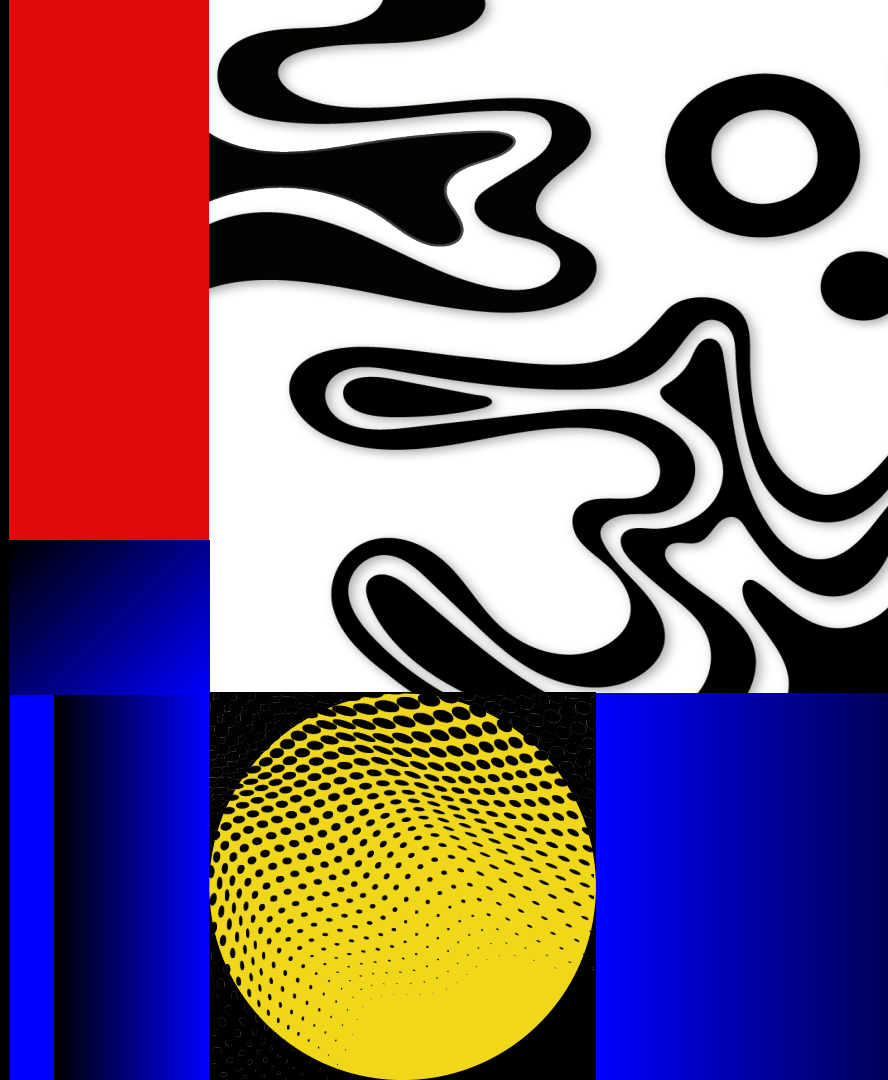


BIG DATA 2022

BECOMING A DATA ROCKSTAR

Katherine García
Andrea Reyes





AGENDA

01

**Descripción
del proyecto**

02

Milestone 1-3

03

**Retos y lecciones
aprendidas**

BECOMING A DATA ROCKSTAR

HOTELES Aa

Una cadena de hoteles busca una solución que permita gestionar toda la información de sus clientes y reservas

- ❑ Operaciones: sponsors del proyecto.
- ❑ Customer experience: experiencia de usuarios antes, durante y después de la estadía en el hotel.
- ❑ Business analytics: preparar información para reportes gerenciales.

OBJETOS DE DATOS



Hotel

Activos, ubicación, etc.



Habitación

Servicios disponibles



Tarifas

Precios base para cada habitación



Ocupación

Disponibilidad de habitaciones.



