

Lab: Split a stream of events into substreams

Problem Statement:

How do you split events in a Kafka topic so that the events are placed into subtopics?

Example use case:

Suppose that you have a Kafka topic representing appearances of an actor or actress in a film, with each event denoting the genre. In this lab, we'll write a program that splits the stream into substreams based on the genre. We'll have a topic for drama films, a topic for fantasy films, and a topic for everything else.

Hands-on code example:

Run it

1. Prerequisites
2. Initialize the project
3. Get Confluent Platform
4. Write the program interactively using the CLI
5. Write your statements to a file

Test it

1. Create the test data
2. Invoke the tests

Run it

Prerequisites

This lab installs Confluent Platform using Docker. Before proceeding:

- Connect with lab environment VM using SSH:

```
ssh USERNAME@YOUR_VM_DNS.courseware.io
```

 - **Username:** Will be provided by Instructor.
 - **Password:** Will be provided by Instructor.
- Verify that Docker is set up properly by ensuring no errors are output when you run `docker info` and `docker compose version` on the command line.

Initialize the project

To get started, make a new directory anywhere you'd like for this project:

```
mkdir split-stream && cd split-stream
```

Then make the following directories to set up its structure:

```
mkdir src test
```

Get Confluent Platform

Next, create the following `docker-compose.yml` file to obtain Confluent Platform:

```
---
version: '2'

services:
  zookeeper:
    image: confluentinc/cp-zookeeper:7.3.0
    hostname: zookeeper
    container_name: zookeeper
    ports:
      - "2181:2181"
    environment:
      ZOOKEEPER_CLIENT_PORT: 2181
      ZOOKEEPER_TICK_TIME: 2000

  broker:
    image: confluentinc/cp-kafka:7.3.0
    hostname: broker
    container_name: broker
    depends_on:
      - zookeeper
    ports:
      - "29092:29092"
    environment:
      KAFKA_BROKER_ID: 1
      KAFKA_ZOOKEEPER_CONNECT: 'zookeeper:2181'
      KAFKA_LISTENER_SECURITY_PROTOCOL_MAP:
PLAINTEXT:PLAINTEXT,PLAINTEXT_HOST:PLAINTEXT
      KAFKA_ADVERTISED_LISTENERS:
PLAINTEXT://broker:9092,PLAINTEXT_HOST://localhost:29092
      KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
      KAFKA_TRANSACTION_STATE_LOG_MIN_ISR: 1
      KAFKA_TRANSACTION_STATE_LOG_REPLICATION_FACTOR: 1
      KAFKA_GROUP_INITIAL_REBALANCE_DELAY_MS: 0

  schema-registry:
    image: confluentinc/cp-schema-registry:7.3.0
    hostname: schema-registry
    container_name: schema-registry
    depends_on:
      - broker
    ports:
      - "8081:8081"
    environment:
      SCHEMA_REGISTRY_HOST_NAME: schema-registry
      SCHEMA_REGISTRY_KAFKASTORE_BOOTSTRAP_SERVERS: 'broker:9092'

  ksqldb-server:
    image: confluentinc/ksqldb-server:0.28.2
```

```
hostname: ksqldb-server
container_name: ksqldb-server
depends_on:
  - broker
  - schema-registry
ports:
  - "8088:8088"
environment:
  KSQL_CONFIG_DIR: "/etc/ksqldb"
  KSQL_LOG4J_OPTS: "-Dlog4j.configuration=file:/etc/ksqldb/log4j.properties"
  KSQL_BOOTSTRAP_SERVERS: "broker:9092"
  KSQL_HOST_NAME: ksqldb-server
  KSQL_LISTENERS: "http://0.0.0.0:8088"
  KSQL_CACHE_MAX_BYTES_BUFFERING: 0
  KSQL_KSQL_SCHEMA_REGISTRY_URL: "http://schema-registry:8081"

ksqldb-cli:
  image: confluentinc/ksqldb-cli:0.28.2
  container_name: ksqldb-cli
  depends_on:
    - broker
    - ksqldb-server
  entrypoint: /bin/sh
  environment:
    KSQL_CONFIG_DIR: "/etc/ksqldb"
  tty: true
  volumes:
    - ./src:/opt/app/src
    - ./test:/opt/app/test
```

