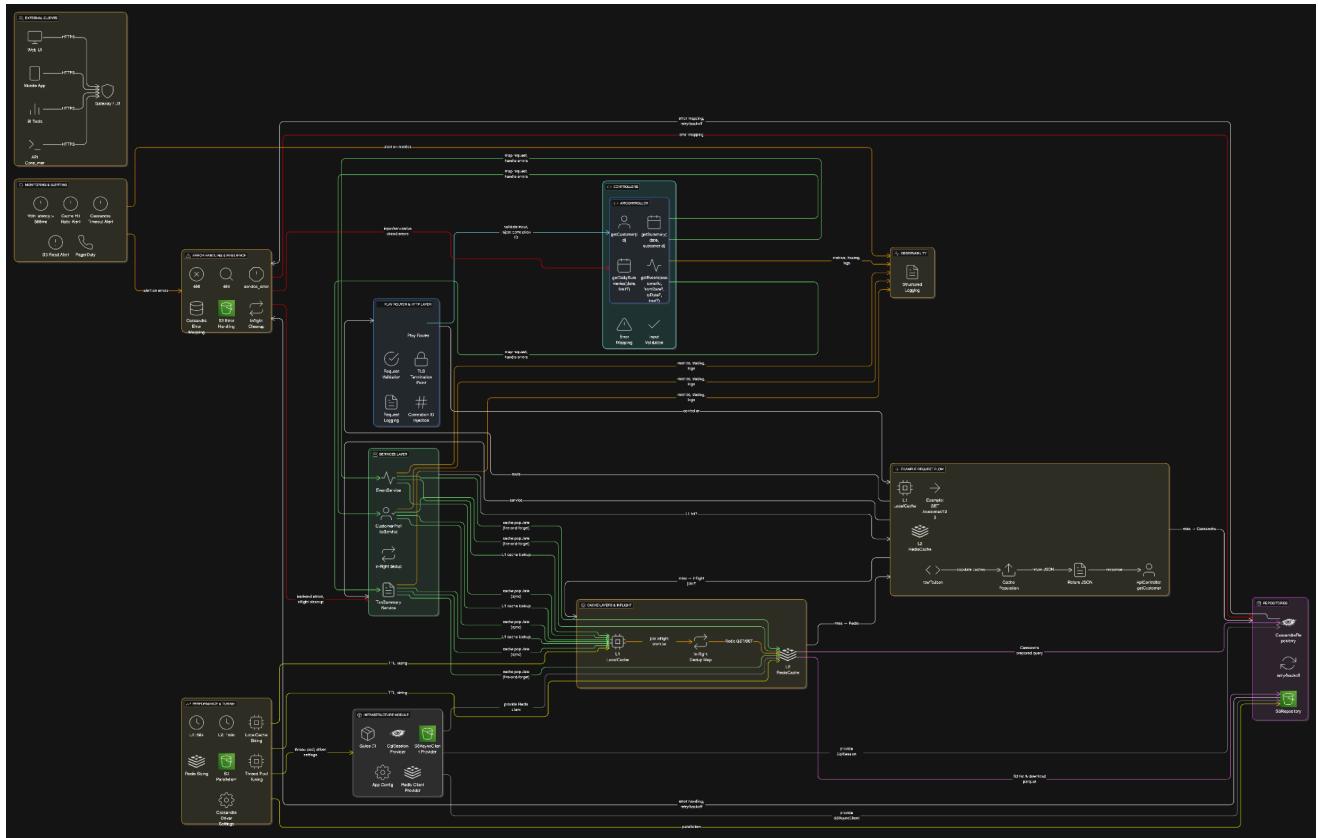
Owner: Sanjeev Kumar V

Date: Dec 12, 2025

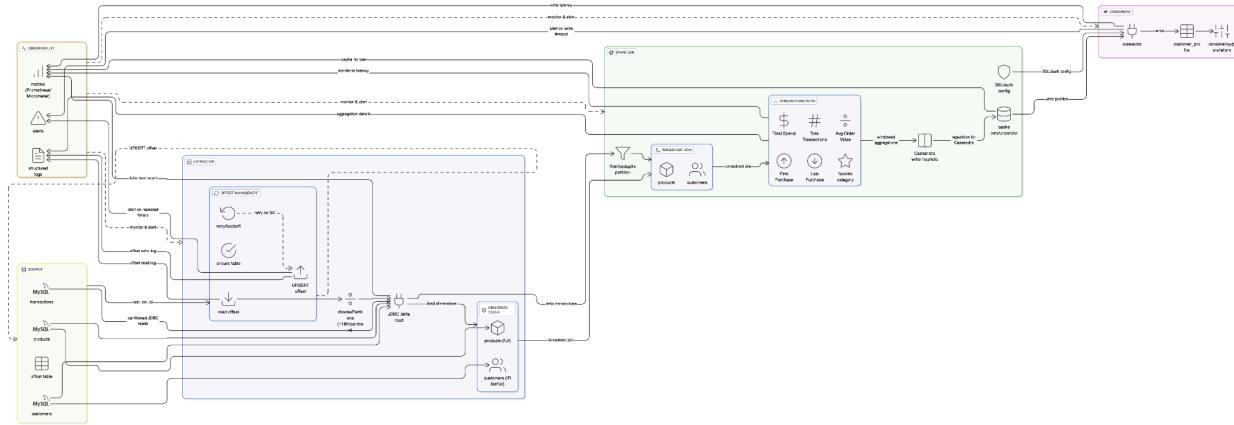
## System Architecture:



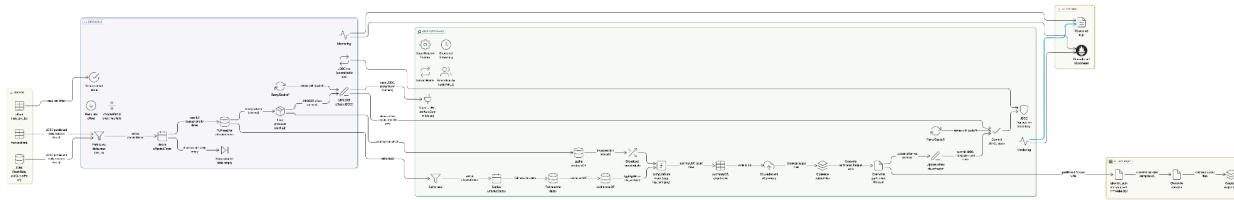
## Play API Architecture



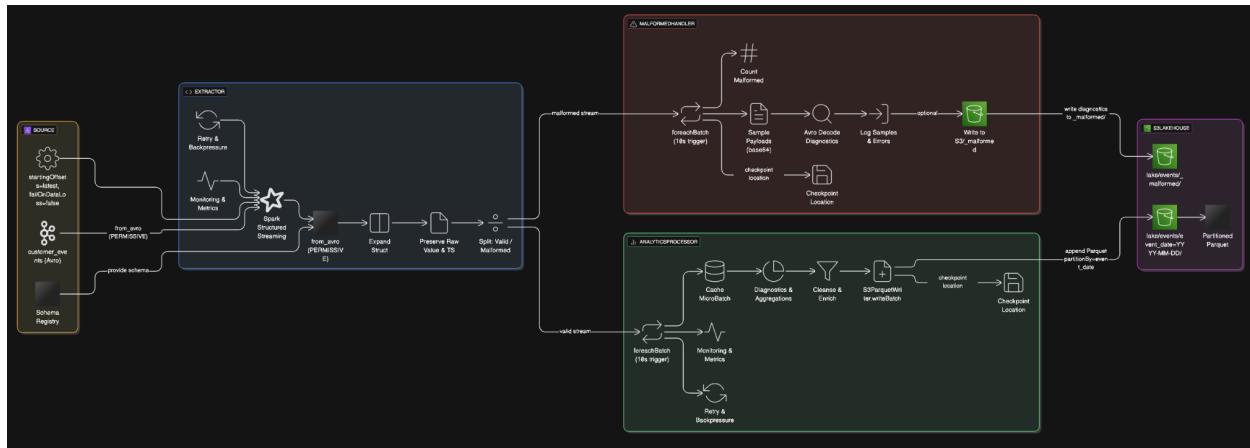
## Pipeline A Architecture:



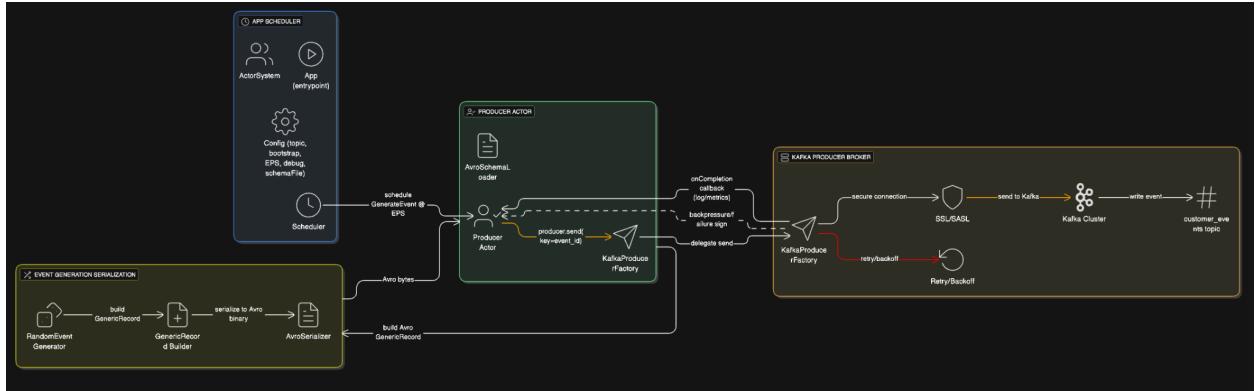
## Pipeline B Architecture:



## Pipeline C Architecture:



## Producer Architecture:



## DATA DICTIONARY

### MySQL - Source System

---

**Table: customers**

This table stores all customer master data.