Cliente

Usuario único con información personal

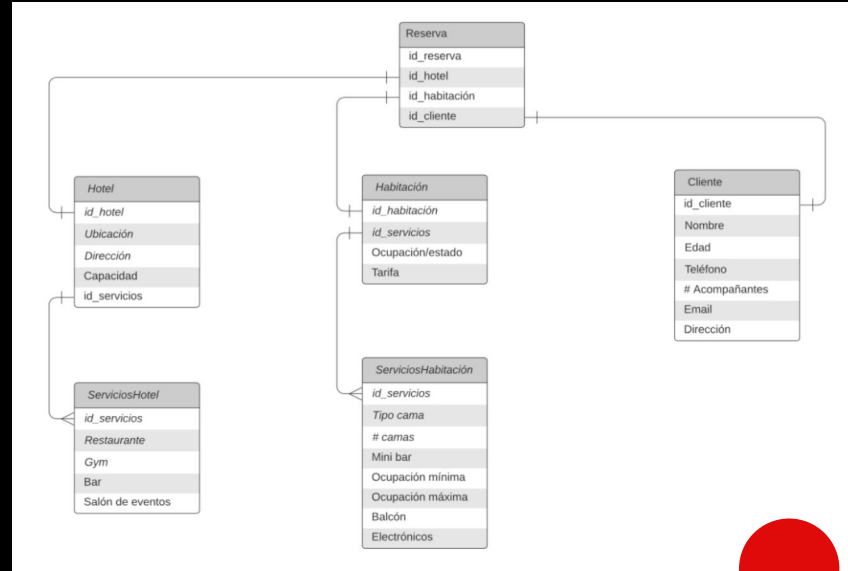
MILESTONE 1

- Diseñar un modelo entidad relación que se adapte al requerimiento
- Instalar una base de datos
- Crear el modelo
- Cargar datos de prueba



MODELO ENTIDAD RELACIÓN

ERD



BASE DE DATOS

- Docker
- postgrSQL

```
You, 21 hours ago | 1 author (You)
1  version: '3.7'           You, 21 hours
2  services:
3    db:
4      image: postgres
5      ports:
6        - 5432:5432
7      mem_limit: 1G
8      mem_reservation: 128M
9      cpus: 1
10     environment:
11       - POSTGRES_USER=rock
12       - POSTGRES_PASSWORD=rock
13       - POSTGRES_DB=rockstar
```

TABLAS

```
1 CREATE TABLE hotel(  
2     id_hotel int primary key,  
3     ubicacion varchar,  
4     direccion varchar,  
5     capacidad varchar,  
6     id_servicios int,  
7     CONSTRAINT id_servicios  
8         FOREIGN KEY(id_servicios)  
9         REFERENCES "serviciosHotel"(id_servicios)  
10        ON DELETE CASCADE  
11 );  
12  
13 CREATE TABLE habitacion(  
14     id_habitacion int primary key,  
15     ocupacion_estado varchar,  
16     tarifa int,  
17     id_servicios int,  
18     CONSTRAINT id_servicios  
19         FOREIGN KEY(id_servicios)  
20         REFERENCES "serviciosHabitacion"(id_servicios)  
21        ON DELETE CASCADE  
22 );
```

```
24 CREATE TABLE cliente(  
25     id_cliente int primary key,  
26     nombre varchar,  
27     edad int,  
28     telefono int,  
29     email varchar,  
30     direccion varchar,  
31     cant_acompañantes int,  
32     id_servicios int  
33 );
```

```
38 CREATE TABLE reserva(  
39     id_reserva int primary key,  
40     id_hotel int,  
41     id_habitacion int,  
42     id_cliente int,  
43     CONSTRAINT id_hotel  
44         FOREIGN KEY(id_hotel)  
45         REFERENCES "hotel"(id_hotel)  
46        ON DELETE CASCADE,  
47  
48     CONSTRAINT id_habitacion  
49         FOREIGN KEY(id_habitacion)  
50         REFERENCES "habitacion"(id_habitacion)  
51        ON DELETE CASCADE,  
52  
53     CONSTRAINT id_cliente  
54         FOREIGN KEY(id_cliente)  
55         REFERENCES "cliente"(id_cliente)  
56        ON DELETE CASCADE  
57 );
```


