

@GAMUSSA

#POSTGRES

@CONFLUENTINC

STREAMING ETL IN PRACTICE: BUILD DATA PIPELINES WITHOUT A SINGLE LINE OF CODE!





@GAMUSSA

Special Thanks!



@RMOFF

RAFFLE, YEAH 🚀

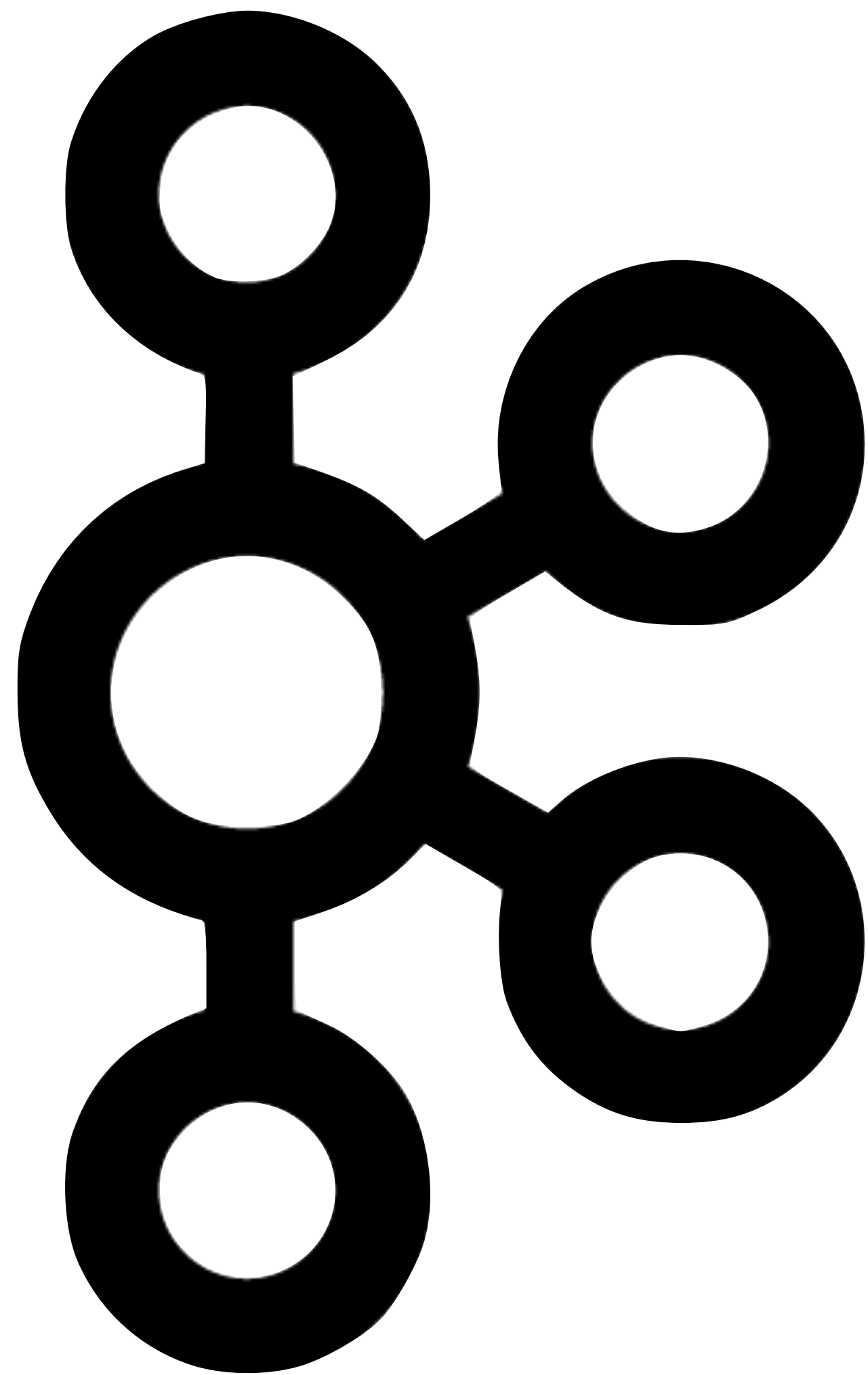


● must follow @gamussa @confluentinc



● Tag @gamussa

● With #peoplepostgresdata



APACHE
kafka



APACHE KAFKA IS AN EVENT STREAMING PLATFORM

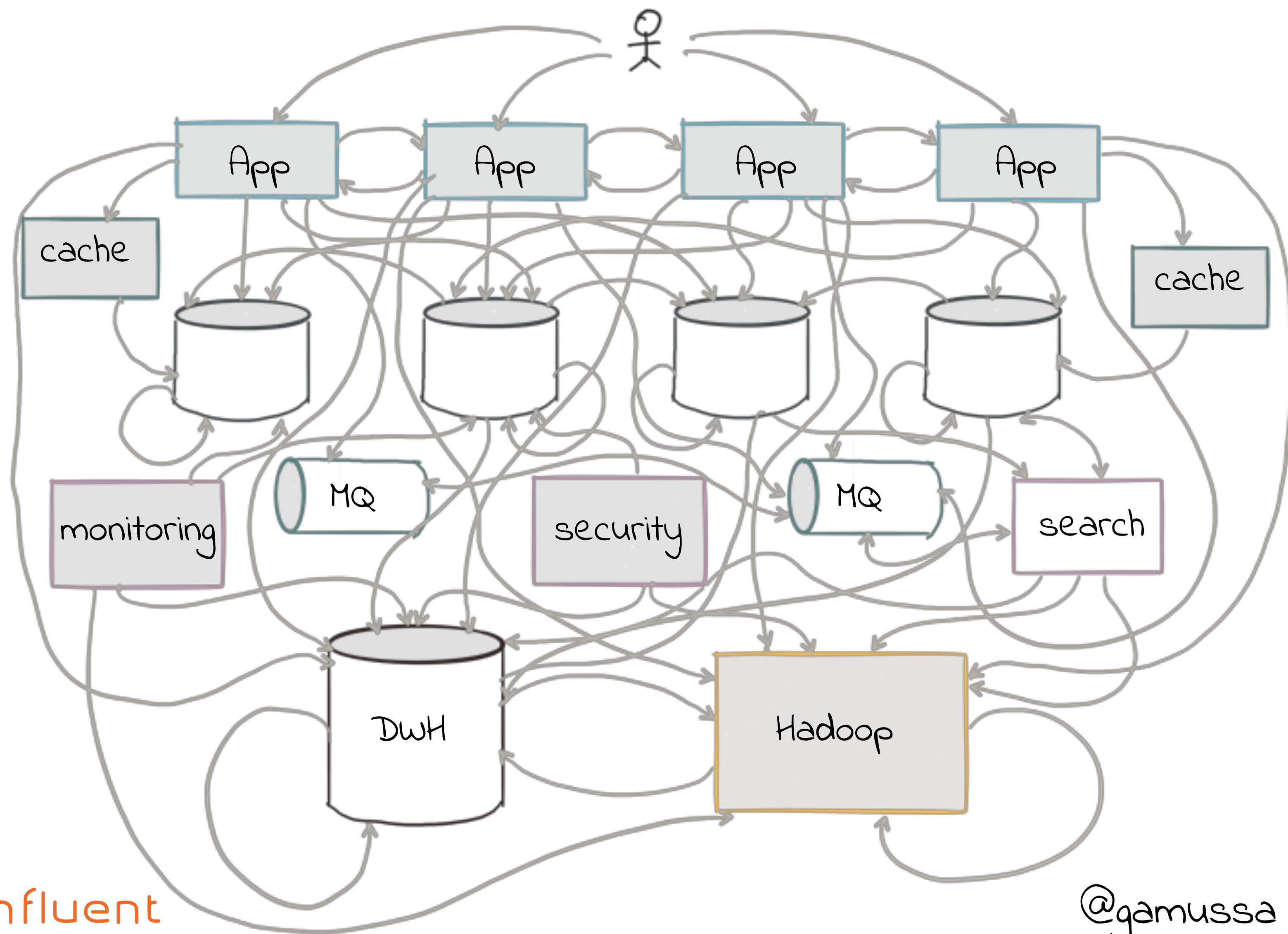
**BUT WHAT IS
AN EVENT STREAMING
PLATFORM?**

@gamussa

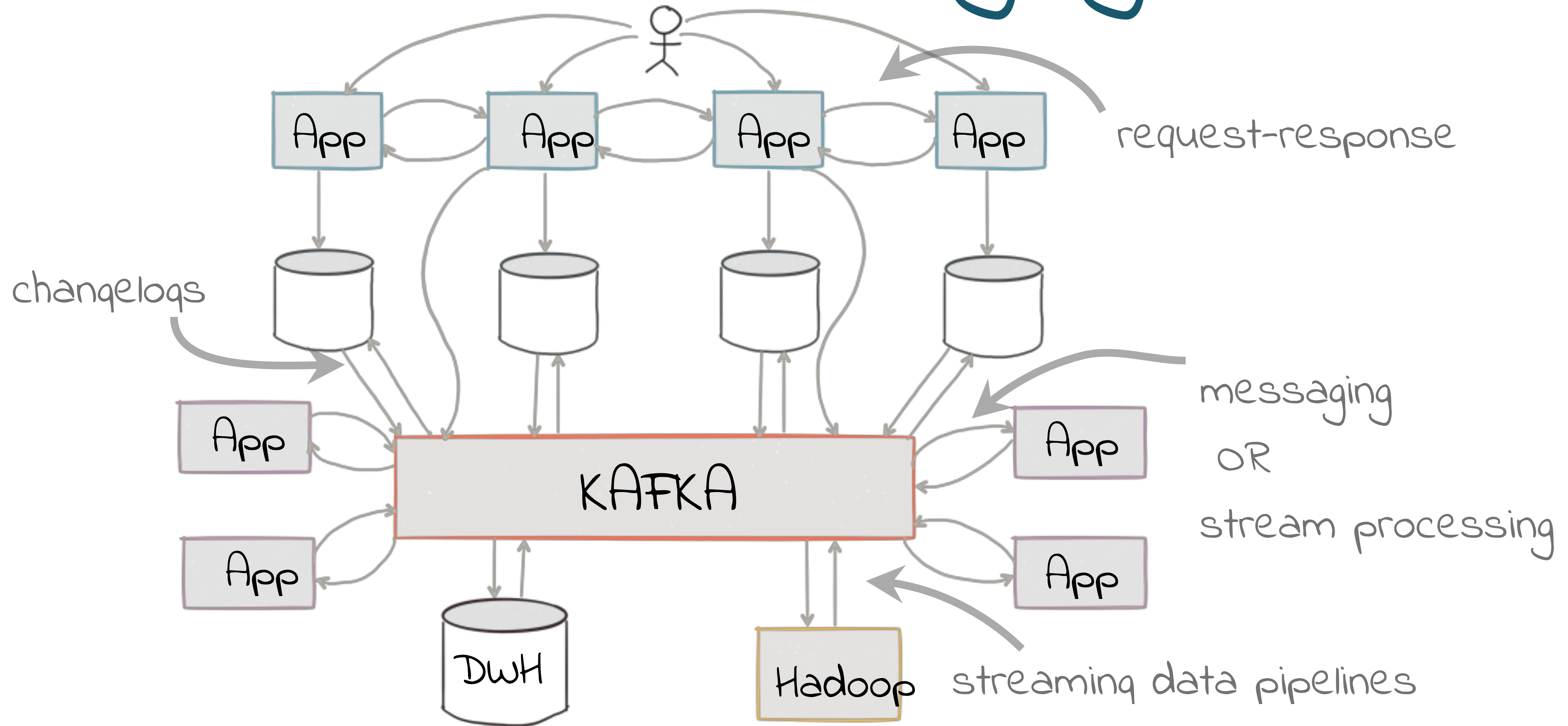
#Postgres

@confluentinc

A bit of a mess...