And launch it by running:

```
docker compose up -d
```

```

ubuntu@ip-172-31-28-38:~/split-stream$ docker compose up -d
[+] Running 38/38
✓ schema-registry 2 layers [###] 0B/0B Pulled
  ✓ 60f0dae92e99 Pull complete
  ✓ 20ff71da0dab Pull complete
✓ ksqldb-cli Pulled
✓ zookeeper 2 layers [###] 0B/0B Pulled
  ✓ fa08a06f385f Pull complete
  ✓ bddb49e2fc4d Pull complete
✓ broker 12 layers [#####] 0B/0B Pulled
  ✓ d5d2e87c6892 Pull complete
  ✓ 008dba906bf6 Pull complete
  ✓ bfeaabe01655 Pull complete
  ✓ 2cb7eb0f5666 Pull complete
  ✓ f70f416c6ce7 Pull complete
  ✓ bc67d000e59b Pull complete
  ✓ d6e744651f37 Pull complete
  ✓ 0427d86fae81 Pull complete
  ✓ 4108e73e61e1 Pull complete
  ✓ ac5563423559 Pull complete
  ✓ d32323e291f3 Pull complete
  ✓ ee69ff430d89 Pull complete
✓ ksqldb-server 17 layers [#####] 0B/0B Pulled
  ✓ 131f1a26eef0 Pull complete
  ✓ 16b78ed2e822 Pull complete
  ✓ e7233e20a08e Pull complete
  ✓ 9b9cbc4e490c Pull complete
  ✓ f7352f1ba78a Pull complete
  ✓ e06c925589f1 Pull complete
  ✓ 3a498372ace6 Pull complete
  ✓ 2ea0fca7e4a4 Pull complete
  ✓ c33c6c0781fc Pull complete
  ✓ 398ad6a64f8a Pull complete
  ✓ 71de1fc320c5 Pull complete
  ✓ 406f0318d1a6 Pull complete
  ✓ 7fd8e6133860 Pull complete
  ✓ ba53aced95c4 Pull complete
  ✓ 6ee1df7caef1 Pull complete
  ✓ 1bf25f440b73 Pull complete
  ✓ 2224172ddedb Pull complete
[+] Running 6/6
✓ Network split-stream_default Created
✓ Container zookeeper Started
✓ Container broker Started
✓ Container schema-registry Started
✓ Container ksqldb-server Started
✓ Container ksqldb-cli Started

```

Write the program interactively using the CLI

To begin developing interactively, open up the ksqldb CLI:

```
docker exec -it ksqldb-cli ksql http://ksqldb-server:8088
```

First, you'll need to create a Kafka topic and stream to represent the actors. The following creates both in one shot:

```
CREATE STREAM actingevents (name VARCHAR, title VARCHAR, genre VARCHAR)
  WITH (KAFKA_TOPIC = 'acting-events', PARTITIONS = 1, VALUE_FORMAT = 'AVRO');
```

Then produce the following events to the stream:

```
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Bill Murray', 'Ghostbusters',
'fantasy');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Christian Bale', 'The Dark
Knight', 'crime');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Diane Keaton', 'The Godfather:
Part II', 'crime');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Jennifer Aniston', 'Office
Space', 'comedy');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Judy Garland', 'The Wizard of
```

```
Oz', 'fantasy');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Keanu Reeves', 'The Matrix',
'fantasy');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Laura Dern', 'Jurassic Park',
'fantasy');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Matt Damon', 'The Martian',
'drama');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Meryl Streep', 'The Iron Lady',
'drama');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Russell Crowe', 'Gladiator',
'drama');
INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Will Smith', 'Men in Black',
'comedy');
```

Now that you have stream with some events in it, let's read them out. The first thing to do is set the following properties to ensure that you're reading from the beginning of the stream:

```
SET 'auto.offset.reset' = 'earliest';
```