Field	Type	Description
<b>customer_id</b>	INT (PK)	Unique customer identifier
<b>name</b>	VARCHAR(100)	Customer full name
<b>email</b>	VARCHAR(150), UNIQUE	Customer email (unique across DB)
<b>gender</b>	ENUM('M','F','O')	Gender of the customer
<b>signup_date</b>	DATE	Date when customer account was created

**Table: products**

Stores details of all products purchased in transactions.

Field	Type	Description
<b>product_id</b>	INT (PK)	Unique product identifier
<b>name</b>	VARCHAR(150)	Product name

<b>category</b>	VARCHAR(50)	Product category (Electronics, Grocery, Sports...)
<b>price</b>	DECIMAL(10,2)	Unit price of the product

#### Table: **transactions**

This is the core **fact table** for Pipelines A & B (profiles + daily summaries).

Field	Type	Description
<b>txn_id</b>	BIGINT (PK)	Unique transaction ID
<b>customer_id</b>	INT (FK)	FK → customers.customer_id
<b>product_id</b>	INT (FK)	FK → products.product_id
<b>qty</b>	INT	Number of items purchased
<b>amount</b>	DECIMAL(10,2)	Total amount of transaction (qty × price)
<b>txn_timestamp</b>	TIMESTAMP	Purchase timestamp

#### Offset Tracking Tables (Used by Pipelines A & B)

#### Table: **profile\_offsets**

Used by **Pipeline A (MySQL → Cassandra)** to ensure incremental processing.

Field	Type	Description
<b>source</b>	VARCHAR(128)	Identifier for the pipeline source (e.g., "transactions")
<b>last_txn_id</b>	BIGINT	Last processed txn_id for profile updates

#### Table: **txn\_summary\_offsets**

Used by **Pipeline B (MySQL → S3 Daily Summaries)** for incremental loads.

Field	Type	Description

<b>source</b>	VARCHAR(128)	Identifier for the pipeline source
<b>last_txnid</b>	BIGINT	Last processed txnid for summary generation

---

### Kafka Event Schema (Pipeline C Input)

(Event Producer → Kafka → Pipeline C → S3)

Field	Type	Description
event_id	STRING	Unique event UUID
customer_id	INT	Customer who performed the event
event_type	STRING	Behaviour type: WISHLIST, LIKE, CART_ADD.
product_id	INT NULL	Optional product reference
event_timestamp	STRING	ISO timestamp assigned at event creation

---

### Cassandra Table (Pipeline A Output)

#### Keyspaces Table: `customer_profile`

Stores aggregated customer analytics.

Column	Type	Description
customer_id	INT (PK)	Unique ID
name	TEXT	Customer name
email	TEXT	Email
gender	TEXT	Gender
total_spend	DECIMAL	SUM(amount)
total_transactions	INT	COUNT(*)
avg_order_value	DECIMAL	total_spend / total_transactions
first_purchase	TIMESTAMP	MIN(txn_timestamp)

last_purchase	TIMESTAMP	MAX(txn_timestamp)
favorite_category	TEXT	Category with highest frequency

---

### S3 Daily Summary Schema (Pipeline B Output)

Partition folder: `lake/txn_summary/date=YYYY-MM-DD/*.parquet`

Column	Type	Description
date	STRING	Partition key
customer_id	INT	Customer ID
total_amount	DOUBLE	Daily spend
total_items	INT	Total quantity
distinct_products	INT	Count of unique products
top_category	STRING	Most purchased category

---

### 2.6 S3 Events Schema (Pipeline C Output)

Partition folder: `lake/events/event_date=YYYY-MM-DD/*.parquet`

Column	Type	Description
event_id	STRING	Event identifier
customer_id	INT	Customer
event_type	STRING	Behaviour
product_id	INT	Nullable
event_timestamp	BINARY (INT96)	Original event timestamp
ingestion_timestamp	BINARY (INT96)	Ingestion time recorded by Spark

# Pipeline-quality checks report

## Pipeline A – MySQL → Cassandra (Customer Profile ETL)

### 1. Input Data Quality Checks

Check Type	Description
Primary Key Integrity	Ensures customer_id is not null and unique.
Email Format Validation	Rejects invalid or malformed email addresses.
Gender Enum Validation	Only accepts M, F, O.
Date Validation	Ensures signup_date is a valid ISO date.

---

### 2. Schema & Transformation Quality Checks

Check	Description
Column completeness	Ensures all required fields exist before transformation.
Type consistency	Converts MySQL types → Cassandra compatible types (e.g., VARCHAR → VARCHAR, DATE → DATE).
Null handling	Replaces unexpected nulls with defaults/logging warnings.

---

### 3. Offset / Incremental Load Checks

Check	Description
Offset table validation (profile_offsets)	Reads last processed txn_id before each run.
Monotonicity check	Ensures new txn_id > stored offset.
Duplicate prevention	Avoids re-processing already consumed rows.

---

#### 4. Write Quality Checks (Cassandra)

Check	Description
Consistency level validation	Ensures write acknowledgements success.
Timeout handling	Retries on Cassandra ReadTimeout, WriteTimeout, Unavailable errors.
Row-level validation after write	Optional: read-back check (configurable).

---

### Pipeline B – MySQL → S3 Lakehouse (Daily Transaction Summary)

#### 1. Input Data Quality Checks

Check Type	Description
Primary key validation	Ensures <code>txn_id</code> exists and is unique.
Foreign key integrity	Valid <code>customer_id</code> and <code>product_id</code> reference.
Amount and Quantity check	Must be positive numeric values.
Timestamp validity	Ensures <code>txn_timestamp</code> is a valid UTC timestamp.

---

#### 2. Aggregation Quality Checks

Check	Description
Group-by consistency	Ensures summaries are calculated correctly per customer per day.
Total amounts validation	No negative totals; amounts match $\text{sum}(\text{qty} * \text{price})$ .

Distinct product count check	Ensures distinct count uses correct product_id.
------------------------------	---

### 3. Offset Checks (`txn_summary_offsets`)

Check	Description
Last processed txn_id validation	Ensures incremental summary updates.
No missing transactions	Detects gaps in txn_id sequence.

### 4. Write Quality Checks (S3 Parquet)

Check	Description
Partition validation	Must write to: lake/txn_summary/date=YYYY-MM-DD/
Schema consistency	Enforces Parquet schema: {customer_id, total_amount, total_items, distinct_products, top_category}
Atomic batch writes	Uses temporary files + rename for safe writes.
File-level validation	Validates Parquet file is readable after write.

---

## Pipeline C – Kafka Avro → S3 Events Lake (Valid + Malformed Streams)

### 1. Kafka Input Quality Checks

Check	Description
Avro schema validation	Deserializes using schema registry with PERMISSIVE mode.
Required field checks	Ensures event_id & event_type must exist.
Event timestamp validation	Parses event timestamp or falls back to Kafka timestamp.

Malformed event detection	Missing keys → automatically routed to <code>malformed</code> stream.
---------------------------	---

## 2. Transformation Quality Checks

Check	Description
Product ID correction	Invalid product_id (<0) → null.
Timestamp normalization	Converts string timestamp → Spark timestamp.
Partition column generation	Adds <code>event_date</code> for partitioning.
Lineage metadata	Adds <code>ingestion_timestamp</code> .

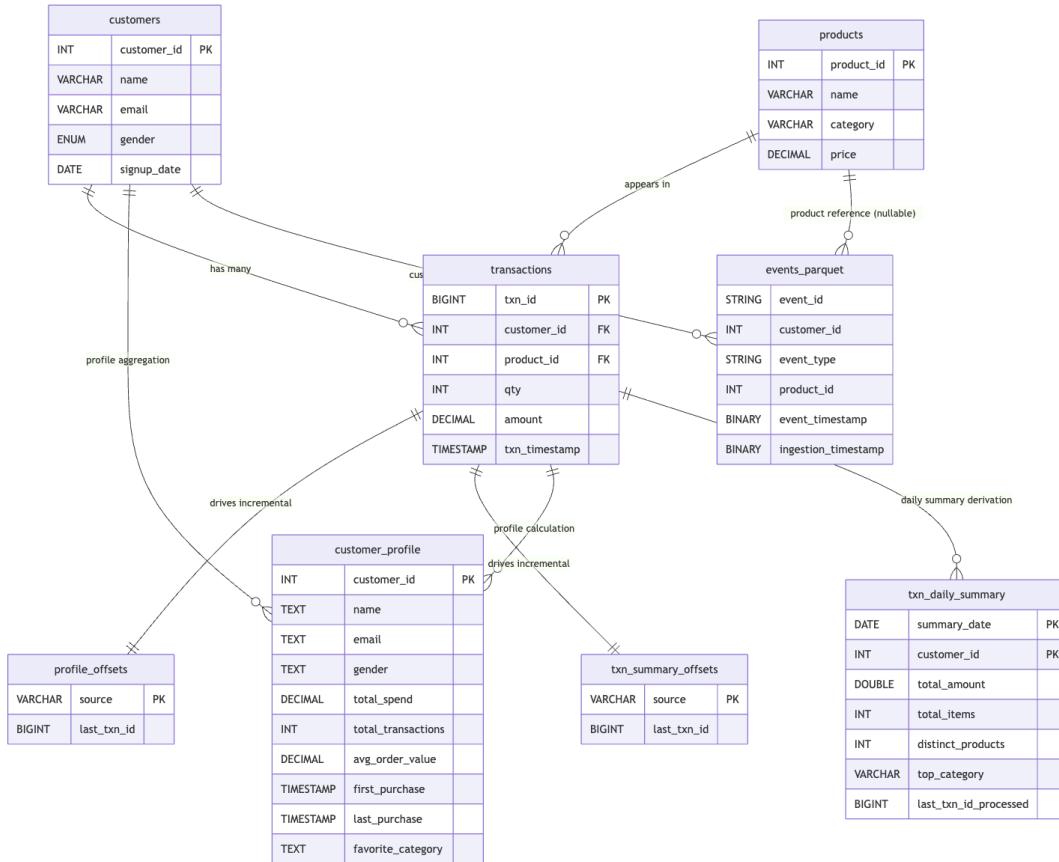
## 3. Malformed Events Checks

Check	Description
Sample logging	Logs up to 5 Base64 samples.
Avro decode attempt	Logs decoding failures for analysis.
Count mismatch detection	Alerts if malformed ratio exceeds threshold.