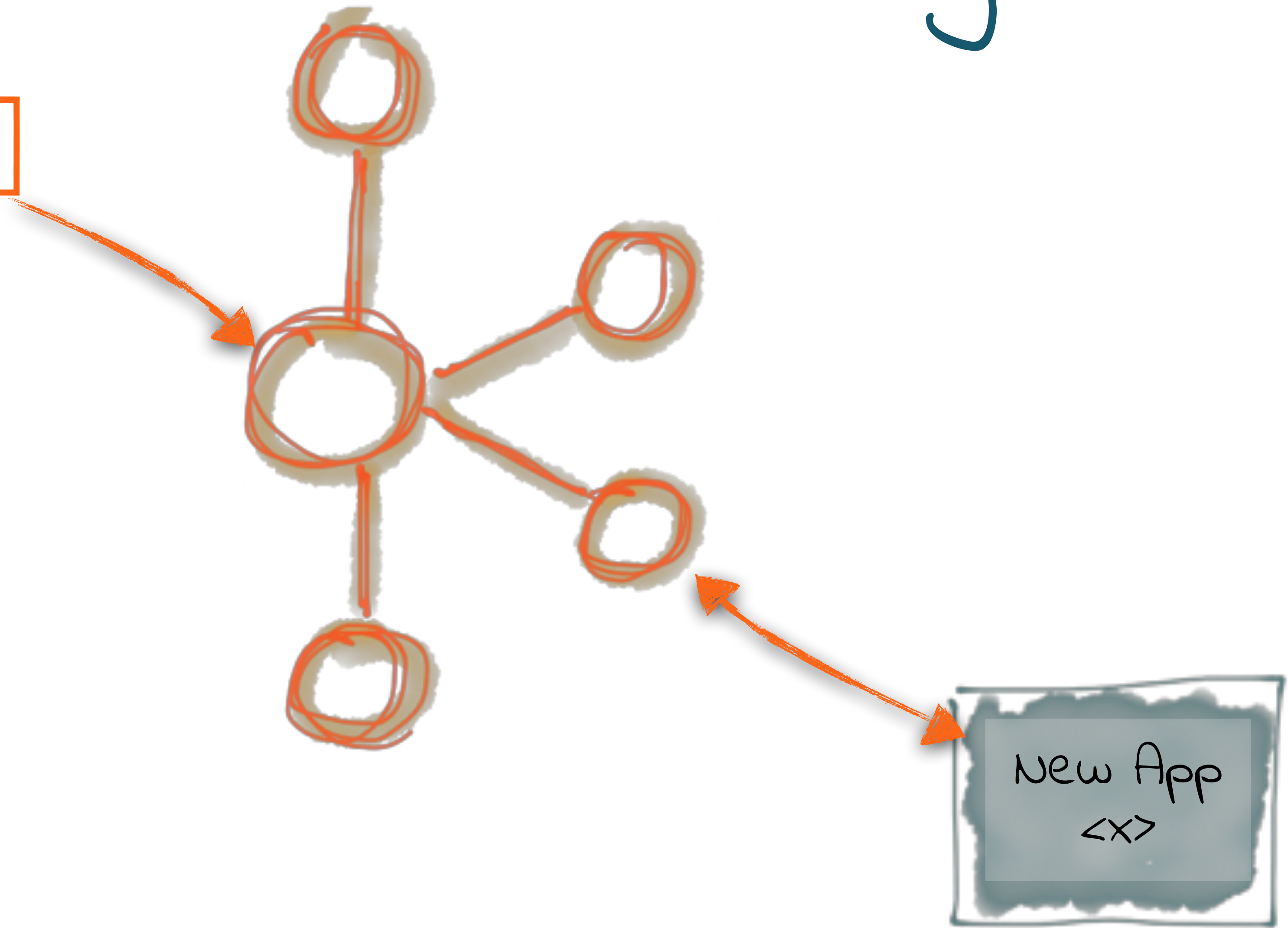
Streams are changing all of this



Message Bus



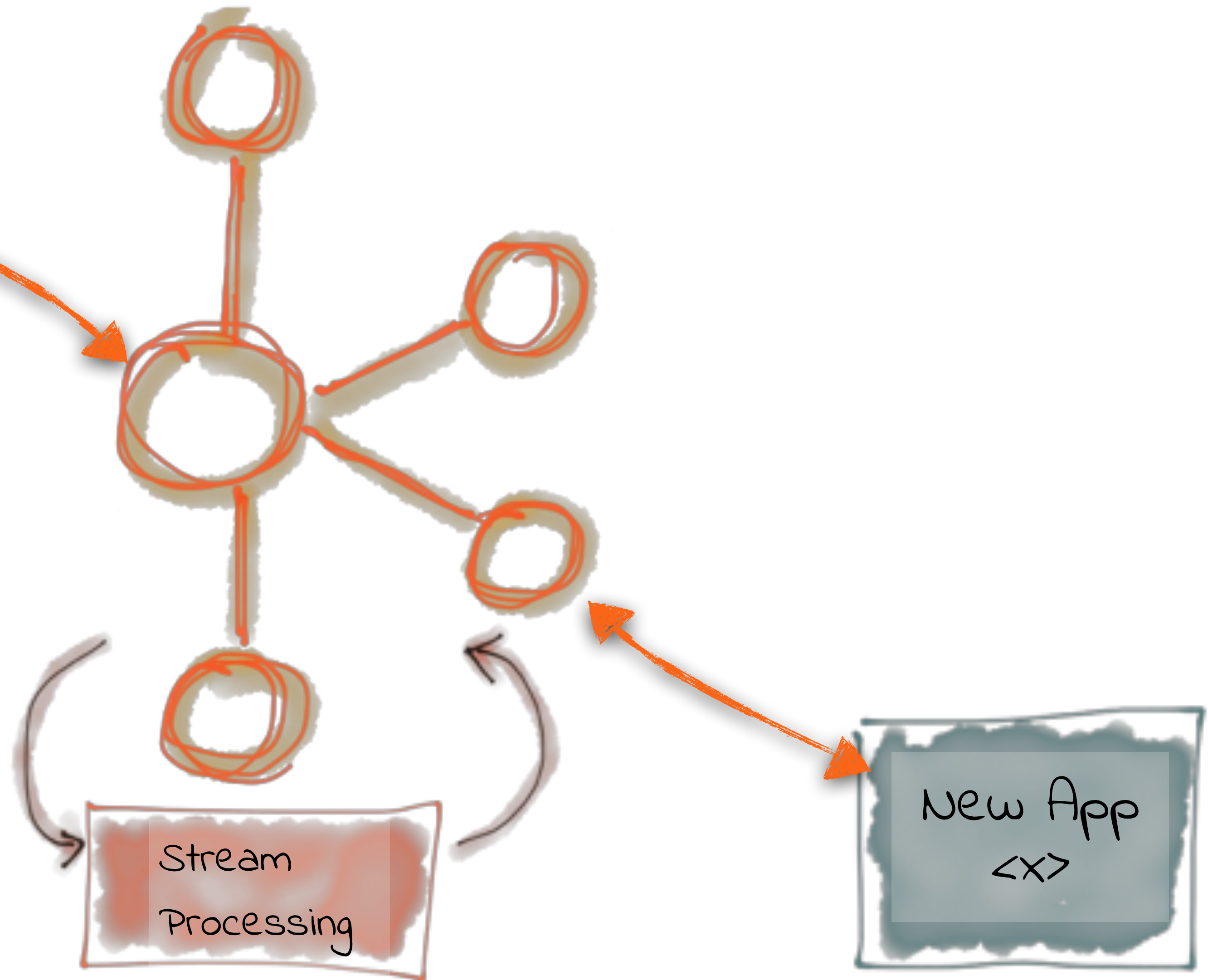
order events



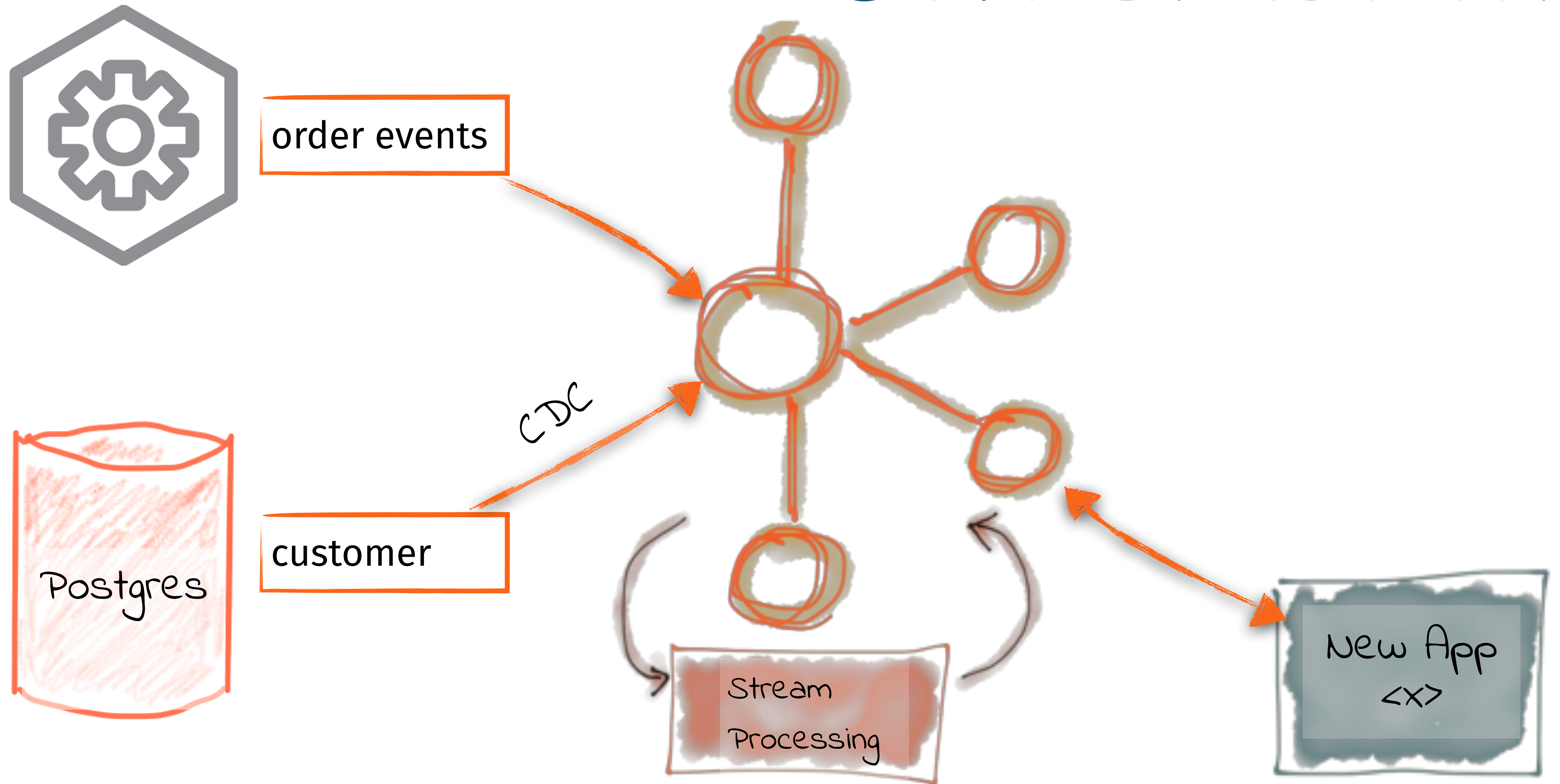
Stream Processing



order events



Data Enrichment



Transform Once, Use Many

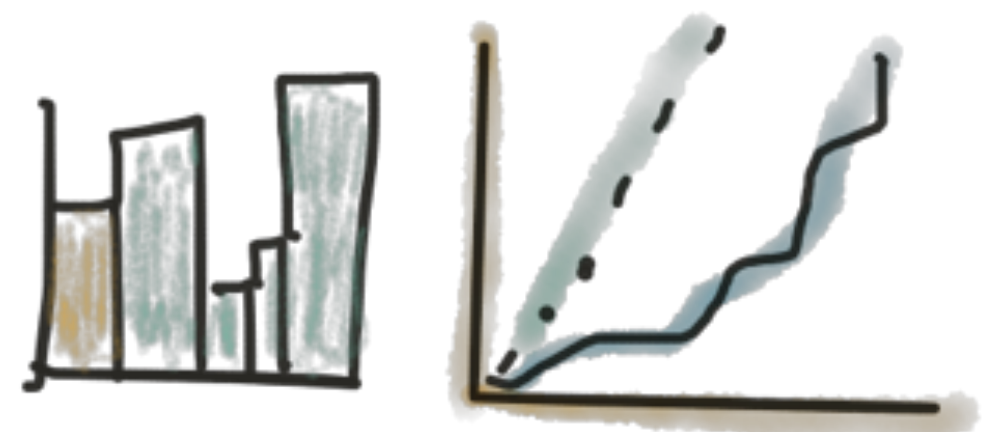


order events



customer

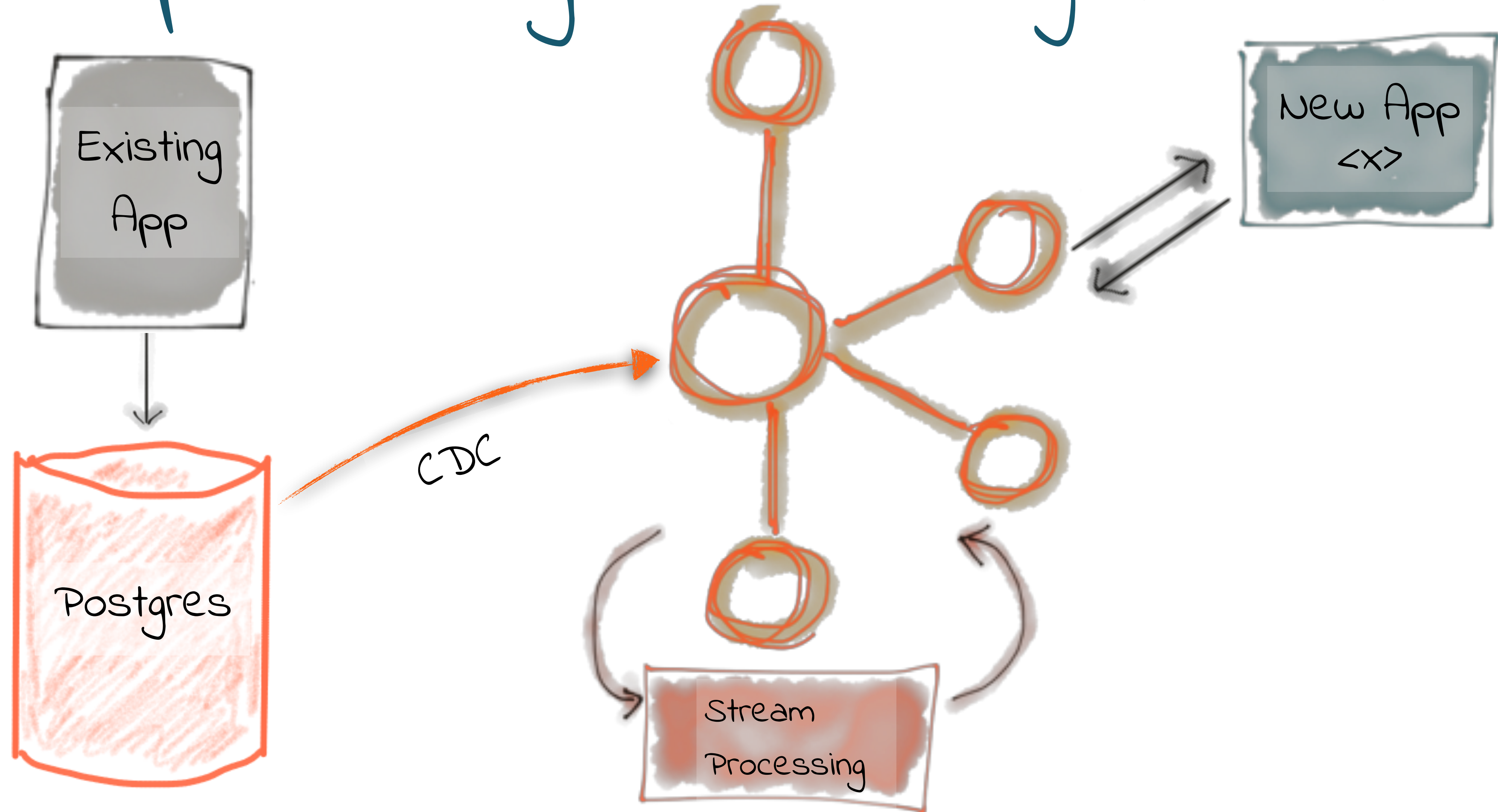
CDC



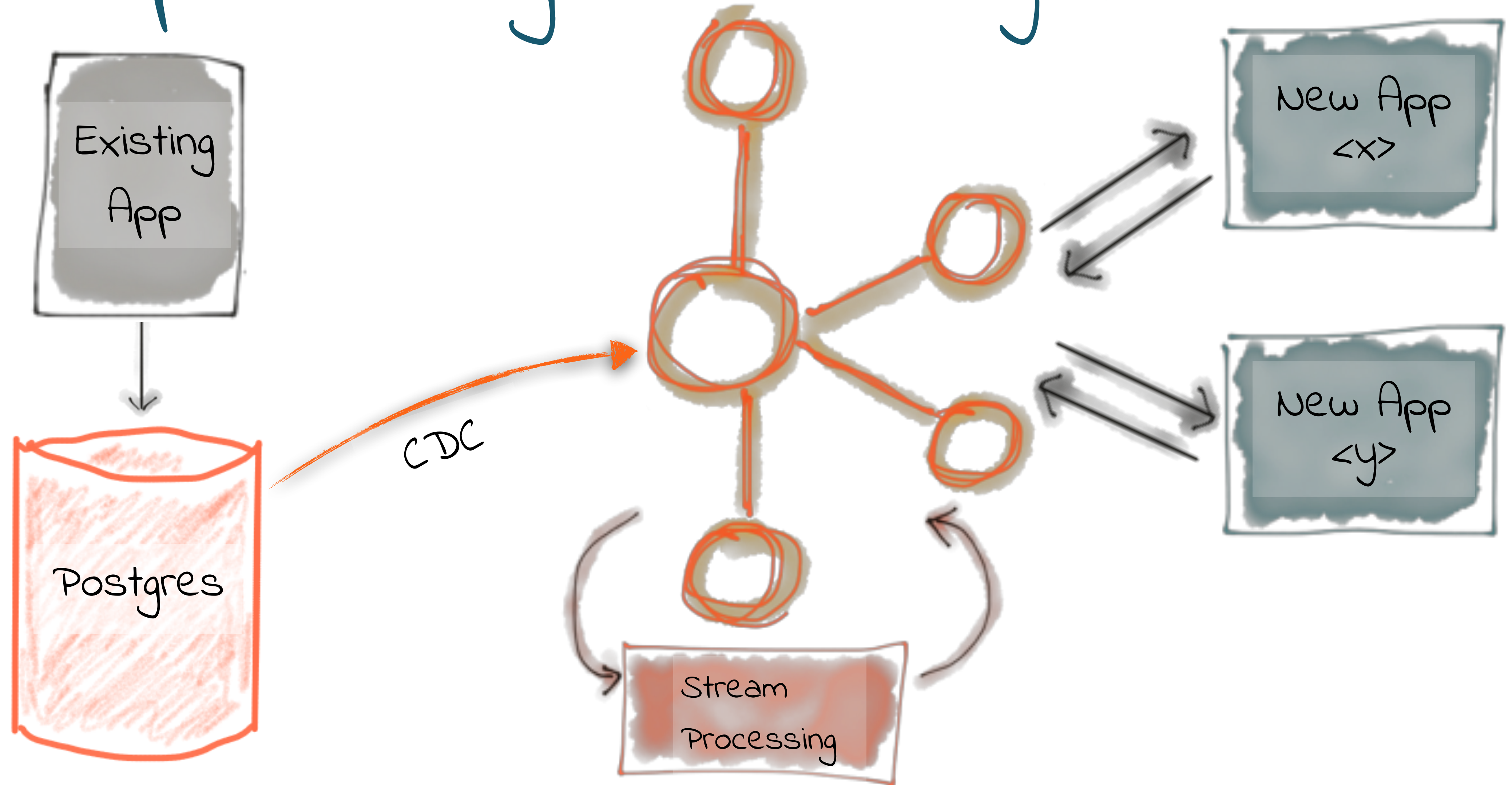
customer orders



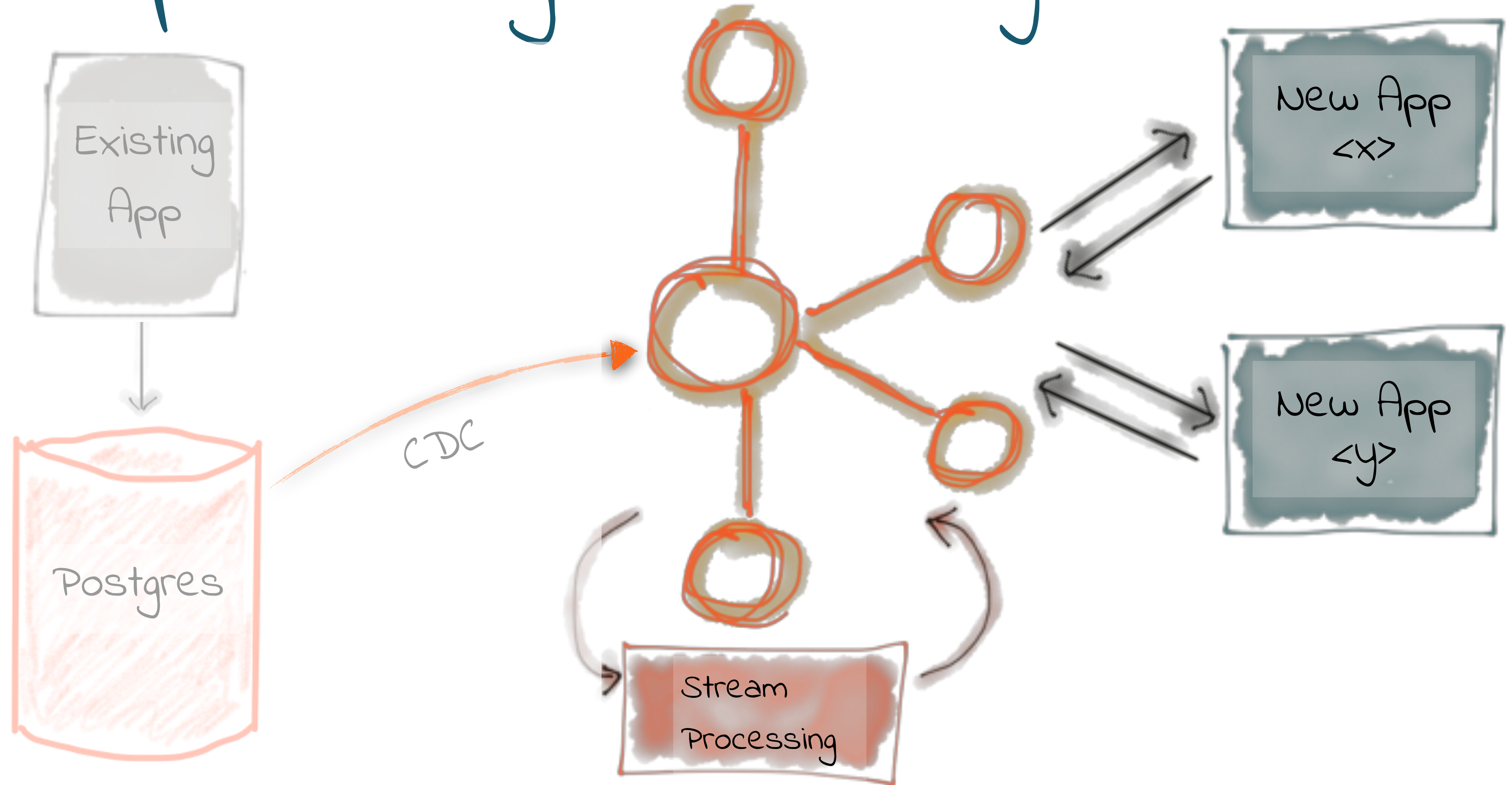
Evolve processing from old systems to new