DATOS DE PRUEBA

WHERE		ORDER BY			
	id_hotel	ubicacion	direccion	capacidad	id_servicios
1	10001	Guatemala	20 calle 8-43 zona 13	100	11001
2	10002	El Salvador	4ta calle 7-76 zona 8	115	11002
3	10003	Costa Rica	Calle Manuel F. Ayau, zona 10	85	11003
4	10004	Honduras	San Cristobal 18 av 1-83	50	11004
5	10005	Belice	Vista Hermosa III, zona 16	98	11005

WHERE		ORDER BY		
	id_habitacion	ocupacion_estado	tarifa	id_servicios
1	10001101	Reserva	120	1000110101
2	10001102	Reserva	150	1000110102
3	10001103	Ocupada	200	1000110103
4	10001104	Cancelada	150	1000110104
5	10001105	Reserva	150	1000110105

WHERE		ORDER BY					
	id_cliente	nombre	edad	telefono	email	direccion	cant_acompañantes
1	12341	Katherine Garcia	20	42186759	Kg@email.com	Guatemala	2
2	12342	Andrea Reyes	21	52030518	Andreareyes@email.com	USA	1
3	12343	Maria Paz	35	56325878	Mpaz@email.com	Guatemala	3
4	12344	Esteban Lopez	66	41076352	lopezEst@email.com	El Salvador	4
5	12345	Monica Sah	25	35447800	monSah@email.com	Panama	1

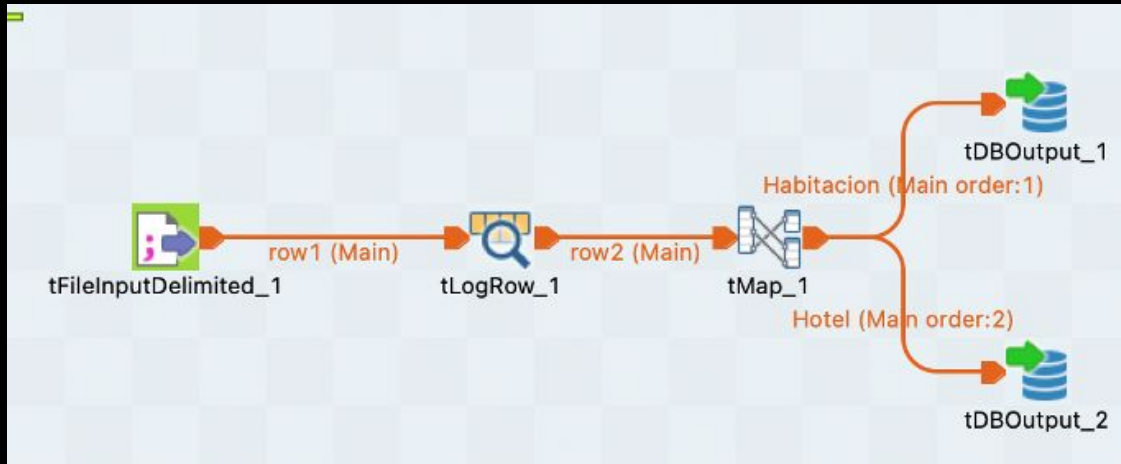
MILESTONE 2

1. Construir un ETL que permita cargar la información de los hoteles desde archivos planos. (CSV) hacia la base de datos relacional.
2. Construir un repositorio en HDFS para los archivos de data analytics.
3. Construir un ETL que tome una captura diaria del estado de reservaciones y clientes y lo cargue en HDFS.
4. Construir una capa en HIVE que permita consultar la información que se ha recibido desde la base relacional.

ETL

Con el archivo proporcionado, se crearon por medio del ETL las dos tablas principales de nuestro modelo.

Las tablas creadas fueron: Habitación y Hotel.