## 4. Write Quality Checks (S3 Parquet)

Check	Description
Partitioning validation	Writes to: <code>lake/events/event_date=YYYY-MM-DD/</code>
Schema consistency	Ensures event schema is stable across writes.
Exactly-once semantics	Achieved via streaming checkpoints.

Fail-safe writes	Writes via <code>foreachBatch</code> with exception handling.
------------------	---



## SQL / DataFrame

Completed Queries: 7

### Completed Queries (7)

Page:

1 Pages. Jump to  . Show  items in a page.

ID	Description	Submitted	Duration	Job IDs	Sub Execution IDs
16	<a href="#">id = 4b78e16e-d9a3-495d-98a4-45bbdb113ba8 runId = a12a1d35-d35f-4a96-aa42-96e4...</a> +details	2025/12/12 16:18:00	3 s		
15	<a href="#">id = 4b78e16e-d9a3-495d-98a4-45bbdb113ba8 runId = a12a1d35-d35f-4a96-aa42-96e4...</a> +details	2025/12/12 16:17:00	4 s		
14	<a href="#">id = 4b78e16e-d9a3-495d-98a4-45bbdb113ba8 runId = a12a1d35-d35f-4a96-aa42-96e4...</a> +details	2025/12/12 16:16:00	3 s		
8	<a href="#">id = 4b78e16e-d9a3-495d-98a4-45bbdb113ba8 runId = a12a1d35-d35f-4a96-aa42-96e4...</a> +details	2025/12/12 16:15:00	27 s	[33][34][35][36][37][38][39][40][41][42][43][44][45][46]	[9][10][11][12][13] +details
7	<a href="#">id = 4b78e16e-d9a3-495d-98a4-45bbdb113ba8 runId = a12a1d35-d35f-4a96-aa42-96e4...</a> +details	2025/12/12 16:14:00	3 s		
6	<a href="#">id = 4b78e16e-d9a3-495d-98a4-45bbdb113ba8 runId = a12a1d35-d35f-4a96-aa42-96e4...</a> +details	2025/12/12 16:13:11	3 s		
0	<a href="#">id = 4b78e16e-d9a3-495d-98a4-45bbdb113ba8 runId = a12a1d35-d35f-4a96-aa42-96e4...</a> +details	2025/12/12 16:12:30	41 s	[7][8][9][10][11][12][13][14][15][16][17][18][19][20][21][22]	[1][2][3][4][5] +details

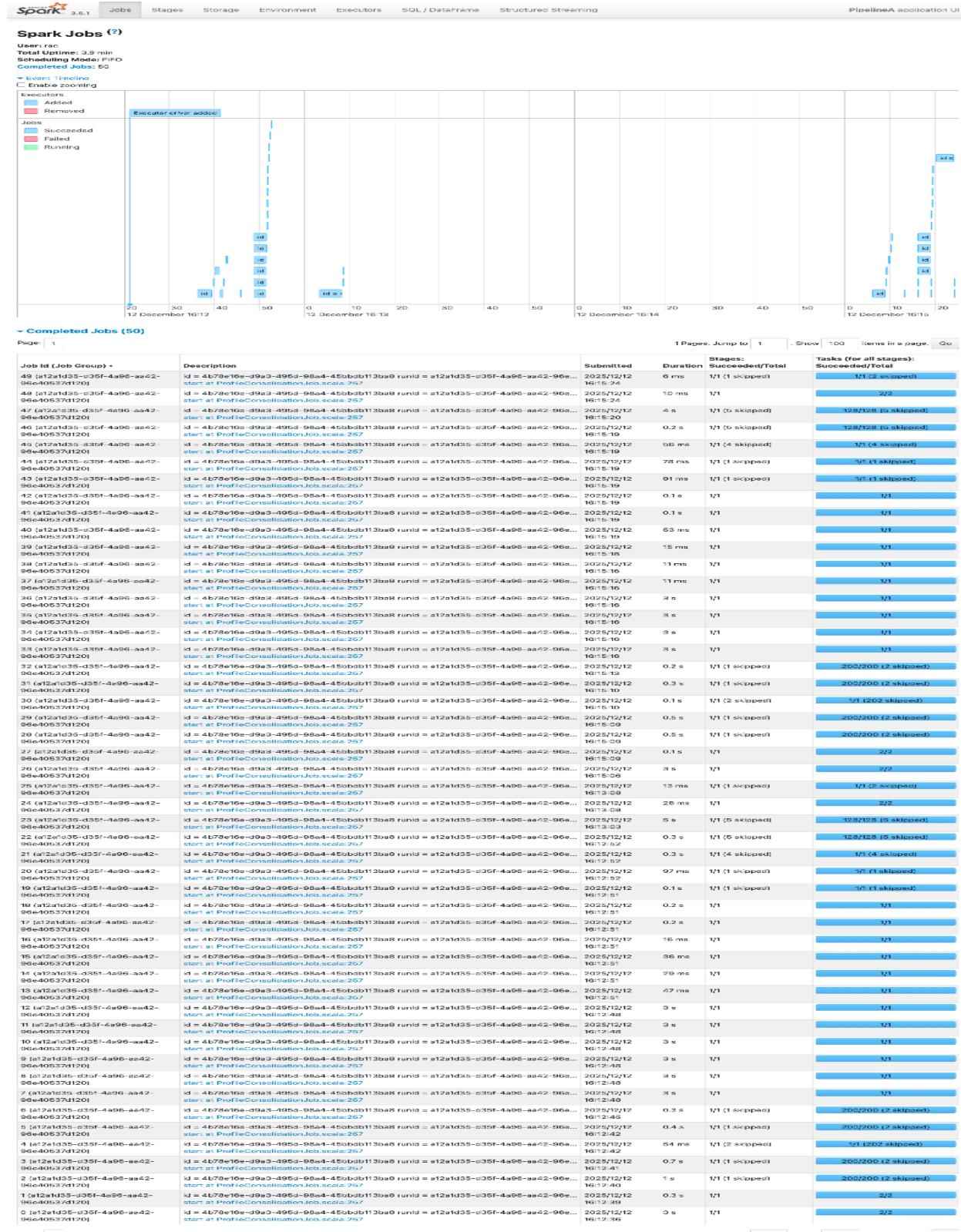
Page:

1 Pages. Jump to  . Show  items in a page.

## Storage

### RDDs

ID	RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory	Size on Disk
44	(*1) Scan JDBCRelation(products) [numPartitions=1] [product_id#0,category#2] PushedFilters: Disk Memory Deserialized 1x Replicated [], ReadSchema: struct<product_id:int,category:string>	Disk Memory Deserialized 1x Replicated	1	100.00%	2.0 KiB	0.0 B



Stages for All Jobs										PipelineA application UI			
Completed Stages: 50		Skipped Stages: 40		Completed Stages (50)		1 Pages Jump to 1 Show 100 Items in a page Go							
Page: 1	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write						
Stage Id	Description												
95	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:24	3 ms	1/1								118.0 kB	
93	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:24	8 ms	2/2		6.7 kB						118.0 kB	
92	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:20	4 ms	128/128		167.1 kB						26.7 kB	
95	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:19	0.1 ms	128/128								26.7 kB	
93	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:19	51 ms	1/1								83.2 kB	24.7 kB
76	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:19	26 ms	1/1								83.6 kB	26.9 kB
95	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:19	86 ms	1/1								98.4 kB	56.3 kB
71	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:19	0.1 ms	1/1		64.6 kB						63.8 kB	
70	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:19	0.1 ms	1/1		64.6 kB						98.4 kB	
93	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:19	50 ms	1/1								10.5 kB	
98	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:16	19 ms	1/1								10.3 kB	
97	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:16	8 ms	1/1								2.0 kB	
88	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:16	10 ms	1/1								2.0 kB	
55	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:16	3 ms	1/1								10.3 kB	
84	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:16	3 ms	1/1								64.6 kB	
63	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:16	3 ms	1/1								9.5 kB	
82	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:16	3 ms	1/1								11.3 kB	
81	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:13	0.2 ms	200/200		32.5 kB						200.200	
99	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:10	0.8 ms	200/200		92.0 kB						92.0 kB	
57	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:10	0.1 ms	1/1								11.3 kB	
54	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:09	0.6 ms	200/200		82.0 kB						82.0 kB	11.3 kB
52	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:09	0.5 ms	200/200								9.5 kB	
50	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:09	0.1 ms	2/2		6.7 kB						9.5 kB	
49	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:08	3 ms	2/2								118.0 kB	
48	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:08	7 ms	1/1								118.0 kB	
48	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:08	19 ms	2/2		6.5 kB						118.0 kB	
45	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:15:08	5 ms	128/128		185.7 kB						185.7 kB	
39	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:02	0.2 ms	128/128								24.7 kB	
33	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:02	0.3 ms	1/1								81.2 kB	24.7 kB
28	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:02	89 ms	1/1								96.7 kB	24.9 kB
26	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	0.1 ms	1/1								61.8 kB	26.3 kB
24	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	0.2 ms	128/128		62.6 kB						85.7 kB	
28	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	0.2 ms	1/1								62.6 kB	81.8 kB
22	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	10 ms	1/1								10.2 kB	
21	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	30 ms	1/1								10.2 kB	
20	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	22 ms	1/1								2.0 kB	
19	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	41 ms	1/1								2.0 kB	
18	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	3 ms	1/1								2.0 kB	
17	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	3 ms	1/1								18.2 kB	
16	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	3 ms	1/1								62.6 kB	
15	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	3 ms	1/1								18.2 kB	
14	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	3 ms	1/1								18.2 kB	
13	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	3 ms	1/1								18.2 kB	
12	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	0.2 ms	200/200		82.0 kB						82.0 kB	
10	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	0.4 ms	200/200		82.0 kB						82.0 kB	
9	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	51 ms	1/1								11.3 kB	
5	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	0.7 ms	200/200		82.0 kB						11.3 kB	
3	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	1 ms	200/200								8.3 kB	
1	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	0.8 ms	2/2		6.5 kB						9.3 kB	
0	Id = 4b78e19c-d9b5-495d-98a4-45bbdb13ba8 runid = e12a1d35-d35f-4a98-aef2-96ea.. +details	2025/12/12 16:12:01	3 ms	2/2									