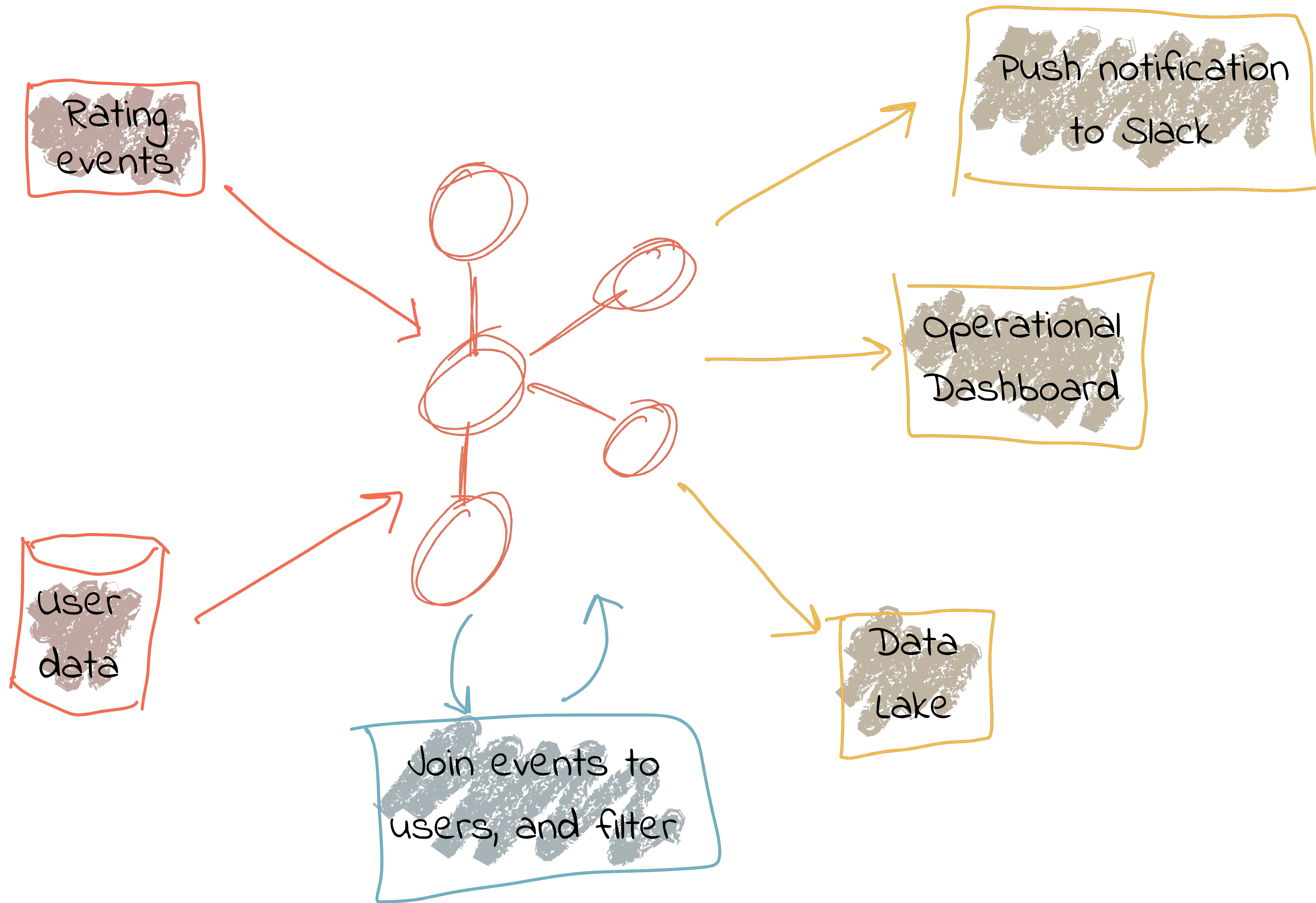


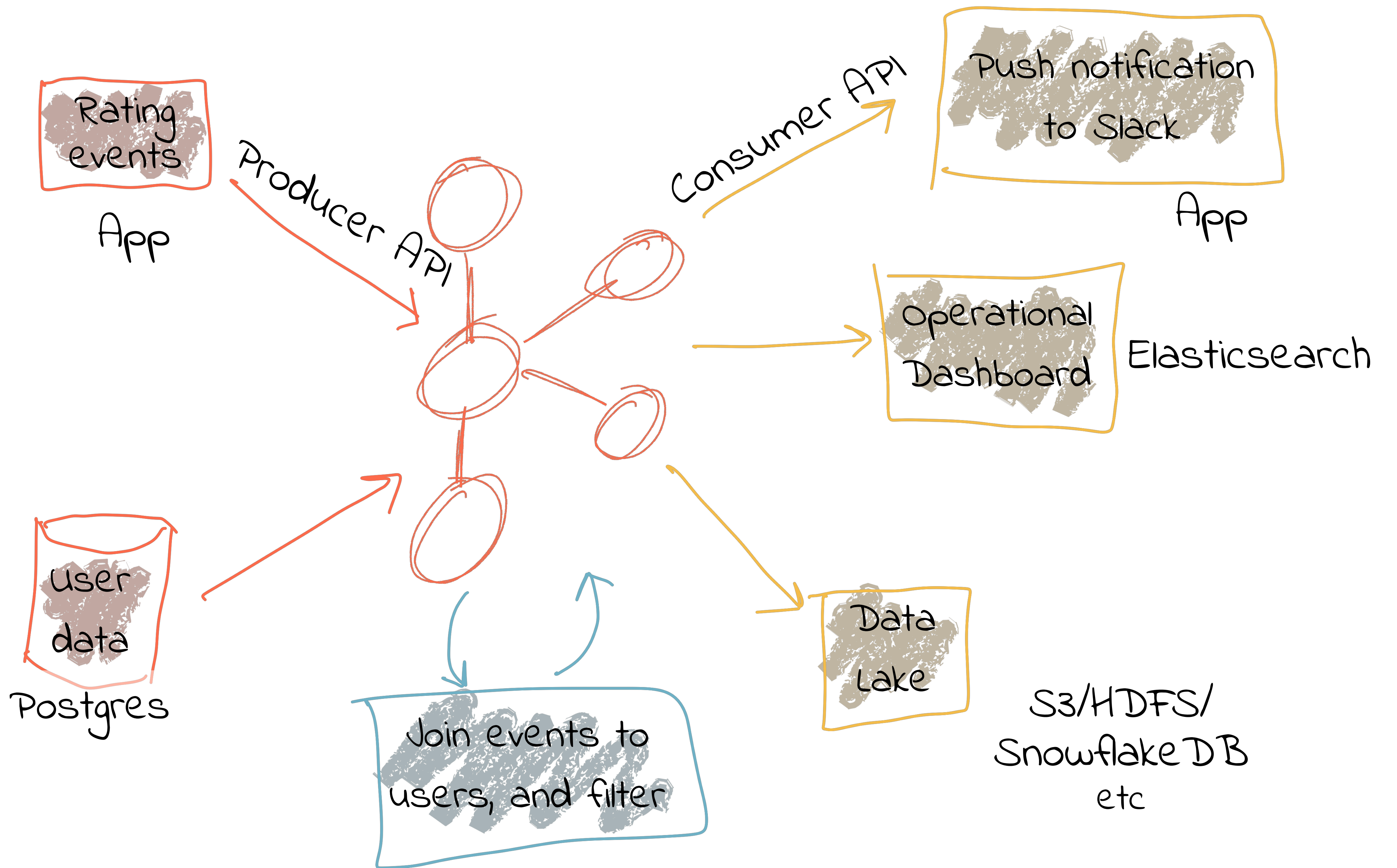
Evolve processing from old systems to new

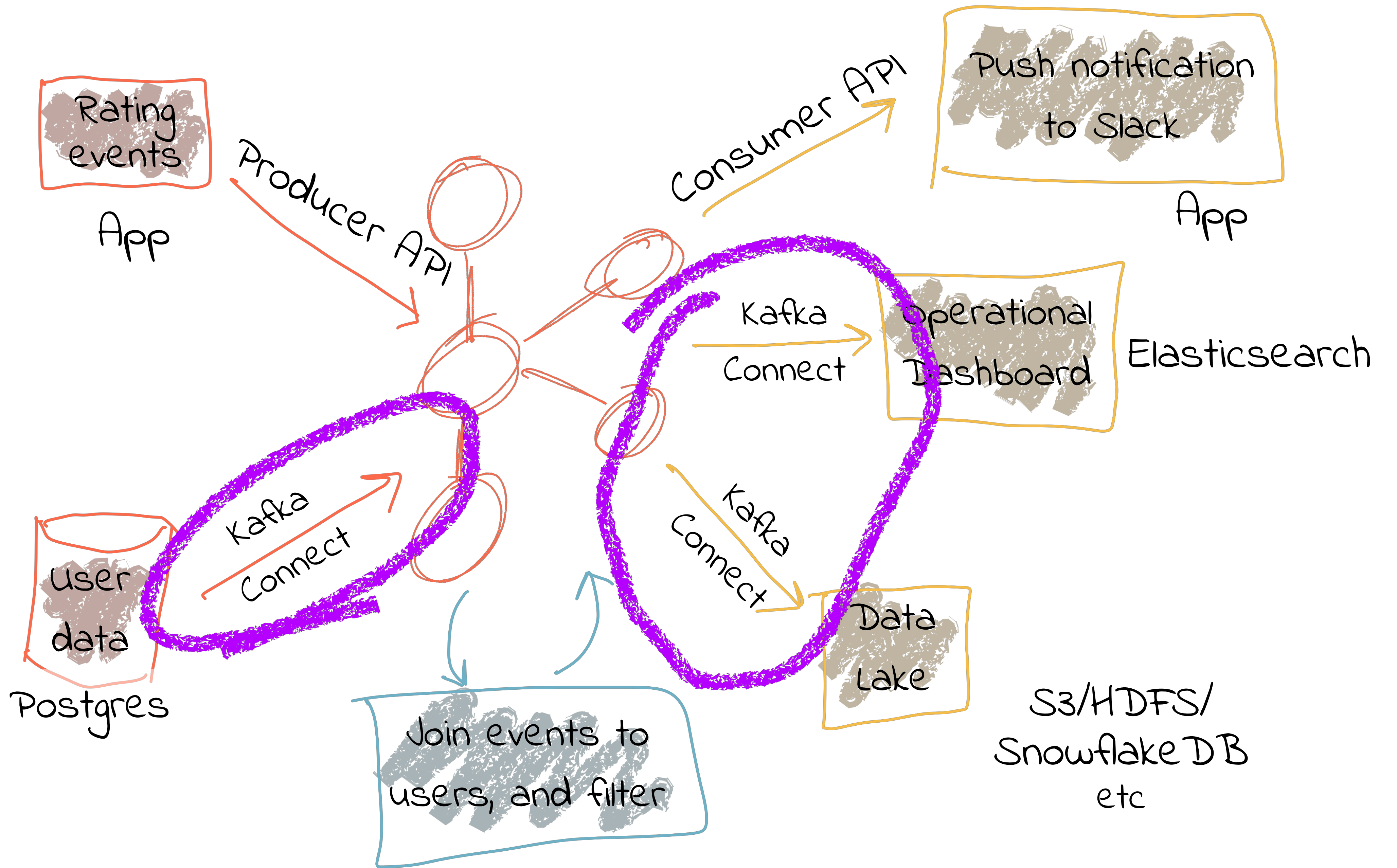


Evolve processing from old systems to new

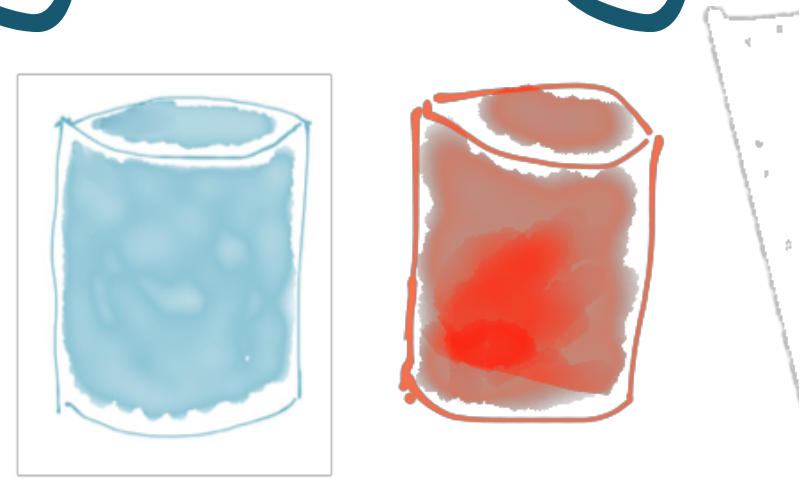




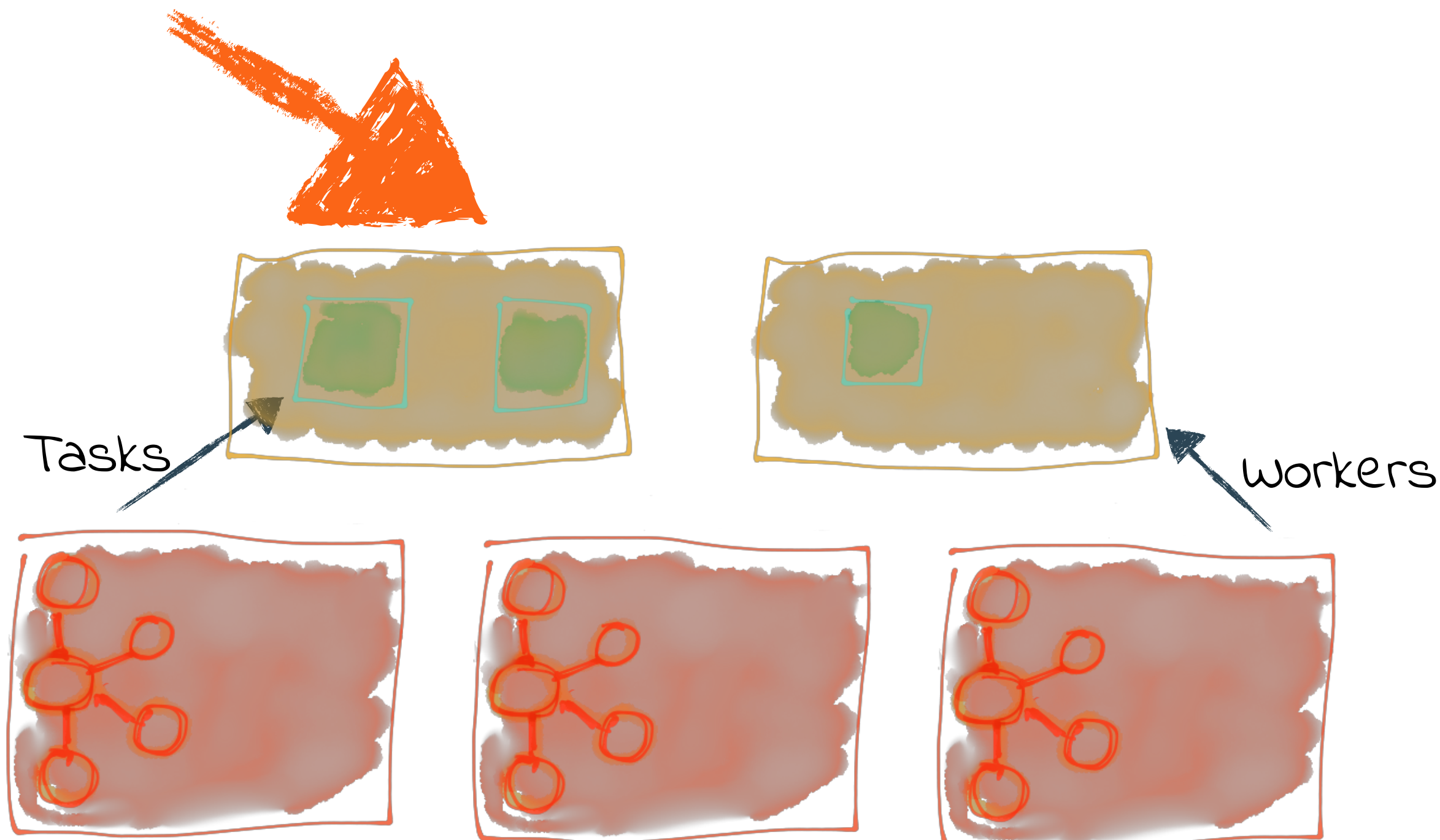




Streaming Integration with Kafka Connect



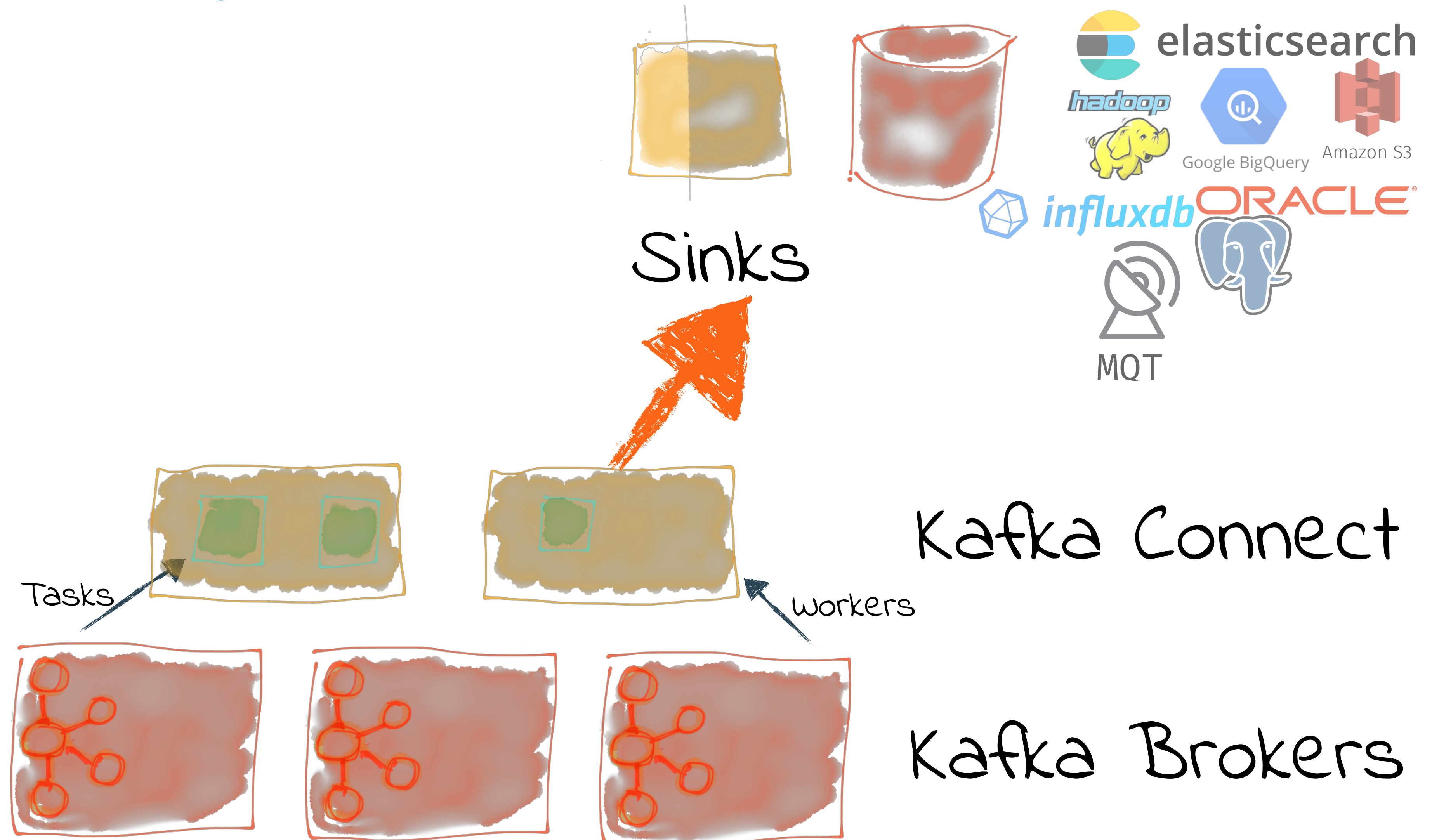
Sources



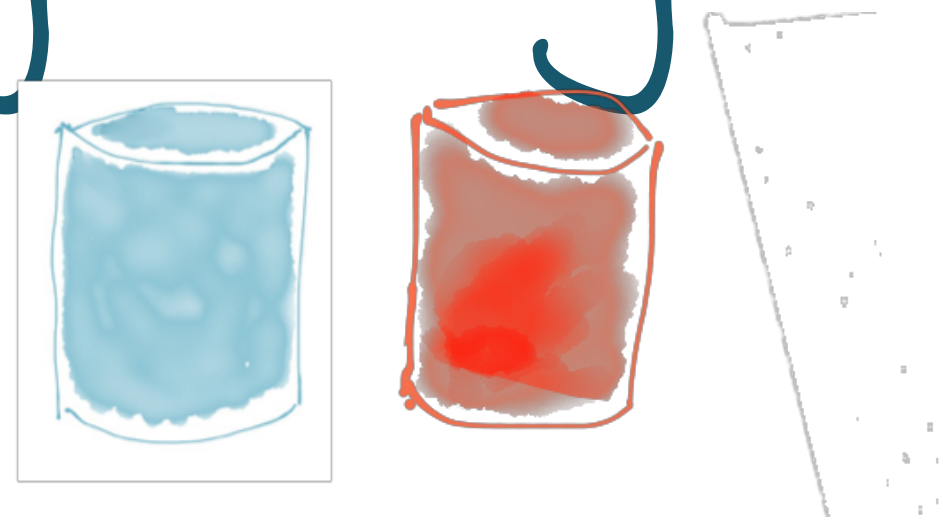
Kafka Connect

Kafka Brokers

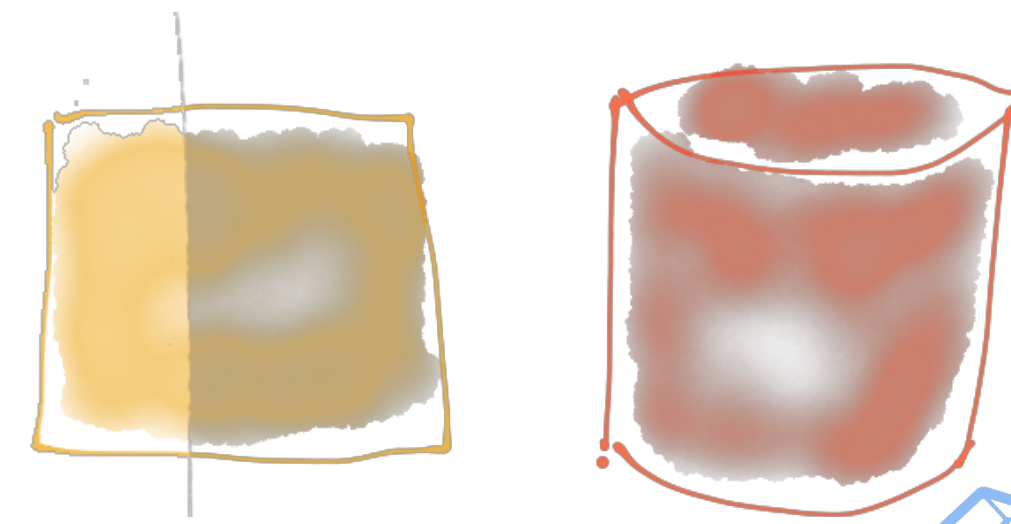
Streaming Integration with Kafka Connect