Let's find all of the drama films. Issue the following transient push query. This will block and continue to return results until it's limit is reached or you tell it to stop.

```
SELECT NAME, TITLE FROM ACTINGEVENTS WHERE GENRE='drama' EMIT CHANGES LIMIT 3;
```

```
ubuntu@ip-172-31-28-38:~/split-stream$ docker exec -it ksqldb-cli ksql http://ksqldb-server:8088
OpenJDK 64-Bit Server VM warning: Option UseConcMarkSweepGC was deprecated in version 9.0 and will likely be removed in a future release.

=====
| ksqldb |
=====
The Database purpose-built
for stream processing apps
=====

Copyright 2017-2022 Confluent Inc.

CLI v0.28.2, Server v0.28.2 located at http://ksqldb-server:8088
Server Status: RUNNING

Having trouble? Type 'help' (case-insensitive) for a rundown of how things work!

ksql> CREATE STREAM actingevents (name VARCHAR, title VARCHAR, genre VARCHAR)
> WITH (KAFKA_TOPIC = 'acting-events', PARTITIONS = 1, VALUE_FORMAT = 'AVRO');

Message
-----
Stream created

ksql> INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Bill Murray', 'Ghostbusters', 'fantasy');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Christian Bale', 'The Dark Knight', 'crime');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Diane Keaton', 'The Godfather: Part II', 'crime');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Jennifer Aniston', 'Office Space', 'comedy');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Judy Garland', 'The Wizard of Oz', 'fantasy');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Keanu Reeves', 'The Matrix', 'fantasy');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Laura Dern', 'Jurassic Park', 'fantasy');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Matt Damon', 'The Martian', 'drama');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Meryl Streep', 'The Iron Lady', 'drama');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Russell Crowe', 'Gladiator', 'drama');
>INSERT INTO ACTINGEVENTS (name, title,genre) VALUES ('Will Smith', 'Men in Black', 'comedy');
>
ksql> SET 'auto.offset.reset' = 'earliest';

Successfully changed local property 'auto.offset.reset' to 'earliest'. Use the UNSET command to revert your change.
ksql> SELECT NAME, TITLE FROM ACTINGEVENTS WHERE GENRE='drama' EMIT CHANGES LIMIT 3;
+-----+-----+
|NAME|TITLE|
+-----+-----+
|Matt Damon|The Martian|
|Meryl Streep|The Iron Lady|
|Russell Crowe|Gladiator|
Limit Reached
Query terminated
ksql>
ksql> SELECT NAME, TITLE, GENRE FROM ACTINGEVENTS WHERE GENRE != 'drama' AND GENRE != 'fantasy' EMIT CHANGES LIMIT 4;
```

This should yield the following output:

```
+-----+-----+
|NAME|TITLE|
+-----+-----+
```

```
|Matt Damon          |The Martian          |
|Meryl Streep        |The Iron Lady        |
|Russell Crowe        |Gladiator            |
Limit Reached
Query terminated
```

You can also use negative matches, that is, messages that *don't* match the condition. Run this query to get a list of all films that aren't drama or fantasy.

```
SELECT NAME, TITLE, GENRE FROM ACTINGEVENTS WHERE GENRE != 'drama' AND GENRE !=
'fantasy' EMIT CHANGES LIMIT 4;
```

This should yield the following output:

```
+-----+-----+-----+
-----+
|NAME          |TITLE                |GENRE
|
+-----+-----+-----+
-----+
|Christian Bale |The Dark Knight      |crime
|
|Diane Keaton  |The Godfather: Part II|crime
|
|Jennifer Aniston|Office Space        |comedy
|
|Will Smith     |Men in Black         |comedy
|
Limit Reached
Query terminated
```

Since the output looks right, the next step is to make the queries continuous. Issue the following to create three new streams that are continuously populated by the queries:

```
CREATE STREAM actingevents_drama AS
  SELECT NAME, TITLE
    FROM ACTINGEVENTS
   WHERE GENRE='drama';

CREATE STREAM actingevents_fantasy AS
  SELECT NAME, TITLE
    FROM ACTINGEVENTS
   WHERE GENRE='fantasy';

CREATE STREAM actingevents_other AS
  SELECT NAME, TITLE, GENRE
    FROM ACTINGEVENTS
   WHERE GENRE != 'drama'
     AND GENRE != 'fantasy';
```

To check that it's working, print out the contents of one of the output stream's underlying topic.

```
PRINT ACTINGEVENTS_FANTASY FROM BEGINNING LIMIT 4;
```