Apache Spark 3.5.1 Jobs Stages Storage Environment Executors SQL / DataFrame Structured Streaming Pipeline B - Incremental Daily T... application UI

**Spark Jobs (39)**

User: rac  
Total Uptime: 4.9 min  
Scheduling Mode: FIFO  
Completed Jobs: 39

Event Timeline  
Enable zooming

Executors  
Added  
Removed

Jobs  
Succeeded  
Failed  
Running

Executor driver added

12 December 16:24 12 December 16:25 12 December 16:26 12 December 16:27

Completed Jobs (39)

Page: 1 1 Pages, Jump to 1 . Show 100 items in a page. Go

Job Id (Job Group) *	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/total
38 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:28:07	3 s	1/1	1/1
37 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:27:06	3 s	1/1	1/1
36 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:26:07	3 s	1/1	1/1
35 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:25:37	2 s	1/1	1/1
34 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:34	15 s	1/1 (2 skipped)	1/1 (201 skipped)
33 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:26	28 ms	1/1 (2 skipped)	1/1 (201 skipped)
32 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:23	2 s	1/2	200/201
31 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:23	1 s	1/1 (1 skipped)	200/200 (1 skipped)
30 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:23	22 ms	1/1	1/1
29 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:18	5 ms	1/1 (3 skipped)	1/1 (202 skipped)
28 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:18	33 ms	1/1 (2 skipped)	1/1 (201 skipped)
27 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:18	34 ms	1/1 (2 skipped)	1/1 (201 skipped)
26 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:15	2 s	1/2	200/201
25 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:15	1 s	1/1 (1 skipped)	200/200 (1 skipped)
24 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:15	15 ms	1/1	1/1
23 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:15	10 ms	1/1 (3 skipped)	1/1 (202 skipped)
22 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:15	41 ms	1/1 (2 skipped)	1/1 (201 skipped)
21 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:15	60 ms	1/1 (2 skipped)	1/1 (201 skipped)
20 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:13	2 s	1/2	200/201
19 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:12	1 s	1/1 (1 skipped)	200/200 (1 skipped)
18 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:12	0.5 s	1/1	1/1
17 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	0.7 s	1/1 (1 skipped)	200/200 (1 skipped)
16 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	0.5 s	1/1 (1 skipped)	200/200 (1 skipped)
15 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	0.1 s	1/1	1/1
14 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	9 ms	1/1 (1 skipped)	1/1 (1 skipped)
13 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	13 ms	1/1	1/1
12 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	11 ms	1/1 (1 skipped)	1/1 (1 skipped)
11 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	19 ms	1/1	1/1
10 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	12 ms	1/1 (1 skipped)	1/1 (1 skipped)
9 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:11	15 ms	1/1	1/1
8 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:08	3 s	1/1	1/1
7 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:06	16 ms	1/1 (1 skipped)	1/1 (1 skipped)
6 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:06	44 ms	1/1	1/1
5 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:03	3 s	1/1	1/1
4 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:24:00	32 ms	1/1 (1 skipped)	1/1 (1 skipped)
3 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:23:57	3 s	1/1	1/1
2 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:23:54	82 ms	1/1 (1 skipped)	1/1 (1 skipped)
1 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:23:54	3 s	1/1	1/1
0 (450b8ce8-43a3-4712-b445-3d0bdffdbaa06)	Id = 8c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdffdbaa06 start at PipelineBJob.scala:187	2025/12/12 16:23:50	3 s	1/1	1/1

## Stages for All Jobs

Completed Stages: 40

Skipped Stages: 27

### - Completed Stages (40)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
66	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:29:06	3 s	1/1				
65	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:28:07	3 s	1/1				
64	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:27:06	3 s	1/1				
63	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:26:07	3 s	1/1				
82	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:25:37	2 s	1/1				
81	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:34	16 s	1/1		19.0 KiB	111.2 KiB	
58	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:26	26 ms	1/1			97.8 KiB	
55	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:23	1 s	200/200	203.0 KiB			97.8 KiB
54	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:23	1 s	200/200	203.0 KiB			111.2 KiB
52	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:23	17 ms	1/1	13.5 KiB			
51	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:18	2 ms	1/1			59.0 B	
47	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:18	30 ms	1/1			78.7 KiB	59.0 B
44	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:18	31 ms	1/1			78.7 KiB	
41	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:15	1 s	200/200	203.0 KiB			78.7 KiB
40	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:15	1 s	200/200	203.0 KiB			78.7 KiB
38	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:15	11 ms	1/1	13.5 KiB			
37	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:15	7 ms	1/1			59.0 B	
33	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:15	38 ms	1/1			78.7 KiB	59.0 B
30	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:15	56 ms	1/1			78.7 KiB	
27	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:13	1 s	200/200	203.0 KiB			78.7 KiB
26	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:12	1 s	200/200	203.0 KiB			78.7 KiB
25	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:12	30 ms	1/1	13.5 KiB			
24	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	0.2 s	200/200	203.0 KiB			
23	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	0.5 s	200/200			58.5 KiB	
21	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	0.1 s	1/1	35.7 KiB			58.5 KiB
20	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	4 ms	1/1			59.0 B	
18	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	9 ms	1/1	13.5 KiB			59.0 B
17	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	6 ms	1/1			59.0 B	
15	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	12 ms	1/1	35.7 KiB			59.0 B
14	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	6 ms	1/1			59.0 B	
12	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:11	9 ms	1/1	13.5 KiB			59.0 B
11	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:08	3 s	1/1				
10	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:06	10 ms	1/1			59.0 B	
8	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:06	33 ms	1/1	35.7 KiB			59.0 B
7	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:03	3 s	1/1				
6	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:24:00	20 ms	1/1			59.0 B	
4	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:23:57	3 s	1/1				59.0 B
3	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:23:57	60 ms	1/1			177.0 B	
1	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:23:54	3 s	1/1				177.0 B
0	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0... start at PipelineJob.scala:187	2025/12/12 16:23:50	3 s	1/1				