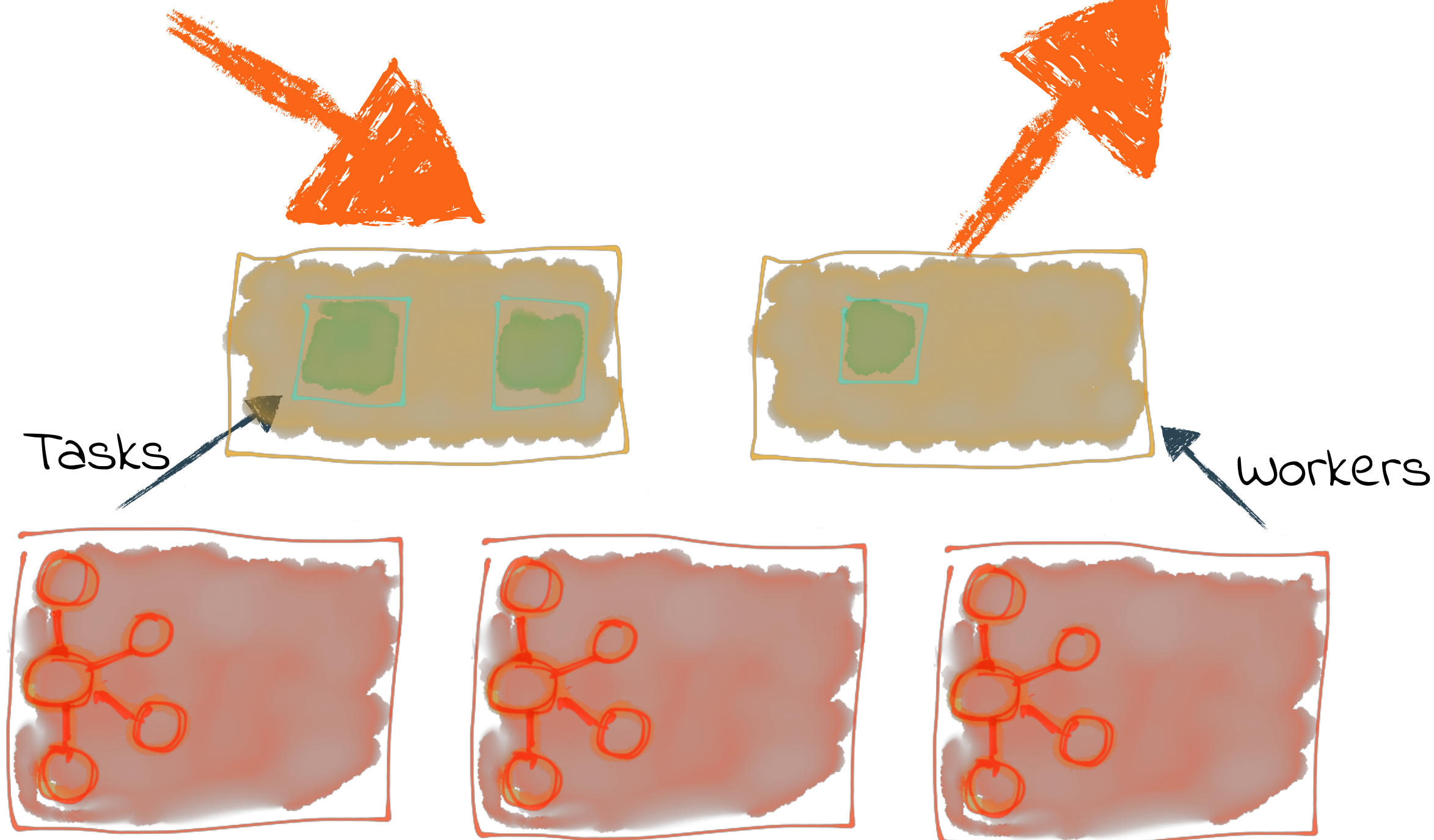
Streaming Integration with Kafka Connect



Sources



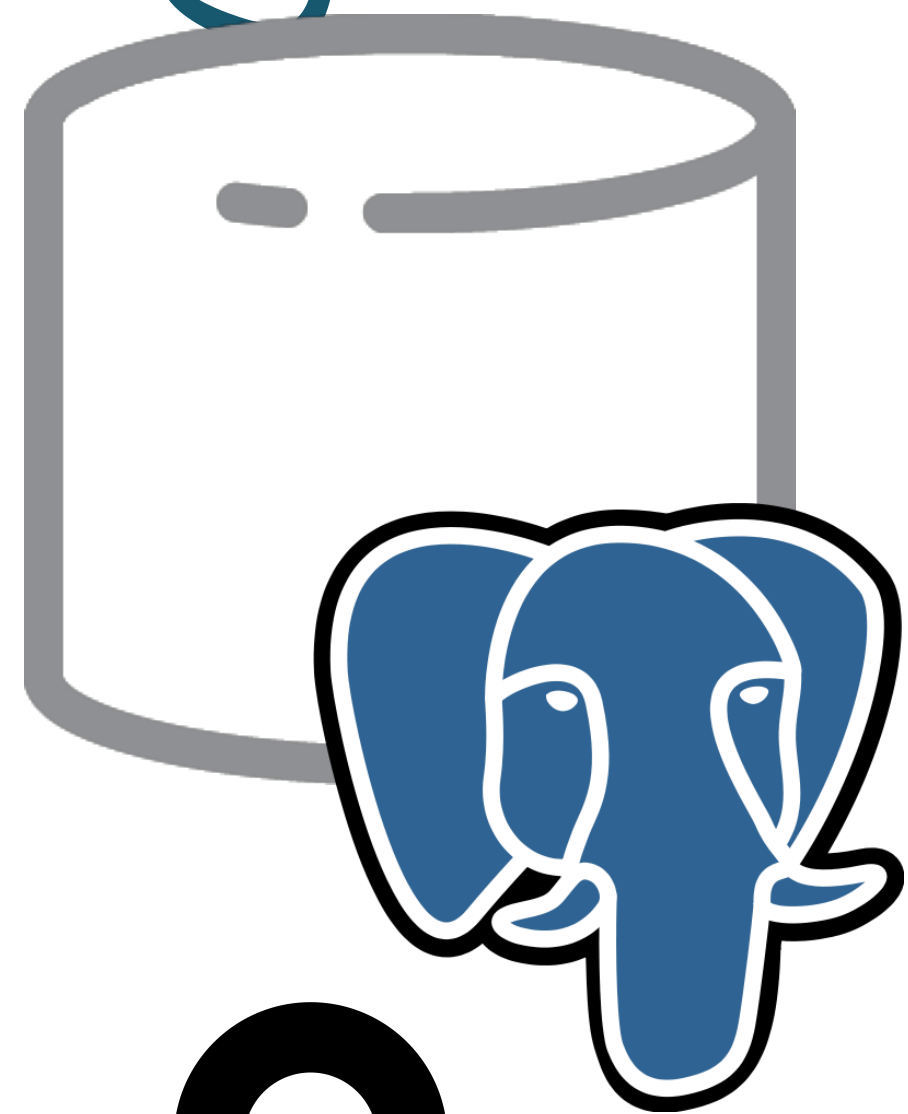
Sinks



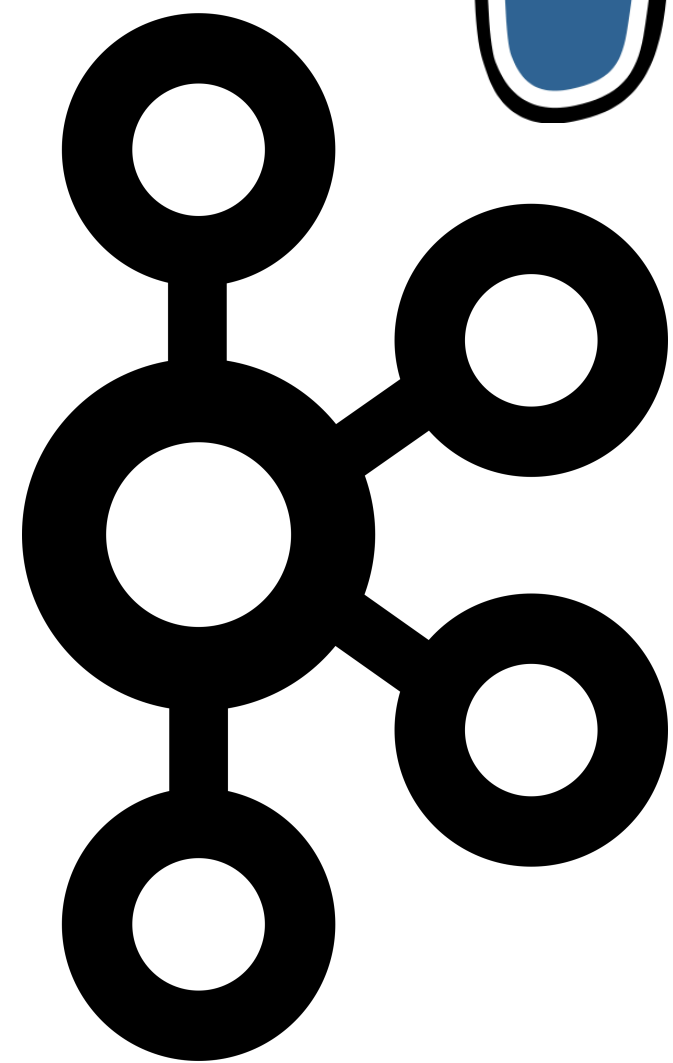
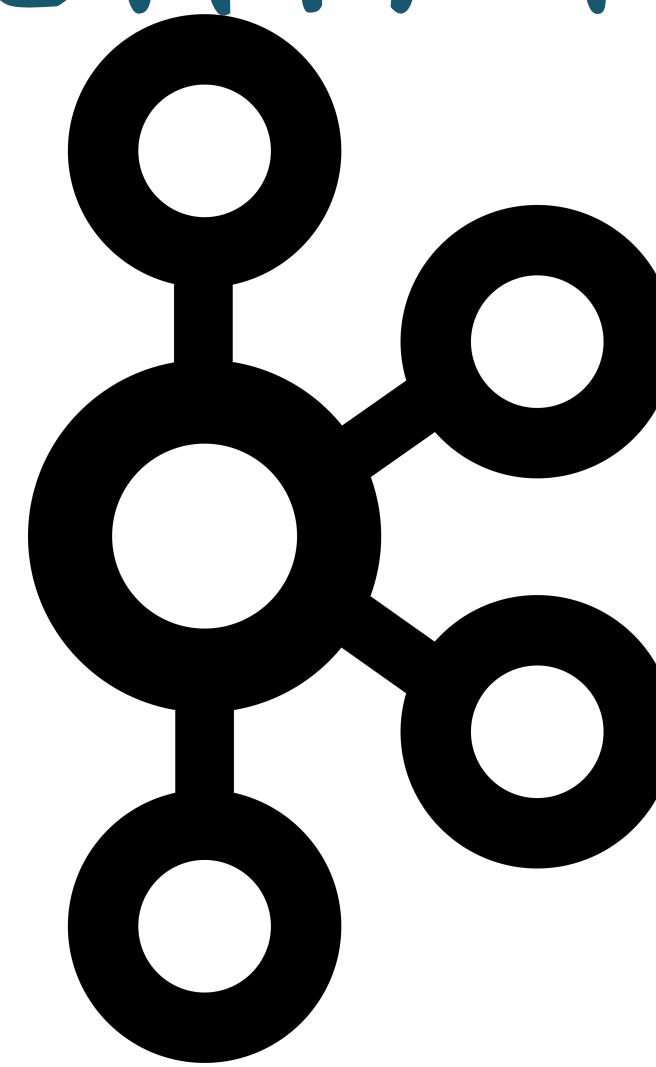
Kafka Connect

Kafka Brokers

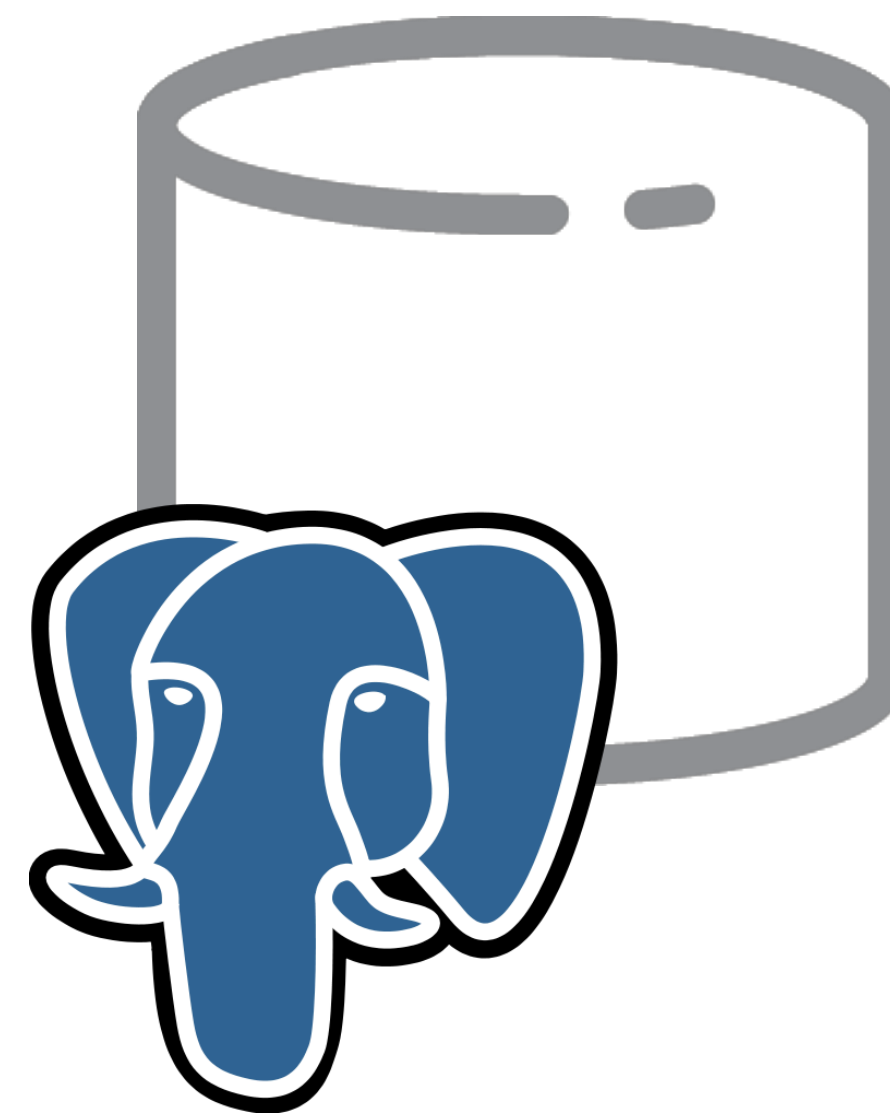
Integrating Postgres with Kafka

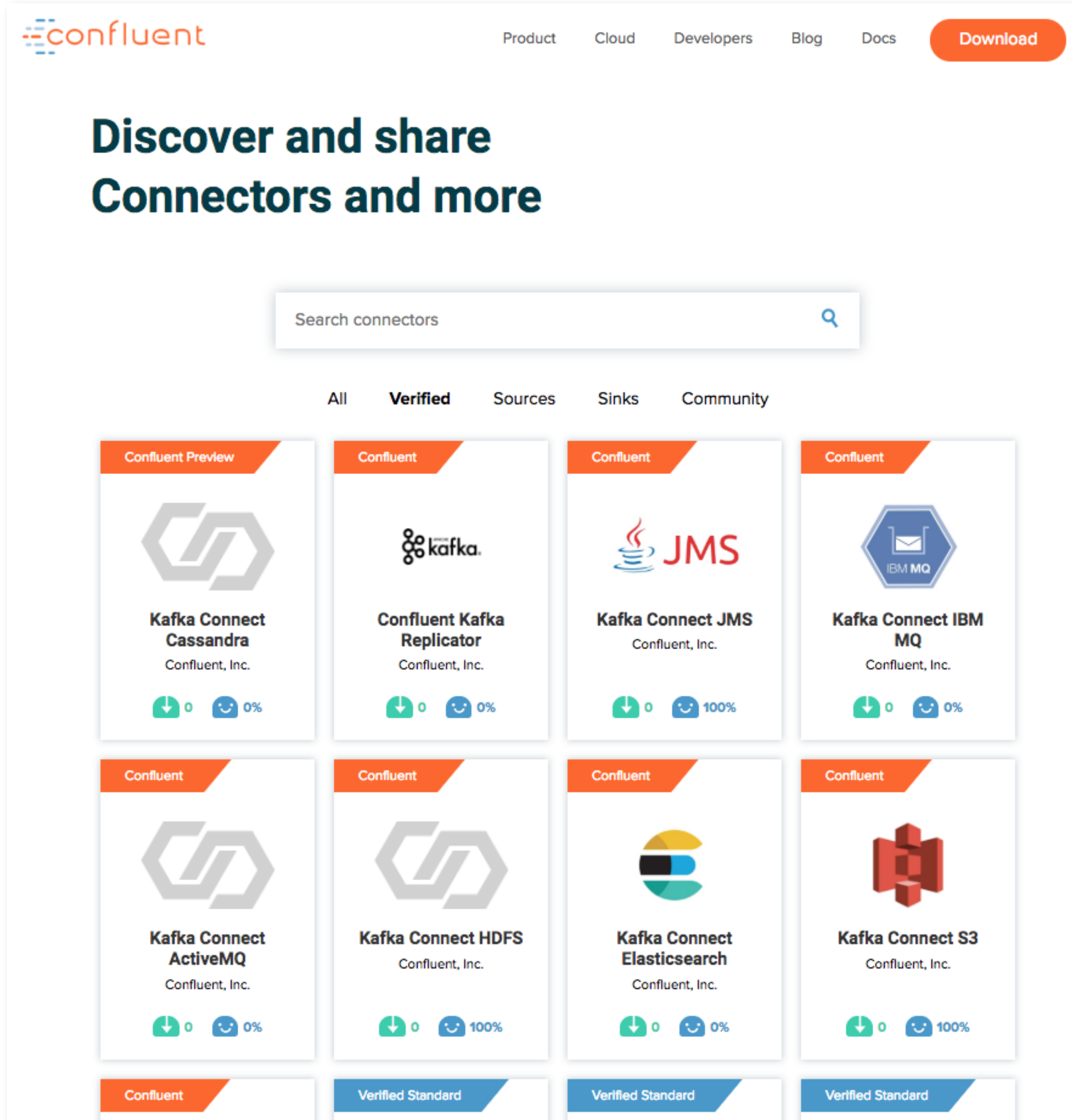


Kafka Connect &
Debezium



Kafka Connect &
JDBC Sink





Confluent Hub

• One-stop place to discover and download :

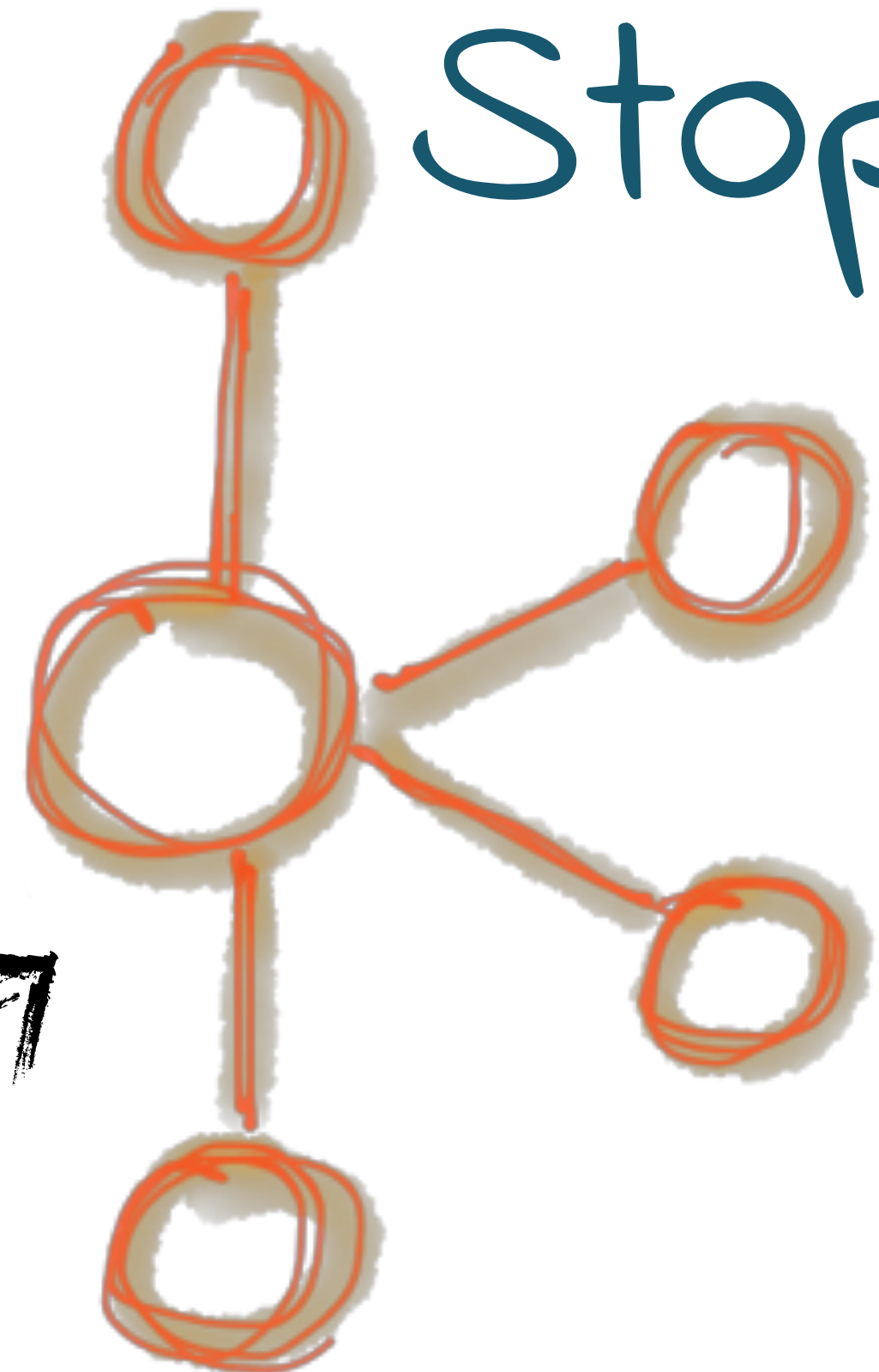
- Connectors
- Transformations
- Converters



hub.confluent.io


```
{
  "rating_id": 2133,
  "user_id": 9,
  "stars": 1,
  "route_id": 7219,
  "rating_time": 1519402815063,
  "channel": "web",
  "message": "worst. flight. ever. #neveragain"
}
```

Stop! Demo time!