This should yield the following output:

```
Key format: \_(`\`)\_/ - no data processed
Value format: AVRO or KAFKA_STRING
rowtime: 2020/05/04 23:40:52.078 Z, key: <null>, value: {"NAME": "Bill Murray",
"TITLE": "Ghostbusters"}, partition: 0
rowtime: 2020/05/04 23:40:52.454 Z, key: <null>, value: {"NAME": "Judy Garland",
"TITLE": "The Wizard of Oz"}, partition: 0
rowtime: 2020/05/04 23:40:52.537 Z, key: <null>, value: {"NAME": "Keanu Reeves",
"TITLE": "The Matrix"}, partition: 0
rowtime: 2020/05/04 23:40:52.607 Z, key: <null>, value: {"NAME": "Laura Dern",
"TITLE": "Jurassic Park"}, partition: 0
Topic printing ceased
```

Try dropping the `LIMIT` from the print command so that it runs indefinitely. To see how any new message on the source stream is automatically routed to the correct target stream, open a new CLI session and insert a record like we did above.

Write your statements to a file

Now that you have a series of statements that's doing the right thing, the last step is to put them into a file so that they can be used outside the CLI session. Create a file at `src/statements.sql` with the following content:

```
CREATE STREAM actingevents (name VARCHAR, title VARCHAR, genre VARCHAR)
  WITH (KAFKA_TOPIC = 'acting-events', PARTITIONS = 1, VALUE_FORMAT = 'AVRO');

CREATE STREAM actingevents_drama AS
  SELECT NAME, TITLE
  FROM ACTINGEVENTS
  WHERE GENRE='drama';

CREATE STREAM actingevents_fantasy AS
  SELECT NAME, TITLE
  FROM ACTINGEVENTS
  WHERE GENRE='fantasy';

CREATE STREAM actingevents_other AS
  SELECT NAME, TITLE, GENRE
  FROM ACTINGEVENTS
  WHERE GENRE != 'drama'
  AND GENRE != 'fantasy';
```

Test it

Create the test data

Create a file at `test/input.json` with the inputs for testing:

```
{
  "inputs": [
    {
      "topic": "acting-events",
      "value": {
        "name": "Bill Murray",
        "title": "Ghostbusters",
        "genre": "fantasy"
      }
    },
    {
      "topic": "acting-events",
      "value": {
        "name": "Christian Bale",
        "title": "The Dark Knight",
        "genre": "crime"
      }
    },
    {
      "topic": "acting-events",
      "value": {
        "name": "Diane Keaton",
        "title": "The Godfather: Part II",
        "genre": "crime"
      }
    },
    {
      "topic": "acting-events",
      "value": {
        "name": "Jennifer Aniston",
        "title": "Office Space",
        "genre": "comedy"
      }
    },
    {
      "topic": "acting-events",
      "value": {
        "name": "Judy Garland",
        "title": "The Wizard of Oz",
        "genre": "fantasy"
      }
    },
    {
      "topic": "acting-events",
      "value": {
        "name": "Keanu Reeves",
        "title": "The Matrix",
        "genre": "fantasy"
      }
    },
    {
      "topic": "acting-events",
```



```

    "value": {
      "name": "Laura Dern",
      "title": "Jurassic Park",
      "genre": "fantasy"
    }
  },
  {
    "topic": "acting-events",
    "value": {
      "name": "Matt Damon",
      "title": "The Martian",
      "genre": "drama"
    }
  },
  {
    "topic": "acting-events",
    "value": {
      "name": "Meryl Streep",
      "title": "The Iron Lady",
      "genre": "drama"
    }
  },
  {
    "topic": "acting-events",
    "value": {
      "name": "Russell Crowe",
      "title": "Gladiator",
      "genre": "drama"
    }
  },
  {
    "topic": "acting-events",
    "value": {
      "name": "Will Smith",
      "title": "Men in Black",
      "genre": "comedy"
    }
  },
  {
    "topic": "acting-events",
    "value": {
      "name": "Barret Oliver",
      "title": "The NeverEnding Story",
      "genre": "fantasy"
    }
  }
]
}