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

### - Skipped Stages (27)

Done 1 Page 1 Items to 1 Show 100 Items in a page Go

## Storage

### ▼ RDDs

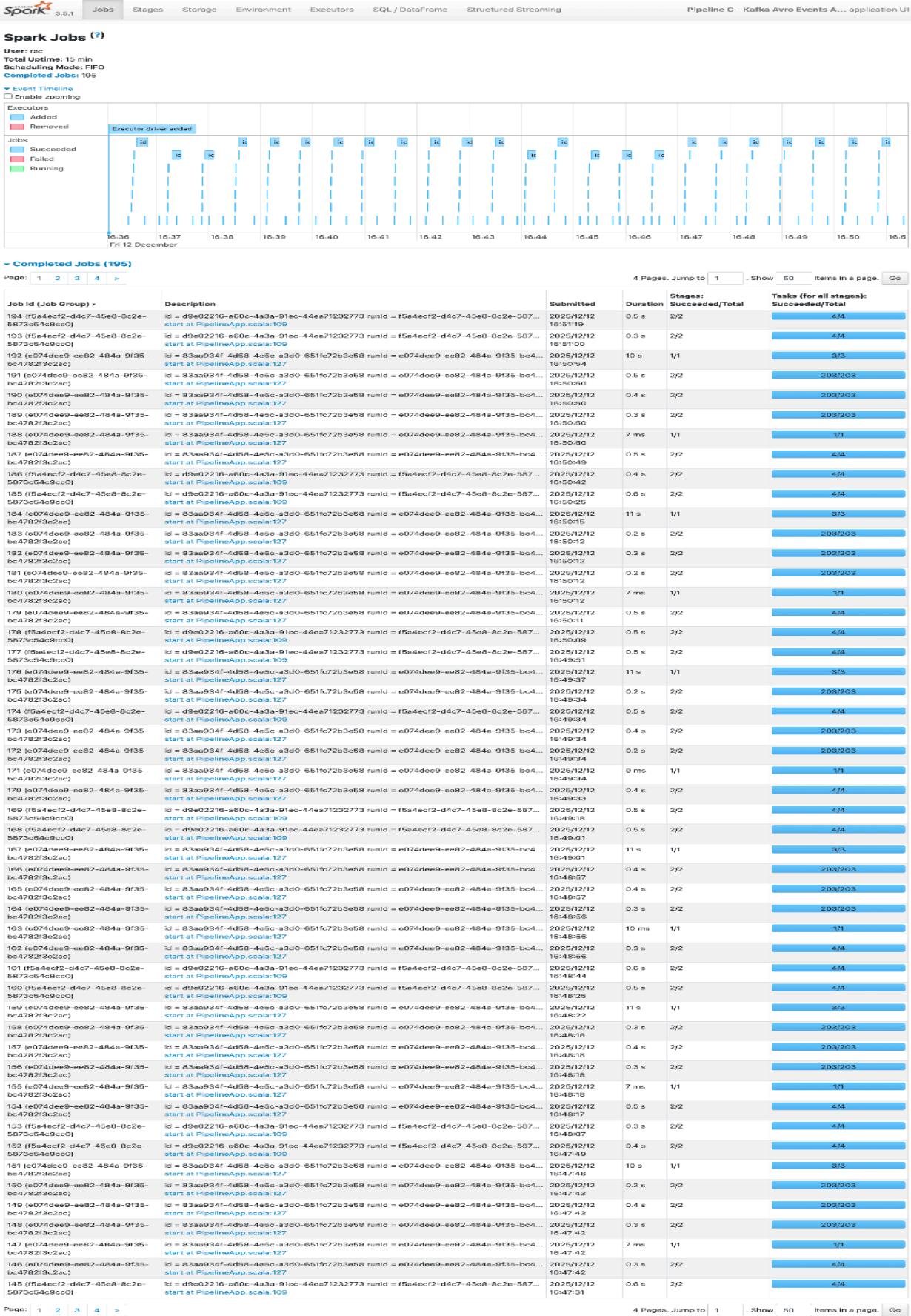
ID	RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory	Size on Disk
22	*(1) Project [txn_id#55L, customer_id#56, product_id#57, qty#58, amount#59, txn_timestamp#60, cast(txn_timestamp#60 as date) AS date#67] ++ *(1) Scan JDBCRelation((SELECT txn_id, customer_id, product_id, qty, amount, txn_timestamp FROM transactions WHERE DATE(txn_timestamp) IN ('2025-12-11','2025-12-12','2025-12-10')) AS tx_full) [numPartitions=1] [txn_id#55L,customer_id#56,product_id#57,qty#58,amount#59,txn_timestamp#60] PushedFilters: [], ReadSchema: struct<txn_id:bigring, customer_id:int, product_id:int, qty:int, amount:decimal(10,2), txn_timestamp :t1...]	Disk Memory Deserialized 1x Replicated	1	100.00%	35.7 KiB	0.0 B
34	*(1) Scan JDBCRelation(products) [numPartitions=1] [product_id#297,name#298,category#299,price#300] PushedFilters: [], ReadSchema: struct<product_id:int,name:string,category:string,price:decimal(10,2)>	Disk Memory Deserialized 1x Replicated	1	100.00%	13.5 KiB	0.0 B
66	AdaptiveSparkPlan isFinalPlan=false +- Exchange hashpartitioning(customer_id#56, 200), REPARTITION_BY_COL, [plan_id=203] +- Project [txn_id#55L, customer_id#56, product_id#57, qty#58, amount#59, txn_timestamp#60, cast(txn_timestamp#60 as date) AS date#434] +- Filter isnotnull(product_id#57) +- InMemoryTableScan [amount#59, customer_id#56, product_id#57, qty#58, txn_id#55L, txm_timestamp#60], [isnotnull(product_id#57)] +- InMemoryRelation [txn_id#55L, customer_id#56, product_id#57, qty#58, amount#59, txm_timestamp#60, date#67], StorageLevel(disk, memory, deserialized, 1 replicas) +- *(1) Project [txn_id#55L, customer_id#56, product_id#57, qty#58, amount#59, txm_timestamp#60, cast(txn_timestamp#60 as date) AS date#67] +- *(1) Scan JDBCRelation((SELECT txn_id, customer_id, product_id, qty, amount, txn_timestamp FROM transactions WHERE DATE(txn_timestamp) IN ('2025-12-11','2025-12-12','2025-12-10')) AS tx_full) [numPartitions=1] [txn...]	Disk Memory Deserialized 1x Replicated	200	100.00%	203.0 KiB	0.0 B

## SQL / DataFrame

Completed Queries: 7

### ▼ Completed Queries (7)

Page:	1	1 Pages. Jump to	1	. Show	100	items in a page.	Go
ID	Description	Submitted	Duration	Job IDs	Sub Execution IDs		
21	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdfdbaa06 batch = 6 +details	2025/12/12 16:30:00	10 s	[22]		+details	
19	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdfdbaa06 batch = 5 +details	2025/12/12 16:29:00	9 s	[20]		+details	
17	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdfdbaa06 batch = 4 +details	2025/12/12 16:28:00	11 s	[18]		+details	
15	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdfdbaa06 batch = 3 +details	2025/12/12 16:27:00	9 s	[16]		+details	
13	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdfdbaa06 batch = 2 +details	2025/12/12 16:26:00	10 s	[14]		+details	
11	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdfdbaa06 batch = 1 +details	2025/12/12 16:25:31	9 s	[12]		+details	
0	id = 6c4dc180-b837-41fc-9890-d0e75ad52e55 runid = 450b8ce8-43a3-4712-b445-3d0bdfdbaa06 batch = 0 +details	2025/12/12 16:23:44	1.8 min	[1][2][3][4][5][6][7][8][9][10]		+details	



## Storage

### • RDDs

ID	RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory	Size on Disk
1853	*(1) Scan ExistingRDD[event_id#28, customer_id#29, event_type#30, product_id#40, event_timestamp#57, event_date#72, ingestion_timestamp#80]	Disk Memory Deserialized 1x Replicated	3	100.00%	25.1 KiB	0.0 B

## SQL / DataFrame

Completed Queries: 84

### • Completed Queries (84)

ID	Description	Submitted	Duration	Job IDs	Sub Execution IDs
301	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 93 +details	2025/12/12 16:53:08	0.6 s	[302]	+details
299	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 92 +details	2025/12/12 16:52:51	0.6 s	[300]	+details
292	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 81 +details	2025/12/12 16:52:41	20 s	[293][294][295][296][297][298]	+details
290	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 91 +details	2025/12/12 16:52:32	0.6 s	[291]	+details
288	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 60 +details	2025/12/12 16:52:12	0.4 s	[289]	+details
281	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 60 +details	2025/12/12 16:52:04	21 s	[282][283][284][285][286][287]	+details
279	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 89 +details	2025/12/12 16:51:54	0.3 s	[280]	+details
277	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 88 +details	2025/12/12 16:51:35	0.6 s	[278]	+details
270	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 49 +details	2025/12/12 16:51:26	20 s	[271][272][273][274][275][276]	+details
268	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 87 +details	2025/12/12 16:51:19	0.5 s	[269]	+details
266	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 86 +details	2025/12/12 16:51:00	0.3 s	[267]	+details
259	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 48 +details	2025/12/12 16:50:49	19 s	[260][261][262][263][264][265]	+details
257	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 85 +details	2025/12/12 16:50:42	0.5 s	[258]	+details
255	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 84 +details	2025/12/12 16:50:25	0.6 s	[256]	+details
248	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 47 +details	2025/12/12 16:50:11	19 s	[249][250][251][252][253][254]	+details
246	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 83 +details	2025/12/12 16:50:09	0.5 s	[247]	+details
244	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 82 +details	2025/12/12 16:49:51	0.6 s	[245]	+details
240	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 81 +details	2025/12/12 16:49:34	0.6 s	[241]	+details
235	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 78 +details	2025/12/12 16:49:33	20 s	[236][237][238][239][242][243]	+details
233	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 80 +details	2025/12/12 16:49:18	0.6 s	[234]	+details
231	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 79 +details	2025/12/12 16:49:01	0.6 s	[232]	+details
224	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 45 +details	2025/12/12 16:48:56	20 s	[225][226][227][228][229][230]	+details
222	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 76 +details	2025/12/12 16:48:43	1 s	[223]	+details
220	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 77 +details	2025/12/12 16:48:25	0.5 s	[221]	+details
213	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 44 +details	2025/12/12 16:48:17	21 s	[214][215][216][217][218][219]	+details
211	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 76 +details	2025/12/12 16:48:07	0.3 s	[212]	+details
209	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 75 +details	2025/12/12 16:47:49	0.4 s	[210]	+details
202	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 43 +details	2025/12/12 16:47:42	19 s	[203][204][205][206][207][208]	+details
200	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 74 +details	2025/12/12 16:47:31	0.6 s	[201]	+details
198	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 73 +details	2025/12/12 16:47:13	0.4 s	[199]	+details
191	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 42 +details	2025/12/12 16:47:06	19 s	[192][193][194][195][196][197]	+details
189	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 72 +details	2025/12/12 16:46:55	0.5 s	[190]	+details
187	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 71 +details	2025/12/12 16:46:38	0.3 s	[188]	+details
180	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 49 +details	2025/12/12 16:46:28	20 s	[181][182][183][184][185][186]	+details
178	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 70 +details	2025/12/12 16:46:19	0.3 s	[179]	+details
176	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 40 +details	2025/12/12 16:46:02	0.5 s	[177]	+details
169	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 68 +details	2025/12/12 16:45:51	20 s	[170][171][172][173][174][175]	+details
167	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 67 +details	2025/12/12 16:45:44	0.3 s	[168]	+details
165	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 67 +details	2025/12/12 16:45:25	0.6 s	[166]	+details
158	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 36 +details	2025/12/12 16:45:15	19 s	[159][160][161][162][163][164]	+details
156	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 66 +details	2025/12/12 16:45:06	0.6 s	[157]	+details
154	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 65 +details	2025/12/12 16:44:49	0.3 s	[155]	+details
147	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 38 +details	2025/12/12 16:44:37	20 s	[148][149][150][151][152][153]	+details
145	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 64 +details	2025/12/12 16:44:30	0.6 s	[146]	+details
143	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 63 +details	2025/12/12 16:44:13	0.6 s	[144]	+details
136	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 37 +details	2025/12/12 16:44:03	18 s	[137][138][139][140][141][142]	+details
134	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 62 +details	2025/12/12 16:43:54	0.6 s	[135]	+details
132	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 61 +details	2025/12/12 16:43:37	0.3 s	[133]	+details
125	id = 83aa934f-4d58-4e5c-a3d0-651fc72b3e58 runid = e074dee9-ee82-484a-9f35-bc4782f3c2ac batch = 36 +details	2025/12/12 16:43:25	19 s	[126][127][128][129][130][131]	+details
123	id = d9e02216-a60c-4a3a-91ec-44ea71232773 runid = f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0 batch = 60 +details	2025/12/12 16:43:20	0.4 s	[124]	+details