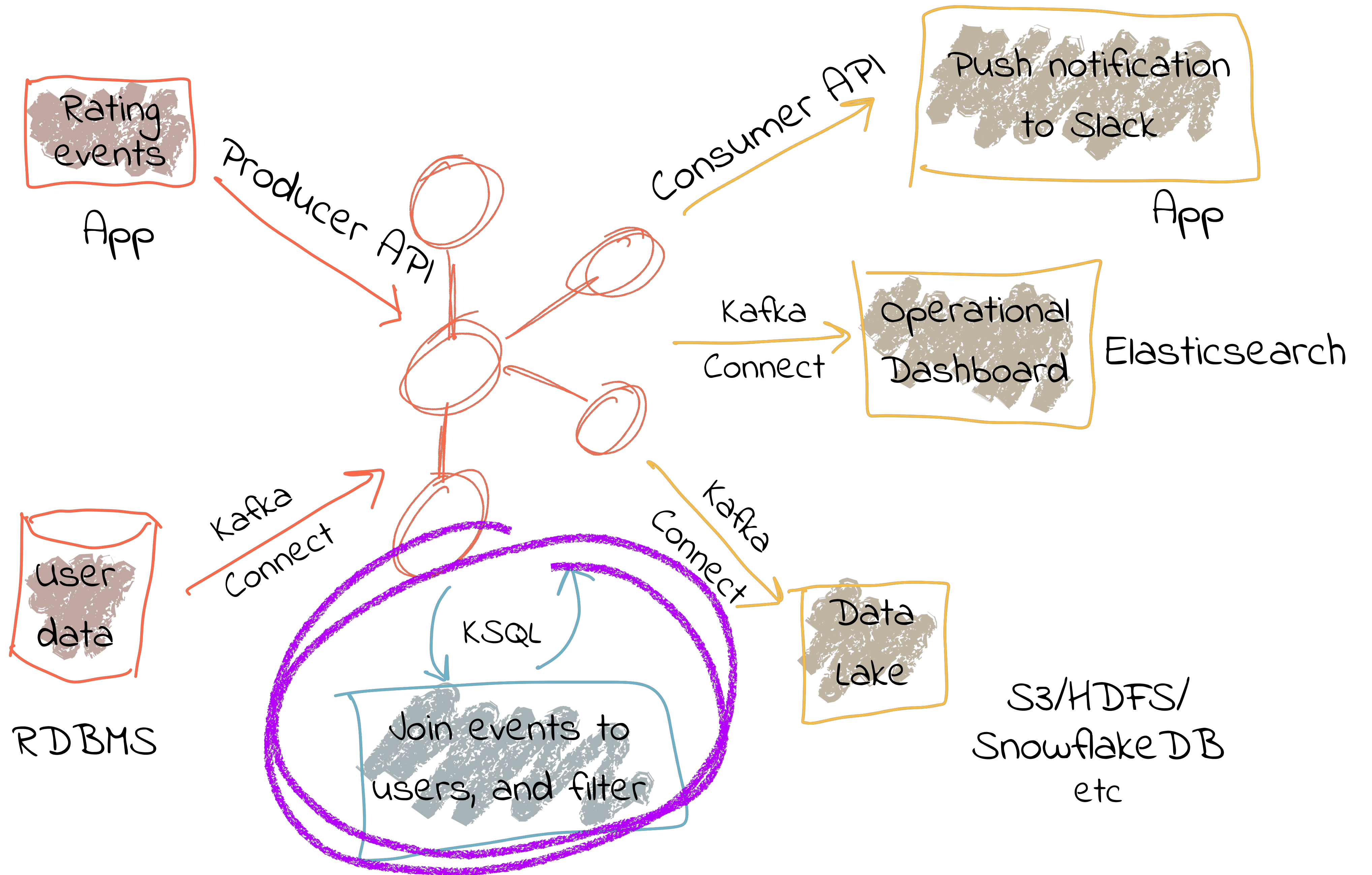
Producer API

Kafka Connect
Debezium

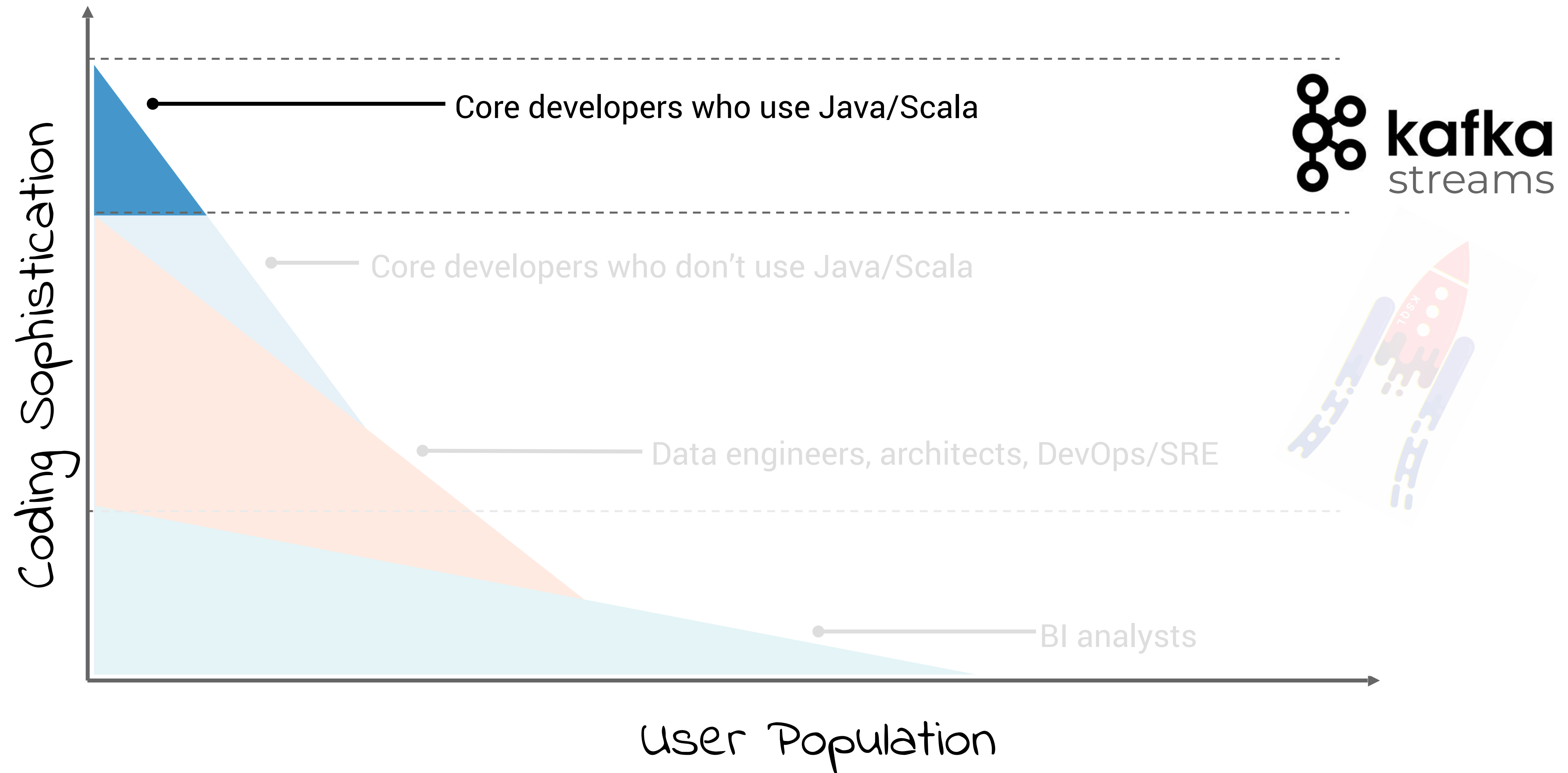


Postgres

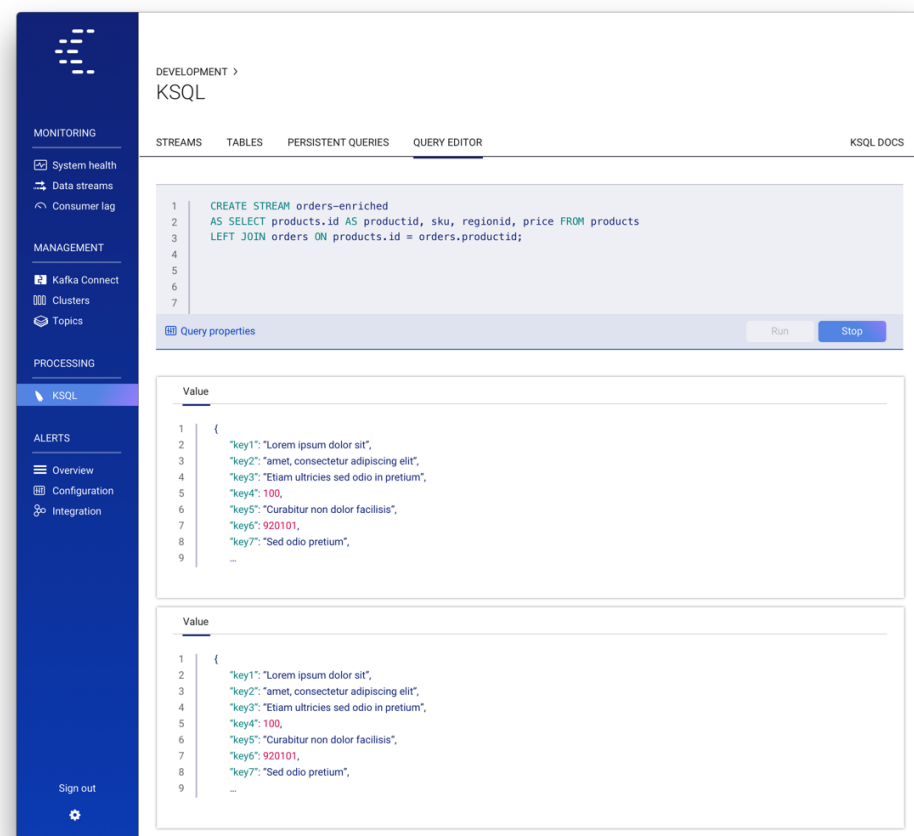
```
{
  "uid": 9,
  "name": "Neha",
  "locale": "en_US",
  "address_city": "Palo Alto",
  "elite": "P"
}
```