```

Similarly, create a file at `test/output.json` with the expected outputs:

```
{
  "outputs": [
    {
      "topic": "ACTINGEVENTS_FANTASY",
      "value": {
        "NAME": "Bill Murray",
        "TITLE": "Ghostbusters"
      }
    },
    {
      "topic": "ACTINGEVENTS_OTHER",
      "value": {
        "NAME": "Christian Bale",
        "TITLE": "The Dark Knight",
        "GENRE": "crime"
      }
    },
    {
      "topic": "ACTINGEVENTS_OTHER",
      "value": {
        "NAME": "Diane Keaton",
        "TITLE": "The Godfather: Part II",
        "GENRE": "crime"
      }
    },
    {
      "topic": "ACTINGEVENTS_OTHER",
      "value": {
        "NAME": "Jennifer Aniston",
        "TITLE": "Office Space",
        "GENRE": "comedy"
      }
    },
    {
      "topic": "ACTINGEVENTS_FANTASY",
      "value": {
        "NAME": "Judy Garland",
        "TITLE": "The Wizard of Oz"
      }
    },
    {
      "topic": "ACTINGEVENTS_FANTASY",
      "value": {
        "NAME": "Keanu Reeves",
        "TITLE": "The Matrix"
      }
    },
    {
      "topic": "ACTINGEVENTS_FANTASY",
      "value": {
        "NAME": "Laura Dern",
        "TITLE": "Jurassic Park"
      }
    }
  ]
}
```

```

    }
  },
  {
    "topic": "ACTINGEVENTS_DRAMA",
    "value": {
      "NAME": "Matt Damon",
      "TITLE": "The Martian"
    }
  },
  {
    "topic": "ACTINGEVENTS_DRAMA",
    "value": {
      "NAME": "Meryl Streep",
      "TITLE": "The Iron Lady"
    }
  },
  {
    "topic": "ACTINGEVENTS_DRAMA",
    "value": {
      "NAME": "Russell Crowe",
      "TITLE": "Gladiator"
    }
  },
  {
    "topic": "ACTINGEVENTS_OTHER",
    "value": {
      "NAME": "Will Smith",
      "TITLE": "Men in Black",
      "GENRE": "comedy"
    }
  },
  {
    "topic": "ACTINGEVENTS_FANTASY",
    "value": {
      "NAME": "Barret Oliver",
      "TITLE": "The NeverEnding Story"
    }
  }
]
}

```

Invoke the tests

Lastly, invoke the tests using the test runner and the statements file that you created earlier:

```

docker exec ksqldb-cli ksql-test-runner -i /opt/app/test/input.json -s
/opt/app/src/statements.sql -o /opt/app/test/output.json

```

Which should pass:

```
>>> Test passed!
```

Cleanup Resources

Delete all the resources by running following command in the `docker-compose.yml` file directory from the terminal:

```
docker compose down

docker container prune
```

```
ubuntu@ip-172-31-28-38:~/split-stream$ docker compose down
[+] Running 4/3
   Container schema-registry   Removed                                10.4s
   Container broker            Removed                                0.3s
   Container zookeeper         Removed                                0.3s
   Network split-stream_default Error                                0.6s
failed to remove network split-stream_default: Error response from daemon: error while removing network: network split-stream_default id 947bab6e7c6b74176aec928edf3567db72c046f65f5397bd1e7fee736cc30b3b has active endpoints
```

Note: If you get above error while running above command. Manually stop the containers and run `docker compose down` again. **Do not delete kafkanew container.**

```
no container to killubuntu@ip-172-31-28-38:~/split-stream$ docker ps
CONTAINER ID   IMAGE                                COMMAND                  CREATED        STATUS        PORTS
9465d6b4aa89   confluentinc/ksqldb-cli:0.28.2     "/bin/sh"               About an hour ago Up About an hour
ksqldb-cli
b6887b23c1ca   confluentinc/ksqldb-server:0.28.2  "/usr/bin/docker/run"   About an hour ago Up About an hour  0.0.0.0:8088->8088/tcp, :::8088->8088/tcp
ksqldb-server
22652a43d640   fenago/kafka-intellij-new-2.8      "/dockerstartup/vnc_..." 8 weeks ago    Up 12 hours    0.0.0.0:80->80/tcp, :::80->80/tcp, 5900/tcp, 0.0.0.0:80->80/tcp, :::80->80/tcp
p_kafkanew
ubuntu@ip-172-31-28-38:~/split-stream$ docker stop 9465d6b4aa89 b6887b23c1ca
9465d6b4aa89
b6887b23c1ca
ubuntu@ip-172-31-28-38:~/split-stream$
ubuntu@ip-172-31-28-38:~/split-stream$ docker compose down
[+] Running 1/1
   Network split-stream_default   Removed
```