## Streaming Query

### Active Streaming Queries (2)

Page: 1 1 Pages. Jump to 1 Show 100 items in a page. Go

Name	Status	ID	Run ID	Start Time	Duration	Avg Input/sec	Avg Process/sec	Latest Batch
<no name>	RUNNING	83aa934f-4d58-4e5c-a3d0-651fc72b3e58	e074dee9-ee82-484a-9f35-bc4782f3c2ac	2025/12/12 16:36:18	18 minutes 1 second	9.62	9.91	52
<no name>	RUNNING	d9e02216-a60c-4a3a-91ec-44ea71232773	f5a4ecf2-d4c7-45e8-8c2e-5873c54c9cc0	2025/12/12 16:36:13	18 minutes 6 seconds	9.81	10.06	96

Page: 1 1 Pages. Jump to 1 Show 100 items in a page. Go

akka-event-gen main

Current File ▾ ▶ ⚡ ⚡ ⚡ ⚡ ⚡ ⚡ Trial: 2 days

Project ▾

- akka-event-gen ~Desktop/ScalaTraining/Phase2Projects/akka-event-gen
- src
- aws
- ...

File Project Terminal Local

dashboard-api PipelineA PipelineB PipelineC akka-event-gen data-gen – AppendYestTodayTxns.

Search Everywhere Double ⌂

Go to File ⌂ ⌂ ⌂

Recent Files ⌂ ⌂ ⌂

Navigation Bar ⌂ ⌂

Drop files here to open them

rac@PTPMR107 akka-event-gen % sbt clean run

```
[info] welcome to sbt 1.11.7 (Homebrew Java 11.0.29)
[info] loading global plugins from /Users/rac/.sbt/1.0/plugins
[info] loading project definition from /Users/rac/Desktop/ScalaTraining/Phase2Projects/akka-event-gen/project
[info] loading settings for project root from build.sbt...
[info] set current project to akka-event-gen (in build file:/Users/rac/Desktop/ScalaTraining/Phase2Projects/akka-event-gen/)
[success] Total time: 0 s, completed Dec 12, 2025, 4:35:41 PM
[info] compiling 7 Scala sources to /Users/rac/Desktop/ScalaTraining/Phase2Projects/akka-event-gen/target/scala-2.12/classes ...
[info] running events.AvroEventProducerApp
```

File Location -- zsh ... 306-u admin -p -- zsh -- zsh ... eeper.properties ... rver1.properties ... rver2.properties ... rver3.properties ... Projects -- zsh

9c234d27-c572-4bd5-ae36-133144fcfa51
56922873-14a3-4846-98b4-5cd878518bc4
4e45909e-5367-4c3e-a43b-fc7e369a067c6
6ca89bd9-c91b-4514-8512-9275624a1e66
5939a2a2-4a09-4a09-9a09-1257b7f8c8b9
bb4cf7146-bc97-4c37-b934-59a884394a96
0794a4f7-749c-492a-adcf-3f2995e4a91af
22c9a23a-4a09-4a09-9a09-1257b7f8c8b9
a836c53e-e6ec-4887-9a8e-7289e52e8951
2945a42a-4a09-4a09-9a09-1257b7f8c8b9
45705dd7-0164-48c0-9720-75ba11c3f9c0
784897dd-1764-4f5e-1083-cde64a13e2ad
1204a23a-4a09-4a09-9a09-1257b7f8c8b9
d549fc37-6848-4fe0-8cd1-1257b7f8c8b9
4282154-a0d5-4387-8b17-2a82c0749ff3
8099a23a-4a09-4a09-9a09-1257b7f8c8b9
188e7769-3553-4f5b-9b94-852c1eeeb8134
1877a23a-4a09-4a09-9a09-1257b7f8c8b9
34a0b113-9375-43b0-93d2d27489b
7857e819 -f012-4db7-8297-b3e4b2fae95d
4794a23a-4a09-4a09-9a09-1257b7f8c8b9
82e4c158-e3c7-4a8e-9364-38233ab1t334
2fc4ddfb-3f6c-4381-a78 -673894ab0862
1204a23a-4a09-4a09-9a09-1257b7f8c8b9
184a4622-6b02-4698-8c8c-bf64240ea0d2
504a4622-6b02-4698-8c8c-bf64240ea0d2
136184d408b8-458b-49d3-1-1d64a4c943d4
ceaf11732-98d6-4f71-9390-f194158722dc
30000000-0000-0000-0000-000000000000
111226f2-42cd-4499-1f7b-ecf227835dc5
7334a79c-b16d-4d89-956d-f88e87922d6
1024a23a-4a09-4a09-9a09-1257b7f8c8b9
c0a212d9-9280-4401-93c0-22a32a45e52b2
608a4622-6b02-4698-8c8c-bf64240ea0d2
3943dcb9-7667-422a-8d8b-5e5d804a6e8e
6e85fb91-ac58-49a3-8611-19578009d9ab
7110ff12-3c36-4cd6-99c1-1edee7d5a64c
d05a2882-cecc-4aceb-4cb5-c-e2258a4a75c
54a55a3a-ae6b-4e15-963e-22a16235d5851
d05a2882-cecc-4aceb-4cb5-c-e2258a4a75c
4df1aca8-9d9d-479a-374f-4d4ebab1f1c9
e3b0a962-2c8b-4337-1b7b-94d6c6613164
8941652a-6e25-42d1-be34-1ab3b3b764e9b
bfca84d4-0b74-4e38-1b24-1c3c6c84a3d9e
7334a79c-b16d-4d89-956d-f88e87922d6
8981652a-6e25-42d1-be34-1ab3b3b764e9b
ad74d3ab-f9c1-4a88-8d84-162334bd6d4
40000000-0000-0000-0000-000000000000
7110ff12-3c36-4cd6-99c1-1edee7d5a64c
cdff16c0-213f-437e-8a81-1-4f793f1ae9e
c111226f2-42cd-4499-1f7b-ecf227835dc5
94f7134d-1c2a-43c3-1e19-1bd8bf4695742
00000000-0000-0000-0000-000000000000
0CP-processed a total of 401777 messages
rac@PTPMR107 ~ %

# Explanation of partitioning logic

## Overview

All three pipelines use partitioning to balance parallelism, reduce network shuffle, and make downstream storage/querying efficient. Partitioning happens at three places:

- (1) how JDBC reads from MySQL are partitioned,
- (2) how Spark repartitions data for groupBy/joins and downstream writes, and
- (3) how final data is partitioned in storage (Cassandra partition key or S3 hive-style date partitions).

---

## Pipeline A — MySQL → Spark → Cassandra (Customer profiles)

### Partitioning steps & heuristics

#### 1. JDBC partitioning (extractor.readDeltaTransactions / readTransactionsForDates)

- Uses `option("partitionColumn", "txn_id")` together with Spark `lowerBound/upperBound` and `numPartitions`.
- `choosePartitions(minId, maxId, cfg)` computes `numPartitions` with the heuristic `~10K rows / partition`:

```
val approxPerPartition = 10000L
val p = Math.max(2, Math.min(cfg, (span / approxPerPartition + 1).toInt))
```

- Purpose: parallelize JDBC reads across multiple connections/threads while avoiding too many tiny partitions.

#### 2. Spark repartition for aggregation

- Pre-aggregation: transactions are `.repartition($"customer_id")` before joins/aggregations. That collocates rows for each customer on the same partition and reduces shuffle during `groupBy(customer_id)`.
- After aggregation, before writing to Cassandra, code computes `repartitions`:

```
val repartitions = Math.max(8, Math.min(writeParallelism, Math.max(1, affectedCountHint * 4)))
```

- This caps/limits the number of output partitions to balance parallel Cassandra writes vs too many small SSTable writes.

#### 3. Cassandra partitioning (storage layout)

- Cassandra table `customer_profile` uses `customer_id` as the partition key (single row per customer). This enables very fast point lookups (API reads).
- Write strategy: repartition by `customer_id` attempts to align Spark partitions to Cassandra partitioning so each Spark task drives writes for a subset of customers.

---

## Pipeline B — MySQL → Spark Structured Streaming (ForeachBatch) → S3 (Daily summaries)

### Partitioning steps & heuristics

1. JDBC partitioning for delta
  - Same `txn_id` partitioning heuristic via `choosePartitions` when reading delta ranges. This parallelizes `readDeltaTransactions`.
2. Derive affected dates
  - The pipeline reads only `affectedDates` (`distinct to_date(txn_timestamp)` from `delta`) and then reads the full transactions for those dates using `readTransactionsForDates(dates)`. This minimizes the scope of re-computation.
3. Spark repartition & aggregation
  - Preprocess transactions with `.repartition($"customer_id")` in the transformer to colocate per-customer rows for the daily summary aggregations.
4. S3 partitioning
  - Final writes to S3 use `partitionBy("date")` and `SaveMode.Overwrite` for idempotent writes of the affected date partition:

```
.partitionBy("date")
.parquet(lakeBase)
```

- The `S3Loader.writeSummary` coalesces output files (`coalesce(coalesceNum)`) prior to write to control number/size of files per partition.