Lower the bar to enter the world of streaming



KSQL #FTW



1 UI



2 CLI



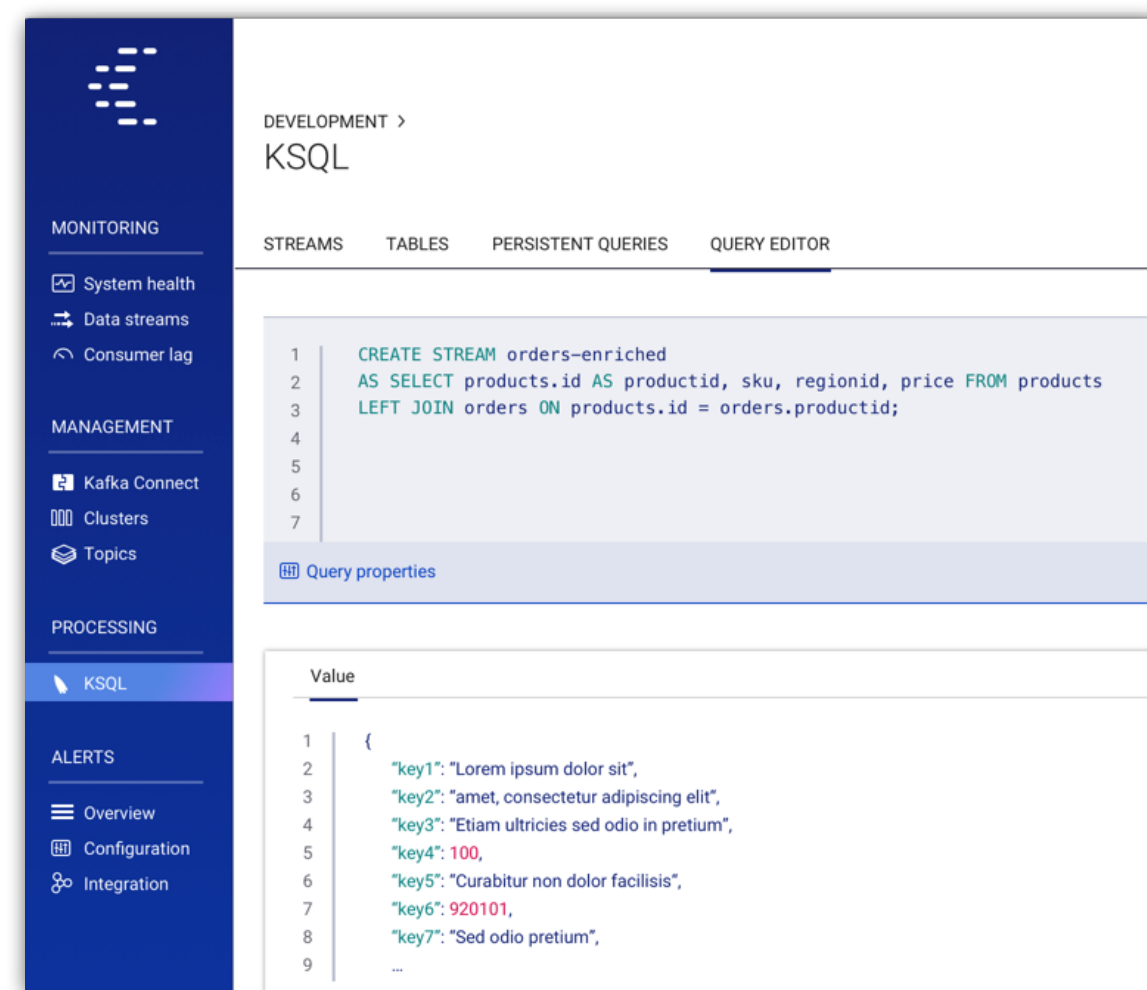
3 REST



4 Headless

Interactive KSQL for development and testing

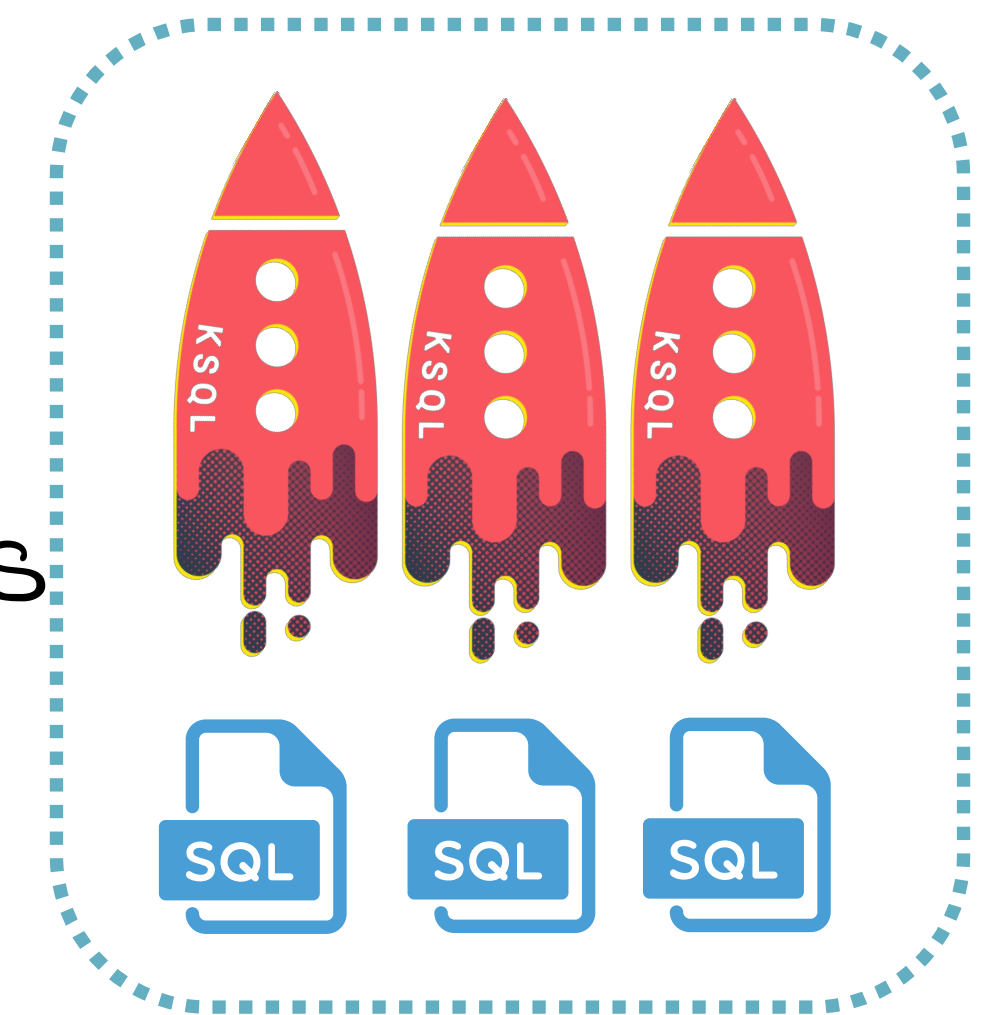
Headless KSQL for Production



REST



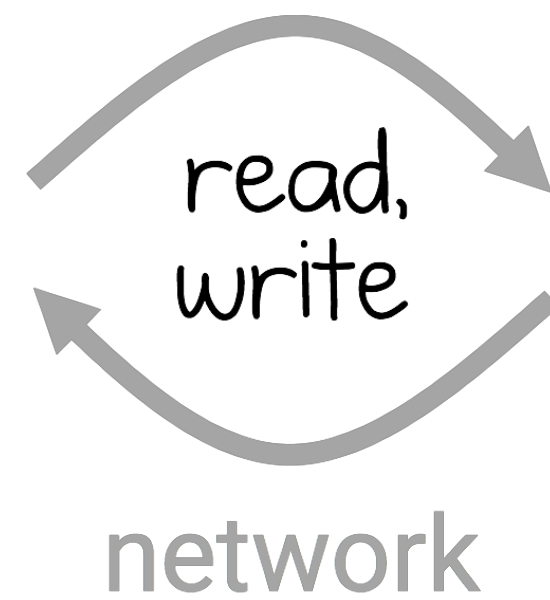
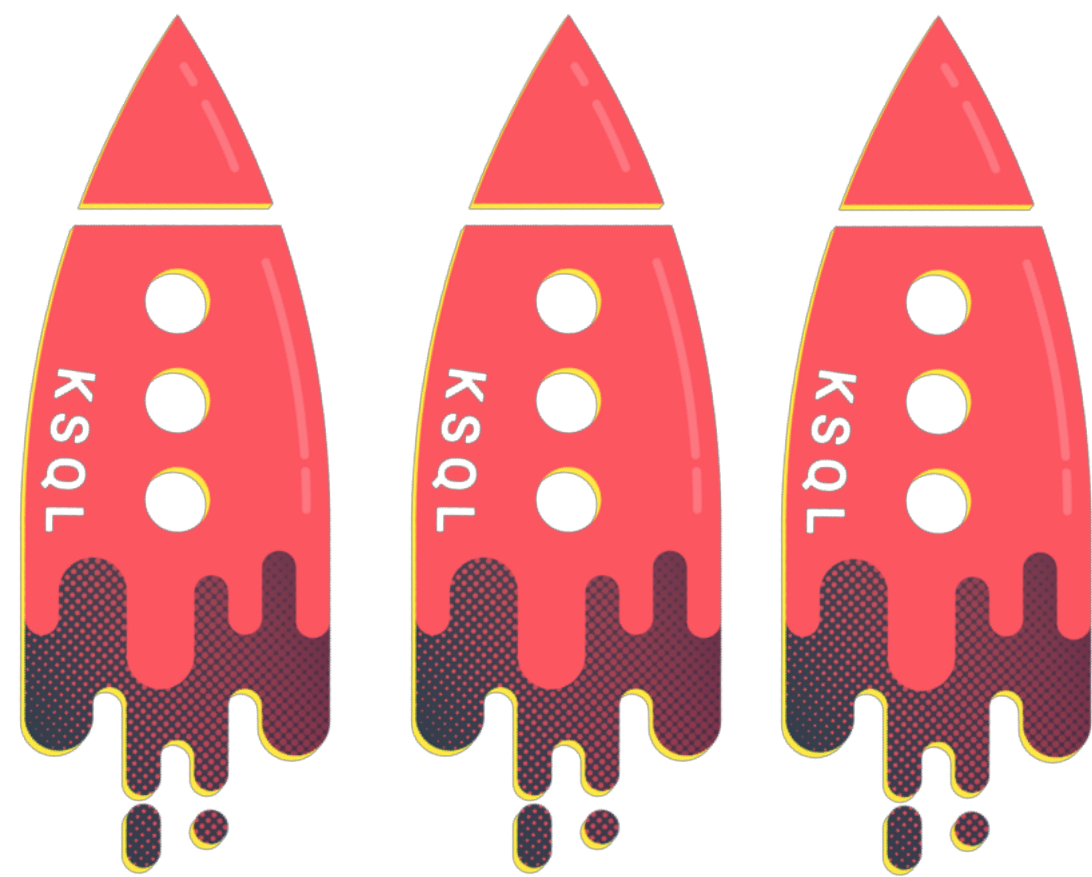
Desired KSQL queries
have been identified



"Hmm, let me try
out this idea.."

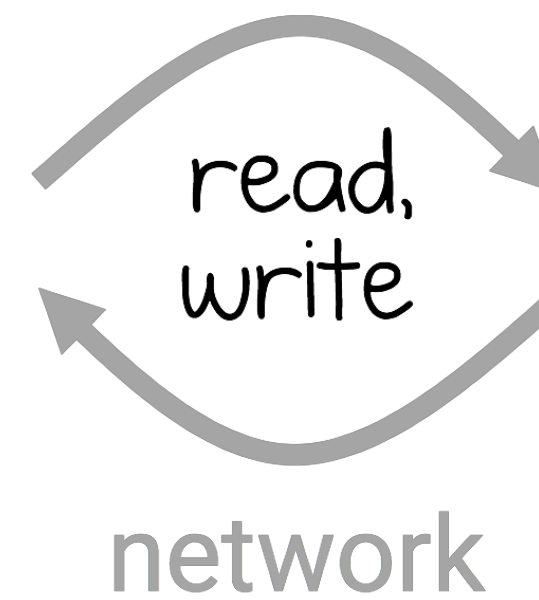
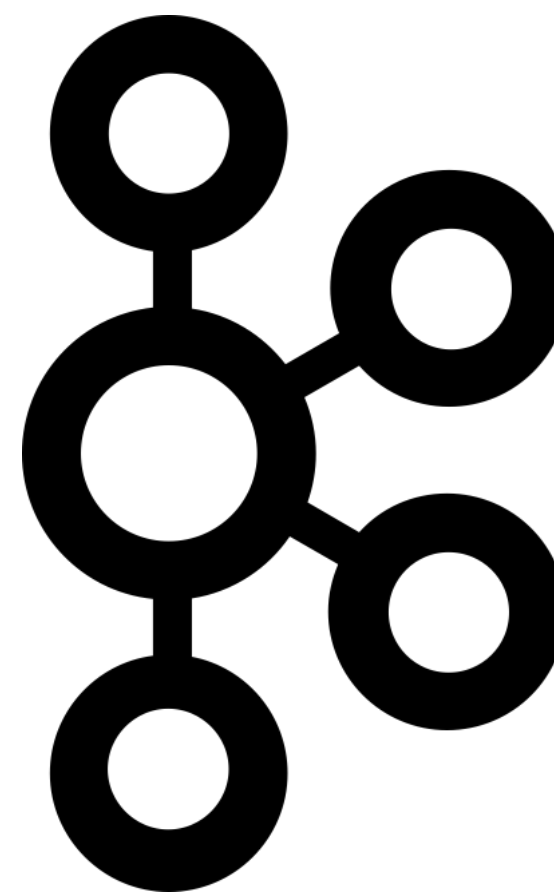
Interaction with Kafka

KSQL
(processing)



Does not run on
Kafka brokers

Kafka
(data)



JVM application
with Kafka Streams (processing)

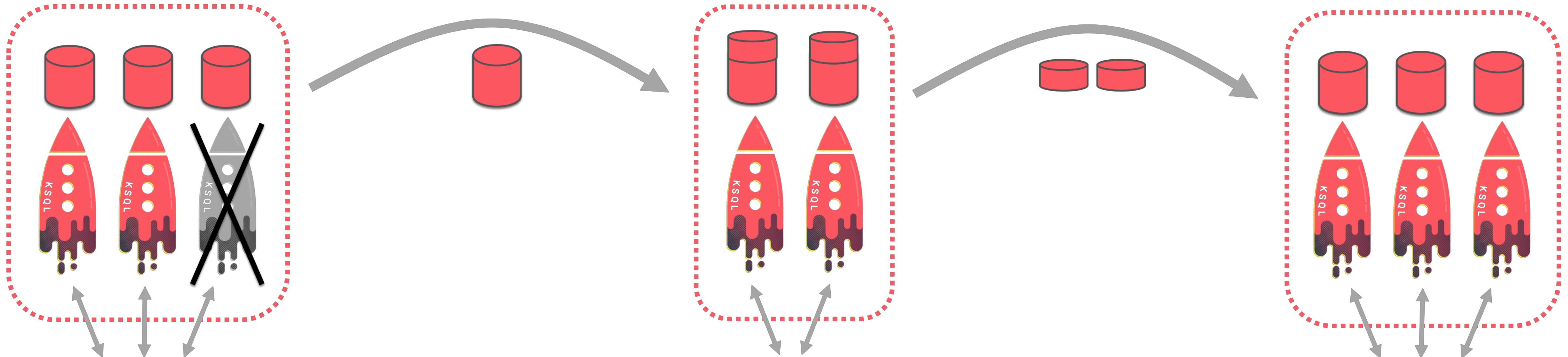


Does not run on
Kafka brokers

Fault-Tolerance, powered by Kafka

#3 died so #1 and #2 take over

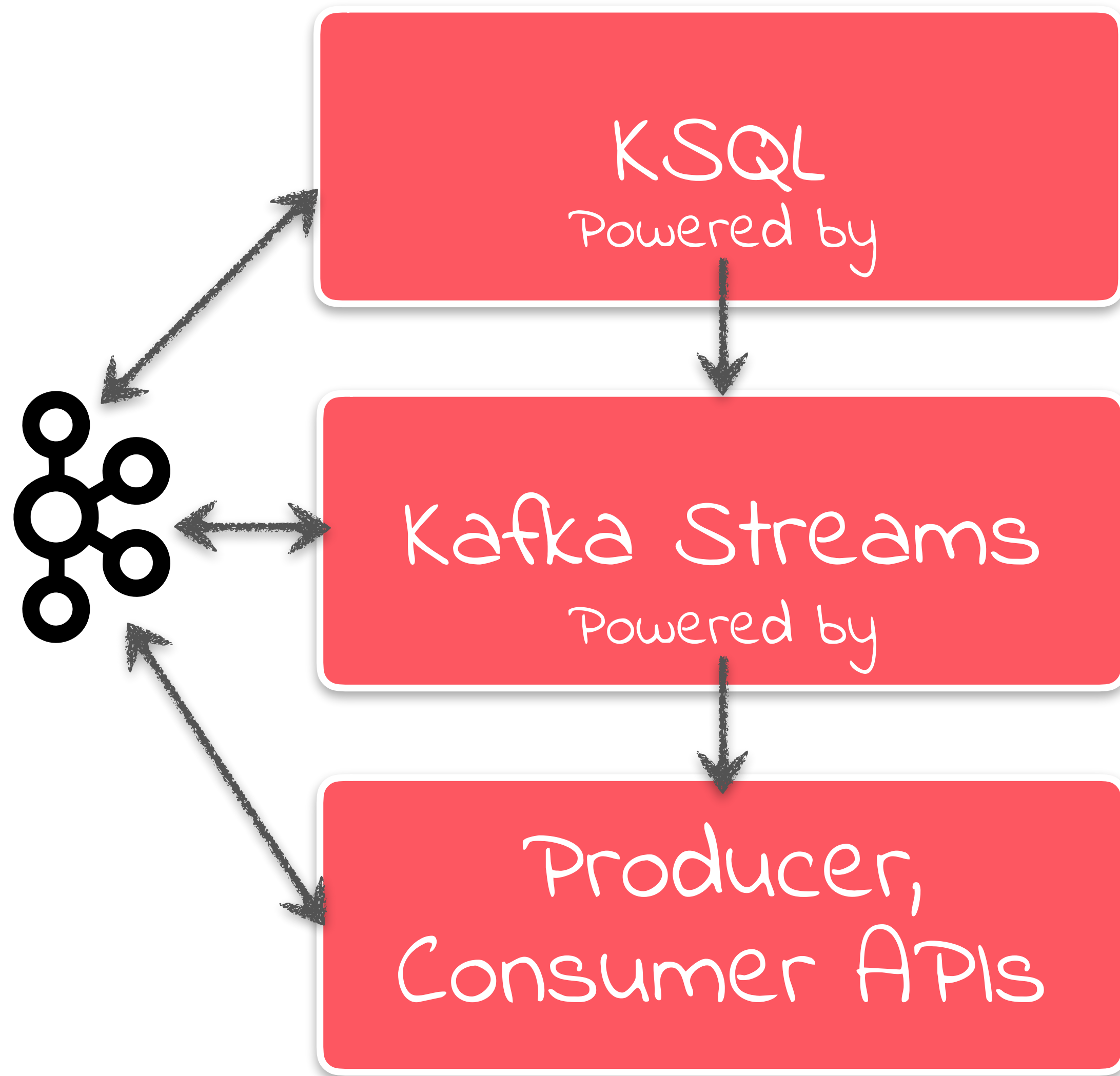
#3 is back so the work is split again



- 1 Kafka consumer group rebalance is triggered
- 2 Processing and **state** of #3 is migrated via Kafka to remaining servers #1 + #2

- 1 Kafka consumer group rebalance is triggered
- 2 Part of processing incl. **state** is migrated via Kafka from #1 + #2 to server #3

Standing on the shoulders of Streaming Giants



Ease of use



Flexibility

```
CREATE STREAM,  
CREATE TABLE, SELECT, JOIN,  
GROUP BY, SUM, ...
```

KSQL UDFs

```
KStream<>, KTable<>, filter(), map(),  
flatMap(), join(), aggregate(),  
transform(), ...
```

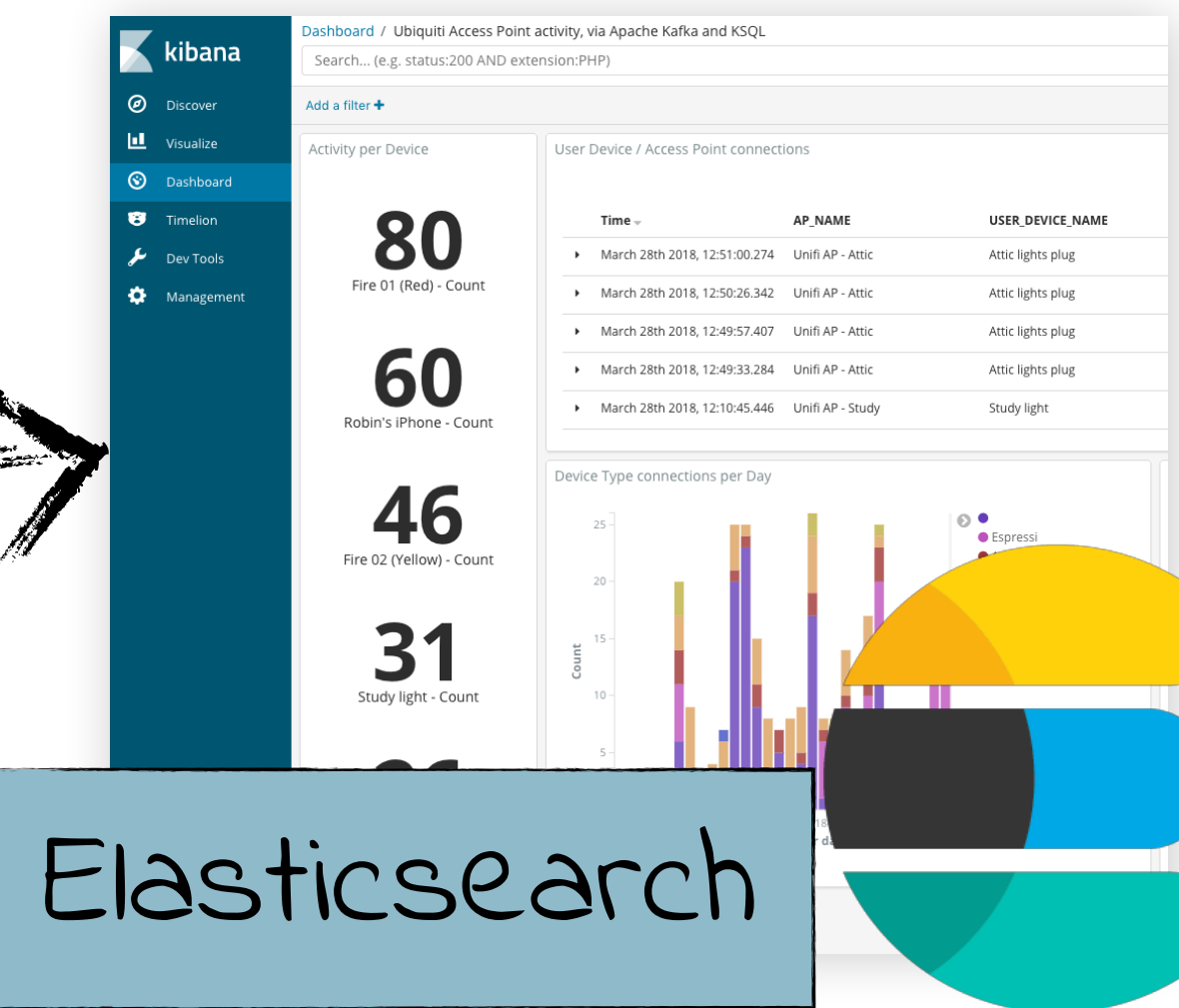
```
subscribe(), poll(), send(),  
flush(), beginTransaction(), ...
```


Demo Time!

```
{  
  "rating_id": 5313,  
  "user_id": 3,  
  "stars": 4,  
  "route_id": 6975,  
  "rating_time": 1519304105213,  
  "channel": "web",  
  "message": "worst. flight. ever. #neveragain"  
}
```

Producer API

Kafka Connect



Elasticsearch

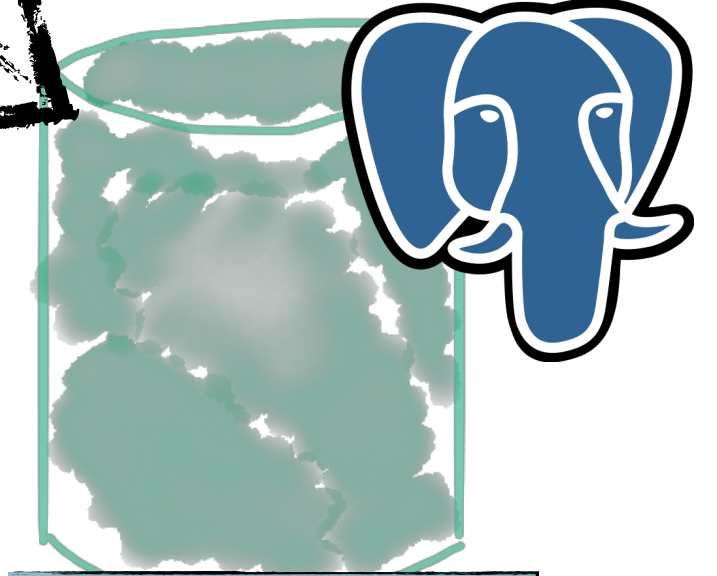
Kafka Connect

Kafka Connect

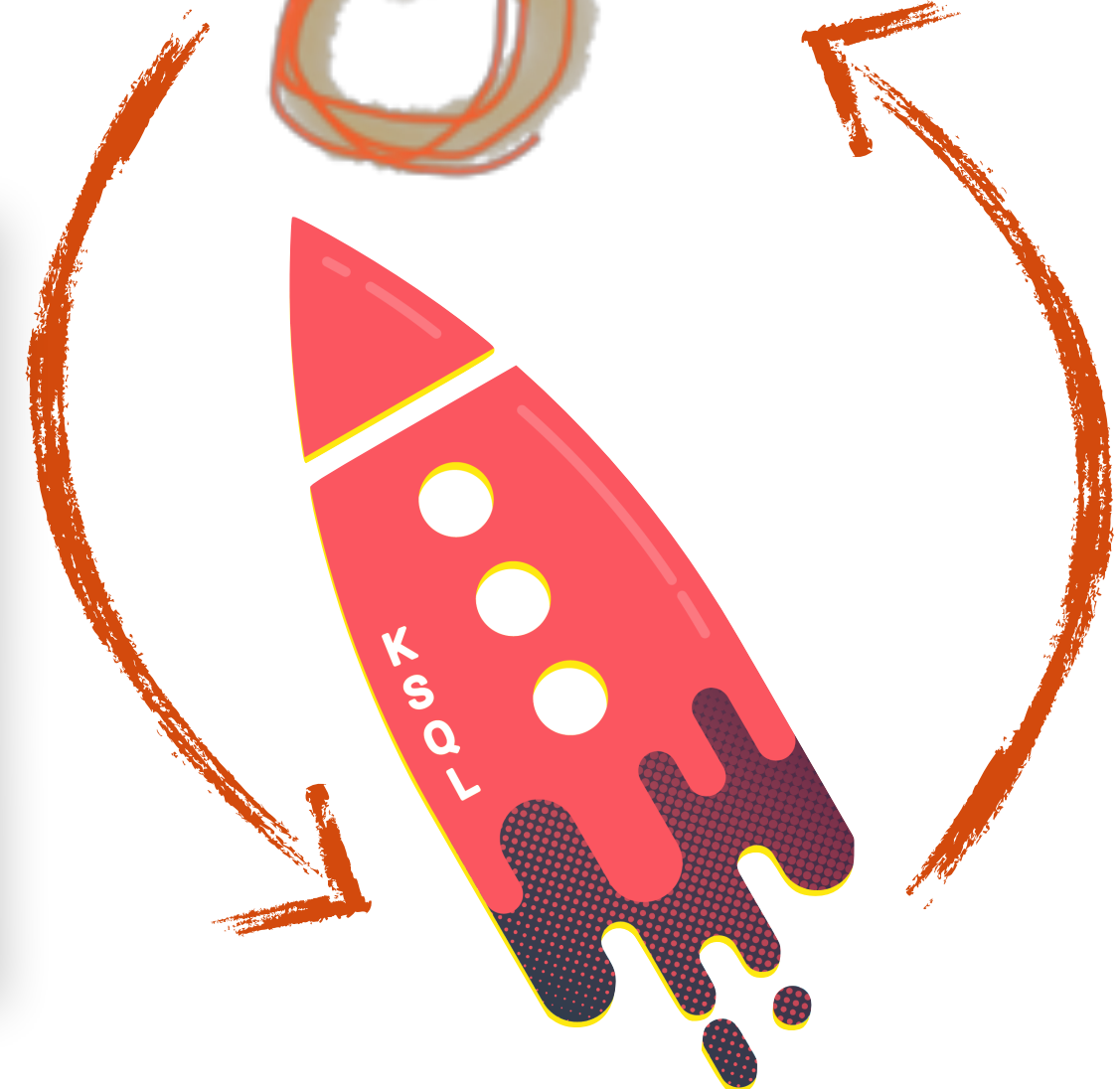
Postgres



```
{  
  "id": 3,  
  "first_name": "Merilyn",  
  "last_name": "Doughartie",  
  "email": "mdoughartie1@dedecms.com",  
  "gender": "Female",  
  "club_status": "platinum",  
  "comments": "none"  
}
```