---

## Pipeline C — Kafka Avro → Spark Structured Streaming → S3 (Events)

### Partitioning steps & heuristics

1. Kafka partitioning (Producer behavior)
  - Your `AvroEventProducerActor` sets the Kafka message key to `eventId`:

```
val record = new ProducerRecord[String, Array[Byte]](topic, eventId,
bytes)
```

- This means events are partitioned by `eventId` hashing. If you need ordering per customer, you'd want to key by `customer_id` instead.
2. Spark streaming partitioning & processing

- Incoming Kafka stream is handled by Spark Structured Streaming. There is no explicit repartitioning in the extractor code beyond usual streaming micro-batch partitioning, but:
  - `from_avro` expands records.
  - `event_date = to_date(event_timestamp)` is added and used as partition column for final writes.

### 3. S3 partitioning

- Final writes use `partitionBy("event_date")` and append mode for streaming:

```
df.write.mode("append").partitionBy("event_date").parquet(basePath)
```

- Malformed records are handled by a separate `malformed` stream and can be written to a dedicated malformed path.

# Screenshots:

Amazon S3 < lake/

**Objects (2)**

Name	Type	Last modified	Size	Storage class
events/	Folder	-	-	-
txns_summary/	Folder	-	-	-

Amazon S3 < lake/ > txns\_summary/

**Objects (999+)**

Name	Type	Last modified	Size	Storage class
date=2024-12-11/	Folder	-	-	-
date=2024-12-12/	Folder	-	-	-
date=2024-12-21/	Folder	-	-	-
date=2024-12-22/	Folder	-	-	-
date=2024-12-23/	Folder	-	-	-
date=2024-12-25/	Folder	-	-	-
date=2024-12-27/	Folder	-	-	-
date=2024-12-28/	Folder	-	-	-
date=2024-12-29/	Folder	-	-	-
date=2025-12-08/	Folder	-	-	-
date=2025-12-09/	Folder	-	-	-
date=2025-12-10/	Folder	-	-	-
date=2025-12-11/	Folder	-	-	-
date=2025-12-12/	Folder	-	-	-

Amazon S3 < lake/ > events/

**Objects (7)**

Name	Type	Last modified	Size	Storage class
_checkpoints/	Folder	-	-	-
_spark_metadata/	Folder	-	-	-
_SUCCESS	Folder	December 12, 2025, 17:43:59 (UTC+05:30)	0 B	Standard
event_date=2025-12-09/	Folder	-	-	-
event_date=2025-12-10/	Folder	-	-	-
event_date=2025-12-11/	Folder	-	-	-
event_date=2025-12-12/	Folder	-	-	-

AWS | Search [Option+S] Account ID: 8069-8551-3100 ▾ sanjeev

Amazon S3 > Buckets > sanjeev-scala-s3 > lake/ > events/ > event\_date=2025-12-12/

### 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

### event\_date=2025-12-12/

Objects (279)

Name	Type	Last modified	Size	Storage class
part-00000-0083e89e-1461-46d4-a96e-f0cf72c82e39.c000.snappy.parquet	parquet	December 12, 2025, 17:21:08 (UTC+05:30)	8.2 KB	Standard
part-00000-0364698c-e26a-4a9c-b618-1b4042a9bb2c.c000.snappy.parquet	parquet	December 12, 2025, 16:47:55 (UTC+05:30)	7.8 KB	Standard
part-00000-03e8a15b-f5c4-4245-875c-c21562b307f5.c000.snappy.parquet	parquet	December 12, 2025, 17:16:05 (UTC+05:30)	39.9 KB	Standard
part-00000-05259e41-4287-4999-872f-b51580bf97c.c000.snappy.parquet	parquet	December 12, 2025, 17:18:00 (UTC+05:30)	8.1 KB	Standard

[Copy S3 URI](#)

AWS | Search [Option+S] Account ID: 8069-8551-3100 ▾ sanjeev

Amazon Keyspaces > CQL editor

### Amazon Keyspaces

- Dashboard
- Keyspaces
- Tables
- CQL editor**
- Configuration

Getting started exercise

Getting started resources [New](#)

Code samples [New](#)

Documentation [New](#)

Execution time: 33 ms

Table view JSON view

### Records returned (5000)

customer_id	avg_order_value	email	favorite_category	first_purchase	gender	last_purcha
2549	20841.68	customer_2549@example.com	Toys	2020-08-16 04:38:54.0+0000	O	2025-12-11
4581	21920.68	customer_4581@example.com	Health	2020-07-05 05:45:34.0+0000	M	2024-10-01
2316	28885.88	customer_2316@example.com	Electronics	2020-03-26 05:20:31.0+0000	M	2024-10-19
2093	20439.44	customer_2093@example.com	Beauty	2020-05-21 00:41:55.0+0000	O	2024-07-22
899	23593.73	customer_899@example.com	Sports	2020-05-08 14:14:16.0+0000	O	2025-12-10
867	22986.93	customer_867@example.com	Grocery	2020-06-27 22:12:38.0+0000	M	2025-12-10
117	31507.66	customer_117@example.com	Beauty	2020-03-28 20:08:03.0+0000	M	2024-04-06
2992	32454.97	customer_2992@example.com	Home	2020-01-13 12:32:59.0+0000	F	2025-12-08
4529	25867.24	customer_4529@example.com	Grocery	2020-10-24 13:29:09.0+0000	O	2024-06-29
4046	20673.37	customer_4046@example.com	Grocery	2020-04-02 03:00:06.0+0000	O	2025-12-09

Download results to CSV

Home Workspaces API Network

file-handle New Import

REST API basics: CRUD, test & variable / Customer Profile

GET {{base\_url}}/customer/4561

Customer Profile

Daily summary

events

Daily summary

Query Params

Key	Value	Description
Key	Value	Description

Body Cookies Headers (9) Test Results (1/1)

200 OK 254 ms 692 B Save Response

{ } JSON Preview Visualize

```
1 {
2   "status": "success",
3   "message": "Customer profile retrieved successfully",
4   "data": [
5     {
6       "customer_id": 4561,
7       "avg_order_value": 27720.72,
8       "email": "customer_4561@example.com",
9       "favorite_category": "Grocery",
10      "first_purchase": "2020-04-16T13:43:12Z",
11      "gender": "F",
12      "last_purchase": "2024-11-09T02:56:37Z",
13      "name": "Customer_4561",
14      "total_spend": 304927.92,
15      "total_transactions": 11
16    }
17 }
```

Save Share

REST API basics: CRUD, test & variable / Daily summary

GET {{base\_url}}/summary/2025-12-10

Daily summary

events

Daily summary

Pre-request

1 Use JavaScript to write tests, visualize response, and more. ⌘P to Ask AI

Post-response

Body Cookies Headers (9) Test Results

200 OK 1.18 s 12.82 KB Save Response

{ } JSON Preview Visualize

```
1 {
2   "status": "success",
3   "message": "Daily transaction summaries retrieved successfully",
4   "data": [
5     {
6       "summaries": [
7         {
8           "date": "2025-12-10",
9           "customer_id": 2142,
10          "total_amount": 38143.2,
11          "total_items": 8,
12          "distinct_products": 1,
13          "top_category": "Electronics"
14        },
15        {
16           "date": "2025-12-10",
17           "customer_id": 1088,
18         }
19       ]
20     }
21   ]
22 }
```

Save Share

Home Workspaces API Network

file-handle New Import

REST API basics: CRUD, test & variable / Daily summary

GET {{base\_url}}/summary/2025-12-10

Daily summary

events

Daily summary

Pre-request

1 Use JavaScript to write tests, visualize response, and more. ⌘P to Ask AI

Post-response

Body Cookies Headers (9) Test Results

200 OK 1.18 s 12.82 KB Save Response

{ } JSON Preview Visualize

```
1 {
2   "status": "success",
3   "message": "Daily transaction summaries retrieved successfully",
4   "data": [
5     {
6       "summaries": [
7         {
8           "date": "2025-12-10",
9           "customer_id": 2142,
10          "total_amount": 38143.2,
11          "total_items": 8,
12          "distinct_products": 1,
13          "top_category": "Electronics"
14        },
15        {
16           "date": "2025-12-10",
17           "customer_id": 1088,
18         }
19       ]
20     }
21   ]
22 }
```

Save Share

Home Workspaces API Network

file-handle New Import

GET Customer | GET Daily summar | GET events | GET events | GET Daily summar | REST API b... | Performance | + | No environment | Upgrade

file-handle Collections Environments History Flows

REST API basics: CRUD, test & variable / events

GET {{base\_url}} /events/1166?limit=10

Docs Params Authorization Headers (6) Body Scripts Settings Cookies

none form-data x-www-form-urlencoded raw binary GraphQL

This request does not have a body

Body Cookies Headers (9) Test Results

200 OK 5 ms 2.72 KB Save Response

{ JSON Preview Visualize }

```
1 {
  "status": "success",
  "message": "Customer events retrieved successfully",
  "data": {
    "events": [
      {
        "date": "2025-12-10",
        "event_id": "8849b88b-7336-46bd-9668-8c81ea01b2b1",
        "customer_id": 1166,
        "event_type": "WISHLIST",
        "product_id": 58,
        "event_timestamp_hex": "2025-12-09T18:50:45.929804Z",
        "ingestion_timestamp_hex": "2025-12-09T18:51:09.789Z"
      },
      {
        "date": "2025-12-11",
        "event_id": "d4b27a2-fcde-406b-a727-aaaae6c3c1d11",
        "customer_id": 1166,
        "event_type": "WISHLIST",
        "product_id": 4,
        "event_timestamp_hex": "2025-12-11T04:45:06.317456Z",
        "ingestion_timestamp_hex": "2025-12-11T04:45:41.484Z"
      }
    ]
  }
}
```