Postgres

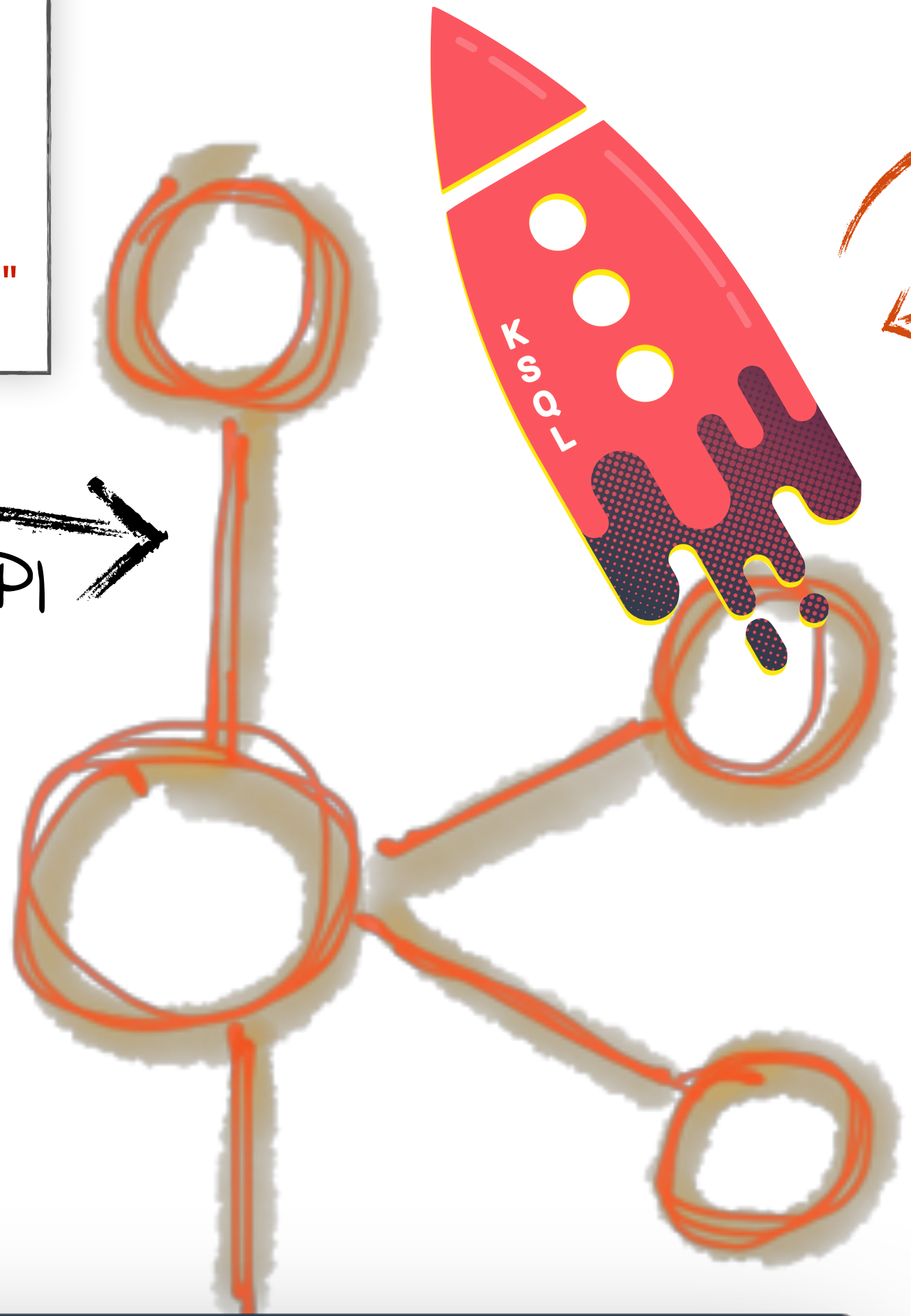


```
{
  "rating_id": 5313,
  "user_id": 3,
  "stars": 4,
  "route_id": 6975,
  "rating_time": 1519304105213,
  "channel": "web",
  "message": "worst. flight. ever. #neveragain"
}
```

Filter all ratings where STARS < 3

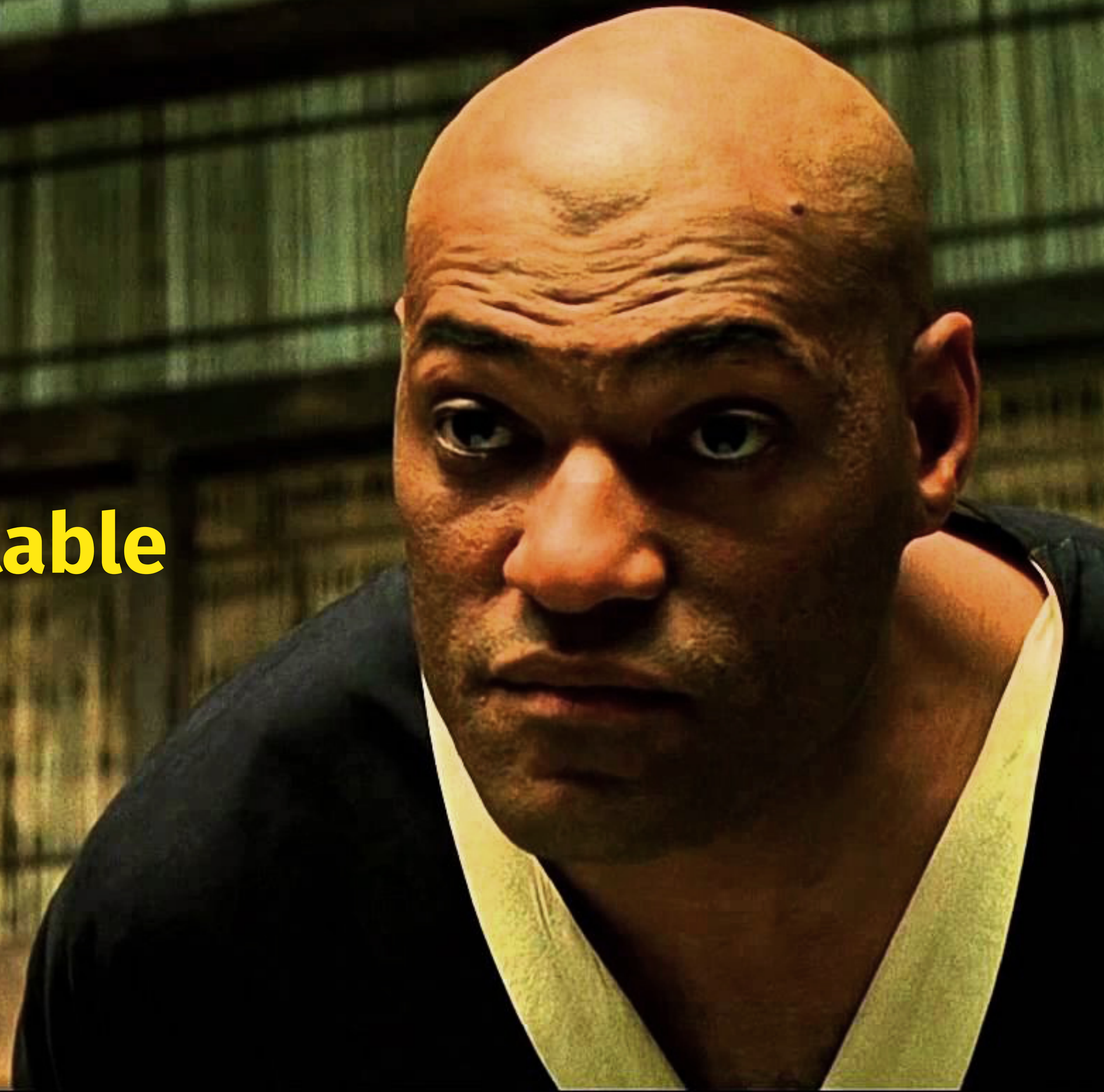
POOR_RATINGS

Producer API



```
CREATE STREAM POOR_RATINGS AS
SELECT * FROM ratings WHERE STARS < 3
```

Do you think that's a **table**
you are querying ?

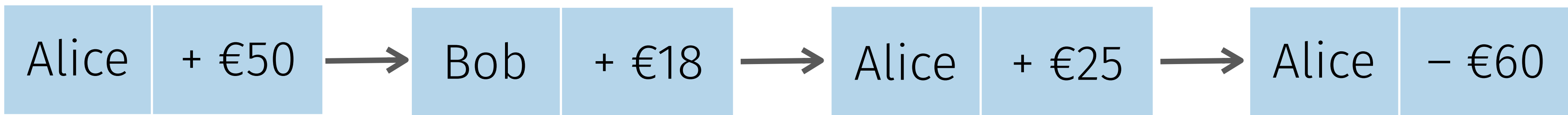


The Stream-Table Duality

Table
(balance)

Alice	€50	Alice	€50	Alice	€75	Alice	€15
		Bob	€18	Bob	€18	Bob	€18

Stream
(payments)



→ time

THE TRUTH IS OUT THERE

The truth is the log.

The database is a cache of a subset of the log.

—Pat Helland

Immutability Changes Everything

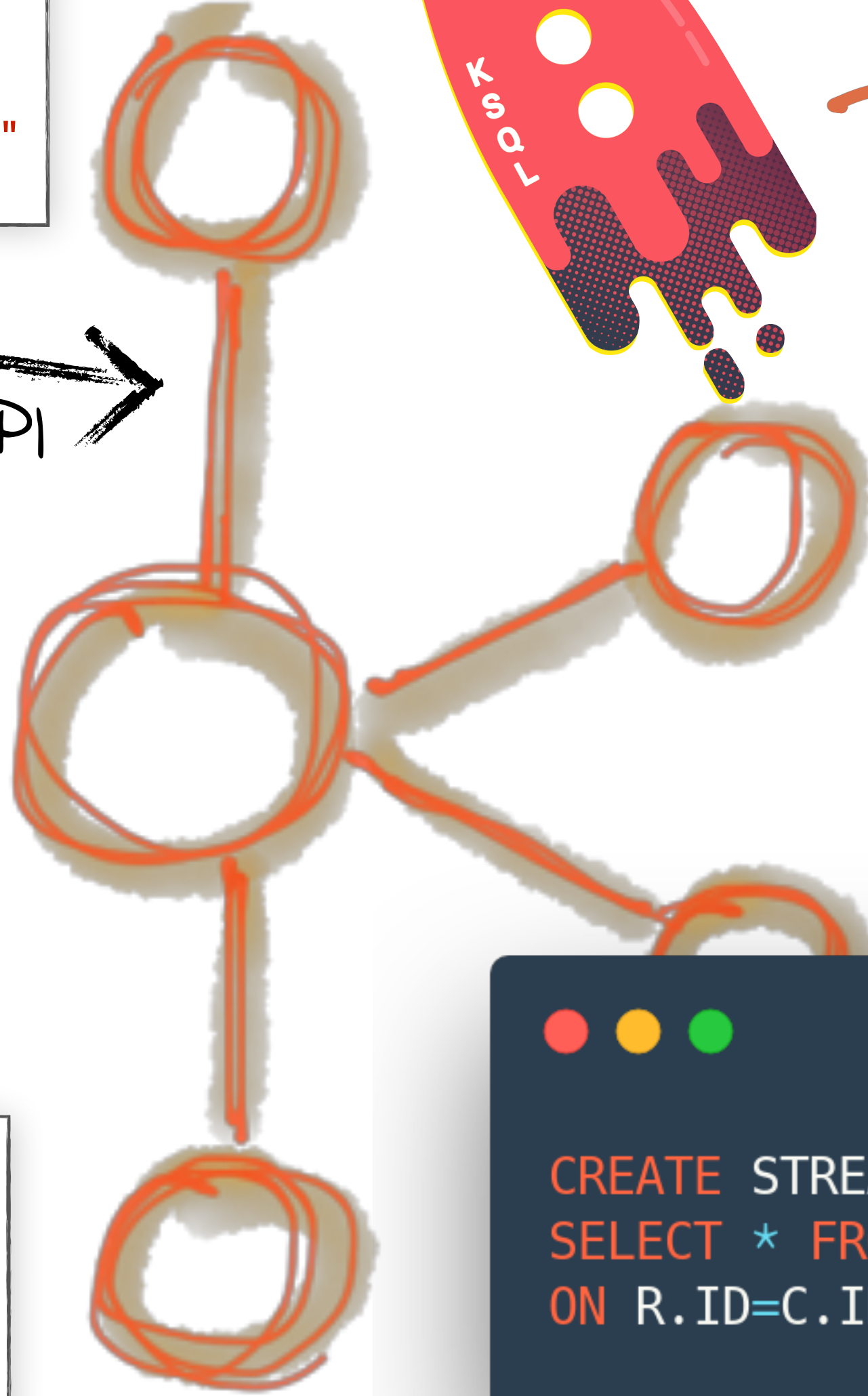
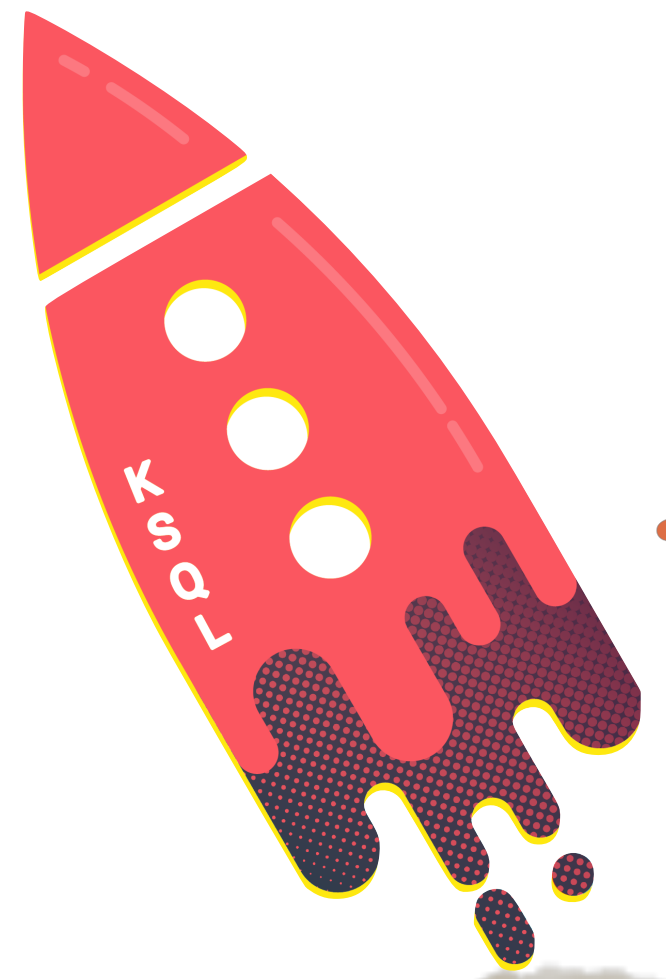
http://cidrdb.org/cidr2015/Papers/CIDR15_Paper16.pdf

```
{
  "rating_id": 5313,
  "user_id": 3,
  "stars": 4,
  "route_id": 6975,
  "rating_time": 1519304105213,
  "channel": "web",
  "message": "worst. flight. ever. #neveragain"
}
```

Producer API

Kafka Connect

```
{
  "id": 3,
  "first_name": "Marilyn",
  "last_name": "Doughartie",
  "email": "mdoughartie1@dedecms.com",
  "gender": "Female",
  "club_status": "platinum",
  "comments": "none"
}
```