Cloud View Find and replace Console Terminal Runner Start Proxy Cookies Vault Trash

Home Workspaces API Network

file-handle New Import

GET Customer | GET Daily summar | GET events | GET events | GET Daily summar | REST API b... | Performance | + | No environment | Upgrade

file-handle Collections Environments History Flows

REST API basics: CRUD, test & variable / events

GET {{base\_url}} /events/1166?fromDate=2025-12-11&toDate=2025-12-11&limit=50

Docs Params Authorization Headers (6) Body Scripts Settings Cookies

none form-data x-www-form-urlencoded raw binary GraphQL

This request does not have a body

Body Cookies Headers (9) Test Results

200 OK 3.32 s 11.92 KB Save Response

{ JSON Preview Visualize }

```
1 {
  "status": "success",
  "message": "Customer events retrieved successfully",
  "data": {
    "events": [
      {
        "date": "2025-12-11",
        "event_id": "d4b27a2-fcde-406b-a727-aaaae6c3c1d11",
        "customer_id": 1166,
        "event_type": "WISHLIST",
        "product_id": 4,
        "event_timestamp_hex": "2025-12-11T04:45:06.317456Z",
        "ingestion_timestamp_hex": "2025-12-11T04:45:41.484Z"
      }
    ]
  }
}
```

Cloud View Find and replace Console Terminal Runner Start Proxy Cookies Vault Trash

# Performance Report of api calls - Dec 11, 2025 (#1)

[Open in Postman](#)

Postman collection: REST API basics: CRUD, test & variable

Report exported on: Dec 11, 2025, 12:29:19 (GMT+5:30)

## Test setup

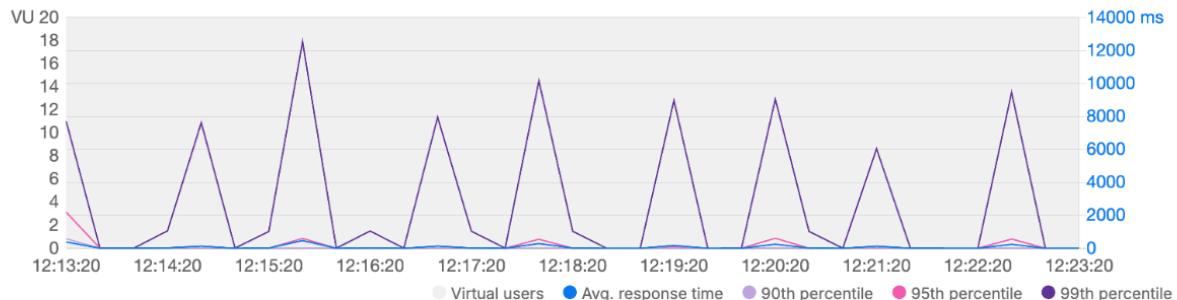
Virtual users	Start time	Load profile
20 VU	Dec 11, 12:13:20 (GMT+5:30)	Fixed
Duration	End time	Environment
10 minutes	Dec 11, 12:23:27 (GMT+5:30)	-

## 1. Summary

Total requests sent	Throughput	Average response time	Error rate
39,968	65.93 requests/second	50 ms	0.00 %

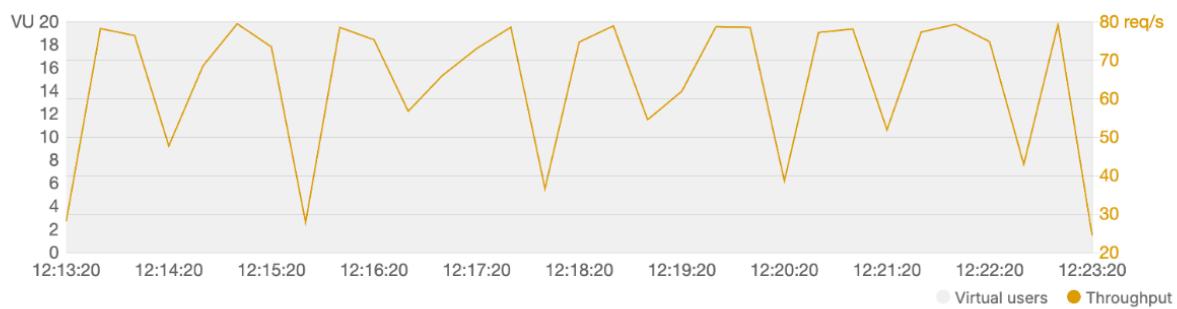
### 1.1 Response time

Response time trends during the test duration.



### 1.2 Throughput

Rate of requests sent per second during the test duration.



### 1.3 Requests with slowest response times

Top 5 slowest requests based on their average response times.

Request	Resp. time (Avg ms)	90th (ms)	95th (ms)	99th (ms)	Min (ms)	Max (ms)
<b>GET events</b> {{base_url}}/events/1166?limit=100	161	3	4	8,441	1	12,943
<b>GET Daily summary</b> {{base_url}}/summary/2025-12-10	23	4	5	1,045	1	3,040
<b>GET Daily summary</b> {{base_url}}/summary/2025-12-10/3156	13	2	3	565	1	1,063
<b>GET Customer Profile</b> {{base_url}}/customer/456	3	3	4	15	1	386

### 2. Metrics for each request

The requests are shown in the order they were sent by virtual users.

Request	Total requests	Requests/s	Min (ms)	Avg (ms)	90th (ms)	Max (ms)	Error %
<b>GET Customer Profile</b> {{base_url}}/customer/456	9,992	16.48	1	3	3	386	0
<b>GET Daily summary</b> {{base_url}}/summary/2025-12-10	9,992	16.48	1	23	4	3,040	0
<b>GET events</b> {{base_url}}/events/1166?limit=100	9,992	16.48	1	161	3	12,943	0
<b>GET Daily summary</b> {{base_url}}/summary/2025-12-10/3156	9,992	16.48	1	13	2	1,063	0

## Sample JSON responses

### 1. Customer Profile API

GET /customer/:id

---

#### Sample Success Response

```
{  
  "status": "success",  
  "message": "Customer profile retrieved successfully",  
  "data": {  
    "customer_id": 456,  
    "name": "Customer_456",  
    "email": "customer_456@example.com",  
    "gender": "M",  
    "total_spend": 143949.5,  
    "total_transactions": 6,  
    "avg_order_value": 23991.58,  
    "first_purchase": "2022-03-07T06:30:00Z",  
    "last_purchase": "2024-05-18T18:17:40Z",  
    "favorite_category": "Grocery"  
  }  
}
```

#### Sample Not-Found Response

```
{  
  "status": "error",  
  "message": "not-found",  
  "data": {}
```

```
}
```

---

## 2. Daily Summary API

**GET /summary/:date/:customerId**

---

### Sample Success Response

```
{
  "status": "success",
  "message": "Daily transaction summary retrieved successfully",
  "data": {
    "summary": {
      "date": "2025-12-10",
      "customer_id": 3156,
      "total_amount": 12895.28,
      "total_items": 2,
      "distinct_products": 1,
      "top_category": "Electronics"
    }
  }
}
```

### Sample Not-Found Response

```
{
  "status": "error",
  "message": "not-found",
  "data": {}
}
```

## Sample Invalid Date Response

```
{  
  "status": "error",  
  "message": "invalid_date_format",  
  "data": {}  
}
```

---

## 3. Events API

**GET /events/:customerId?fromDate=&toDate=&limit=**

---

## Sample Success Response

```
{  
  "status": "success",  
  "message": "Customer events retrieved successfully",  
  "data": {  
    "events": [  
      {  
        "date": "2025-12-09",  
        "event_id": "dd4b27a2-fcde-406b-a727-aeee6c3c1d11"  
        "customer_id": 2942,  
        "event_type": "CART_ADD",  
        "product_id": 491,  
        "event_timestamp": "2025-12-09T18:18:50.850395Z",  
        "ingestion_timestamp": "2025-12-09T18:18:59.149Z"  
      }  
    ]  
  }
```

```
}
```

### Sample Invalid fromDate Response

```
{
  "status": "error",
  "message": "invalid_fromDate_format",
  "data": {}
}
```

### Sample Invalid toDate Response

```
{
  "status": "error",
  "message": "invalid_toDate_format",
  "data": {}
}
```

### Sample Invalid Limit Response

```
{
  "status": "error",
  "message": "invalid_limit",
  "data": {}
}
```

### Sample Customer Not Found Response

```
{
  "status": "error",
  "message": "not-found",
  "data": {}
}
```

---

## 4. Daily Summaries API

**GET /summary/:date**

---

### Sample Success Response

```
{  
  "status": "success",  
  "message": "Daily transaction summaries retrieved successfully",  
  "data": {  
    "summaries": [  
      {  
        "date": "2025-12-10",  
        "customer_id": 1088,  
        "total_amount": 57714.9,  
        "total_items": 6,  
        "distinct_products": 1,  
        "top_category": "Sports"  
      },  
      {  
        "date": "2025-12-10",  
        "customer_id": 3156,  
        "total_amount": 12895.28,  
        "total_items": 2,  
        "distinct_products": 1,  
        "top_category": "Electronics"  
      }  
    ]  
  }  
}
```

## **Sample Not-Found Response**

```
{  
  "status": "error",  
  "message": "not-found",  
  "data": {}  
}
```

## **Sample Invalid Date Response**

```
{  
  "status": "error",  
  "message": "invalid_date_format",  
  "data": {}  
}
```