Join each rating to customer data

RATINGS_WITH_CUSTOMER_DATA

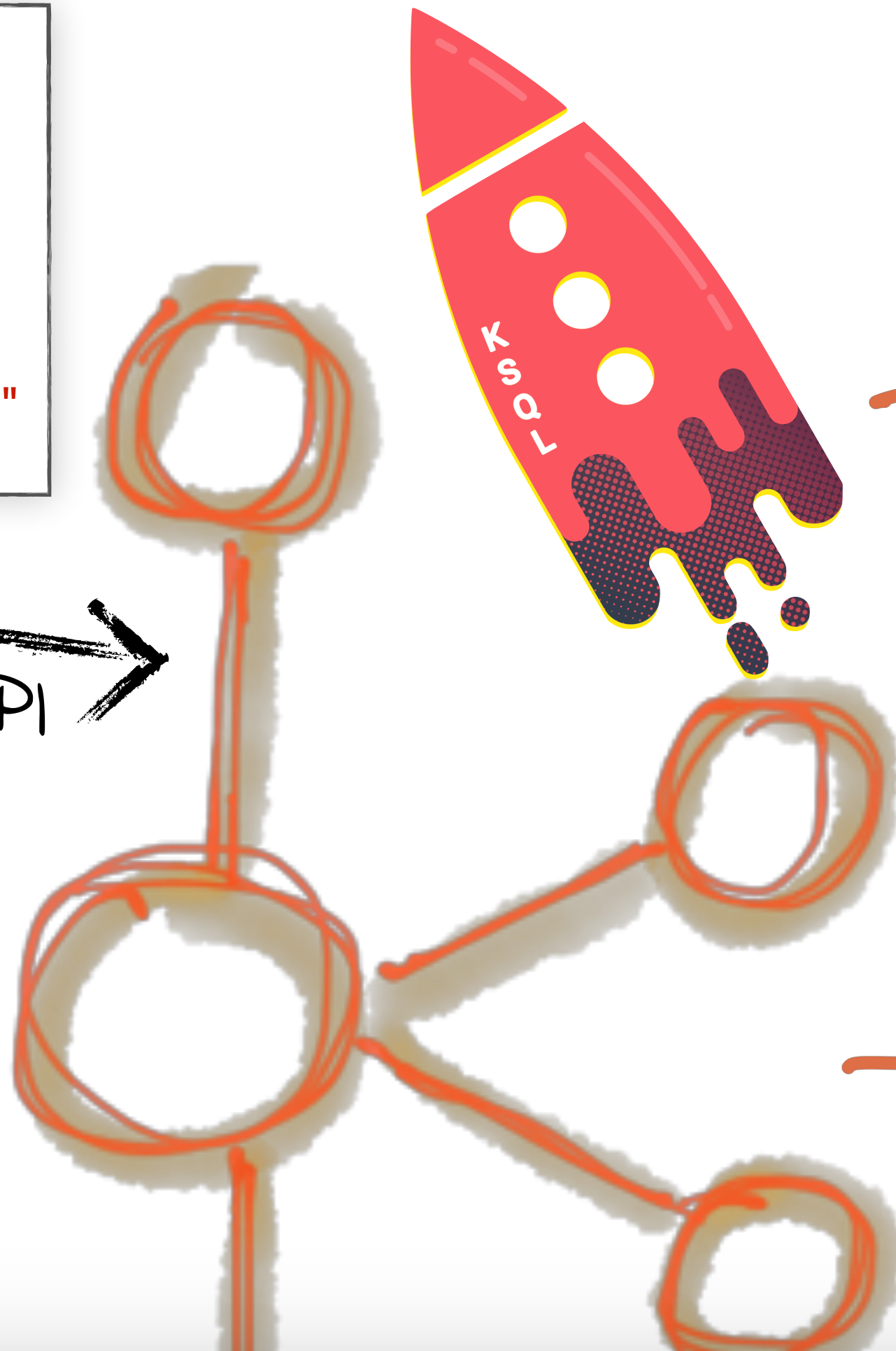
```
CREATE STREAM RATINGS_WITH_CUSTOMER_DATA AS
SELECT * FROM RATINGS LEFT JOIN CUSTOMERS
ON R.ID=C.ID;
```



```
{
  "rating_id": 5313,
  "user_id": 3,
  "stars": 4,
  "route_id": 6975,
  "rating_time": 1519304105213,
  "channel": "web",
  "message": "worst. flight. ever. #neveragain"
}
```

Producer API

Kafka Connect



Join each rating to customer data

```
RATINGS_WITH_CUSTOMER_DATA
```

Filter for just PLATINUM customers

```
UNHAPPY_PLATINUM_CUSTOMERS
```

```
{
  "id": 3,
  "first_name": "Marilyn",
  "last_name": "Doughartie",
  "email": "mdoughartie1@ded",
  "gender": "Female",
  "club_status": "platinum",
  "comments": "none"
}
```

```
CREATE STREAM UNHAPPY_PLATINUM_CUSTOMERS AS
SELECT * FROM RATINGS_WITH_CUSTOMER_DATA
WHERE STARS < 3
```

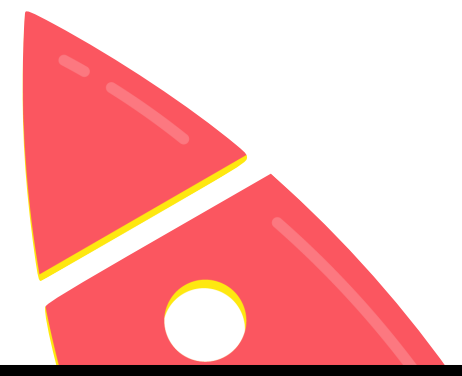


@gamussa

#Postgres

@confluentinc

```
{
  "rating_id": 5313,
  "user_id": 3,
  "stars": 4,
  "route_id": 6975,
  "rating_time": 15193041052,
  "channel": "web",
  "message": "worst. flight"
}
```



```
CREATE TABLE RATINGS_BY_CLUB_STATUS AS
SELECT CLUB_STATUS, COUNT(*)
FROM RATINGS_WITH_CUSTOMER_DATA
WINDOW TUMBLING (SIZE 1 MINUTES)
GROUP BY CLUB_STATUS;
```

7



Kafka Connect

```
{
  "id": 3,
  "first_name": "Marilyn",
  "last_name": "Doughartie",
  "email": "mdoughartie1@dedecms.com",
  "gender": "Female",
  "club_status": "platinum",
  "comments": "none"
}
```



Aggregate per-minute by CLUB_STATUS

RATINGS_BY_CLUB_STATUS_1MIN

@gamussa #Postgres @confluentinc



Resources and Next Steps



<https://github.com/confluentinc/examples>



<http://confluent.io>



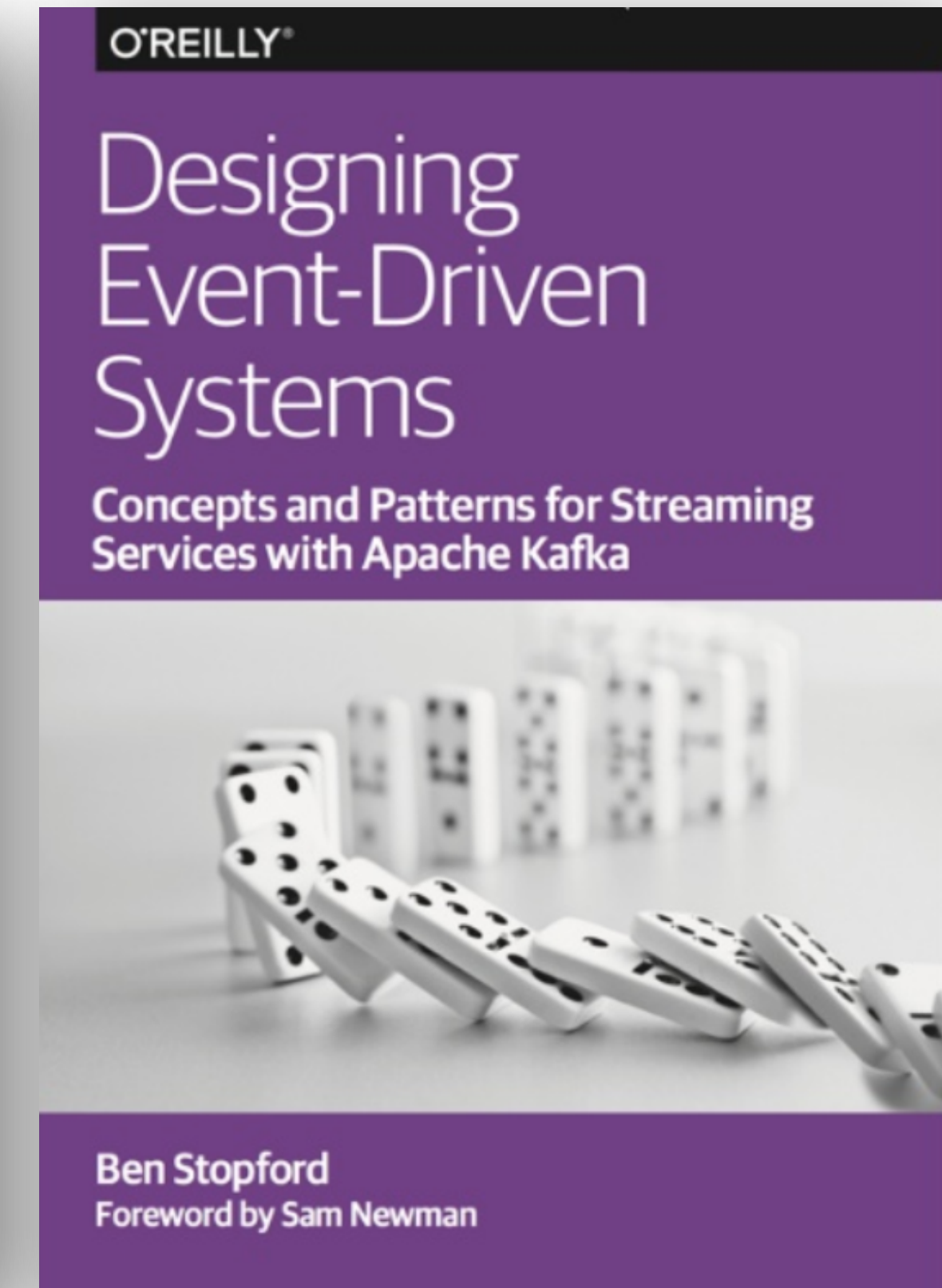
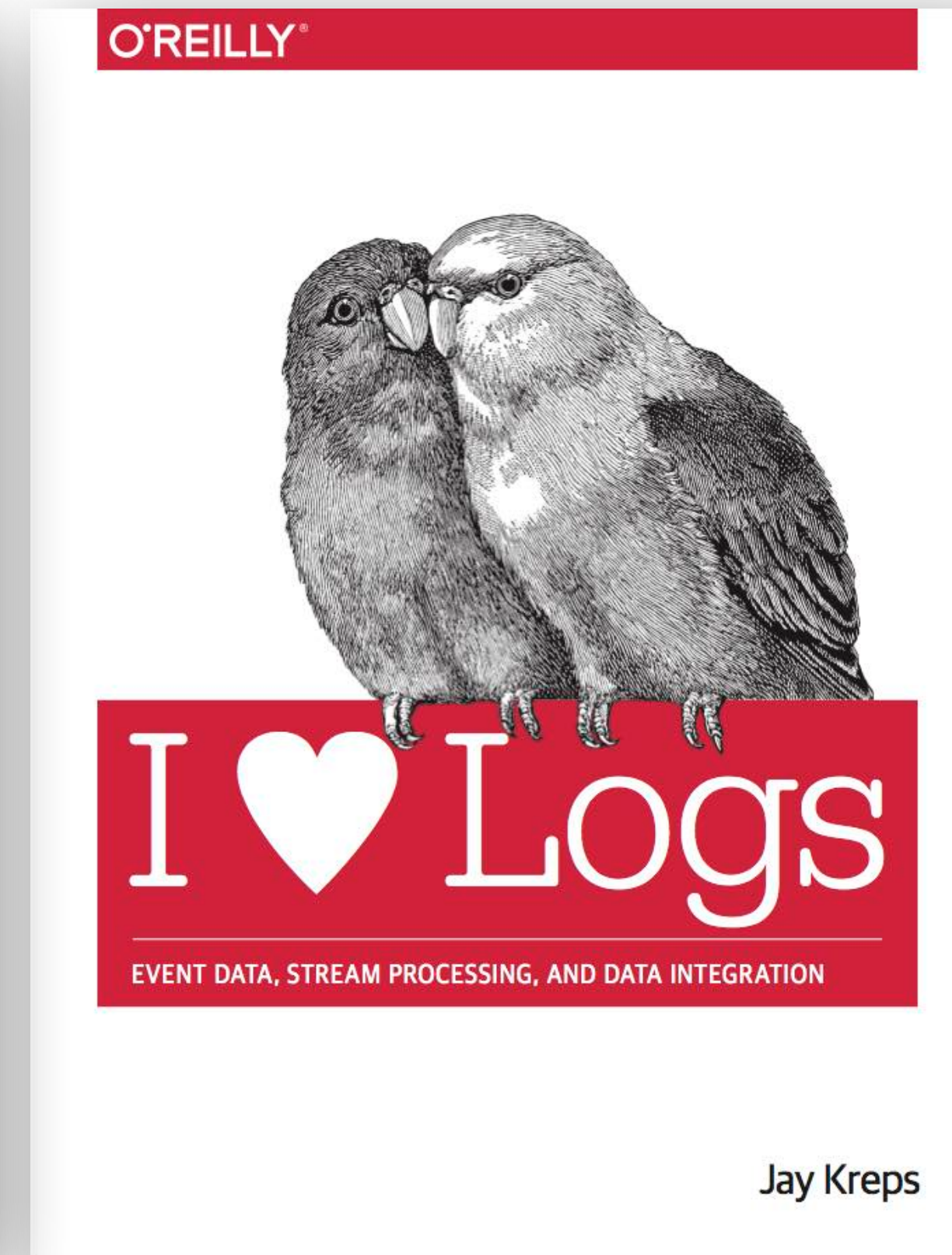
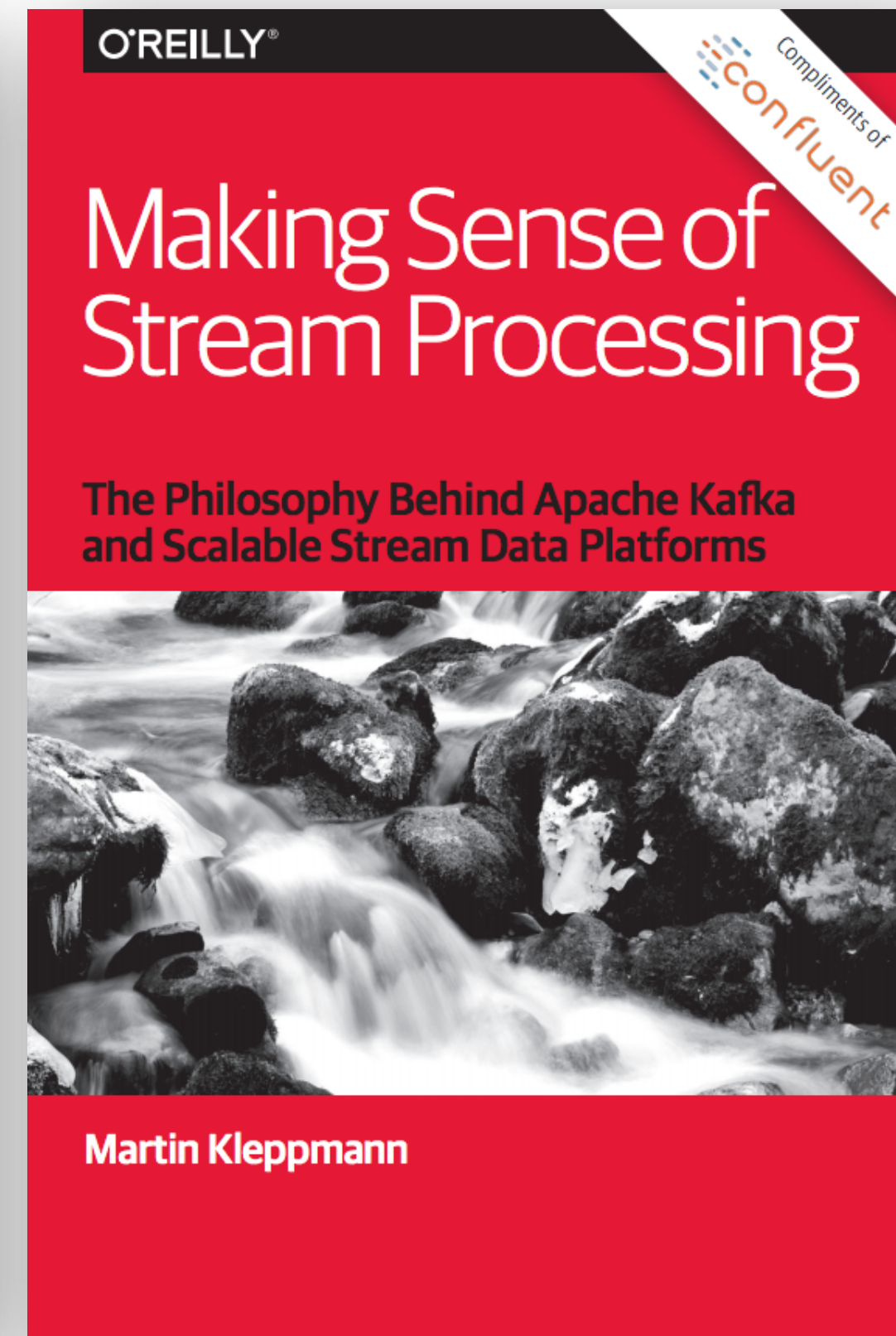
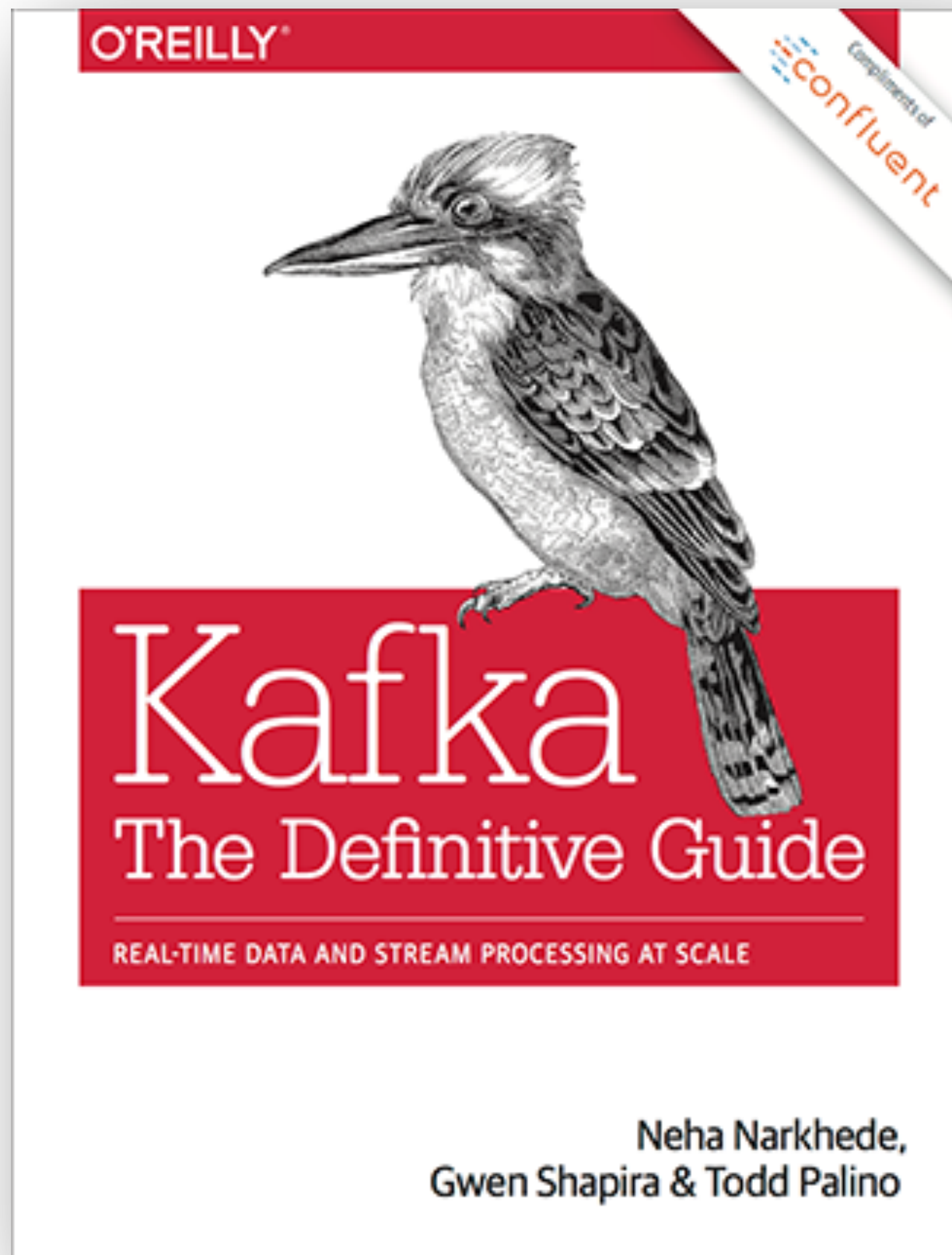
<https://slackpass.io/confluentcommunity>

#ksql #connect

Free Books!



<https://www.confluent.io/apache-kafka-stream-processing-book-bundle>



@gamussa

#Postgres

@confluentinc

ONE LAST THING...

kafka summit

<https://kafka-summit.org>

Gamov30

APRIL 2, 2019

NEW YORK CITY

REGISTER

LEARN MORE

THANKS!



@gamussa

viktor@confluent.io