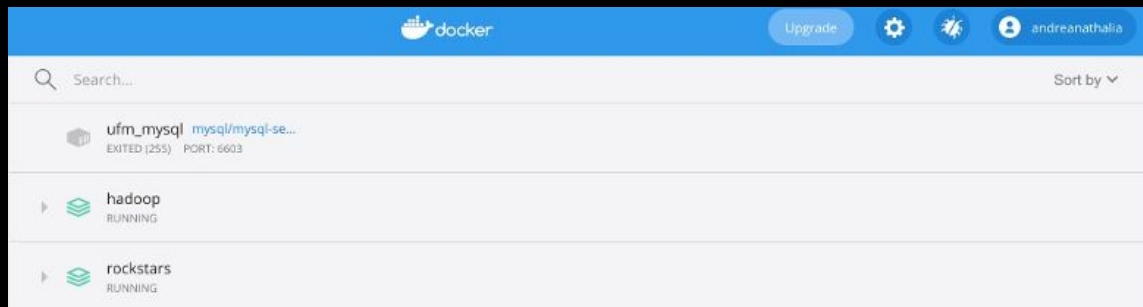
TALLENDO

HDFS

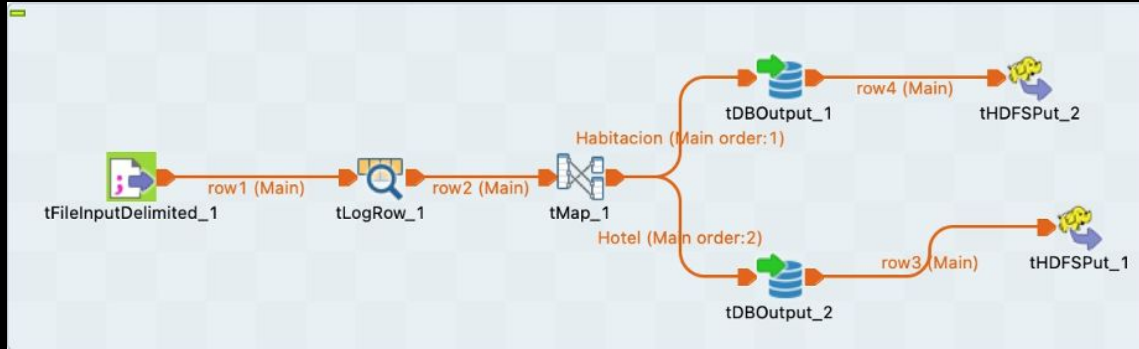
01

Docker

Un Docker Compose que nos proporcionaba las herramientas necesarias para poder construir el repositorio.



CAPTURA DIARIA



CONSULTA HIVE

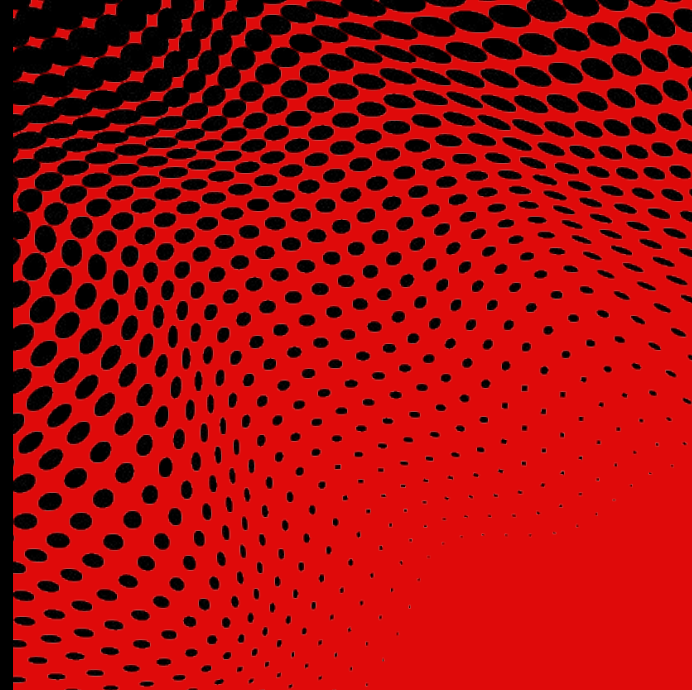
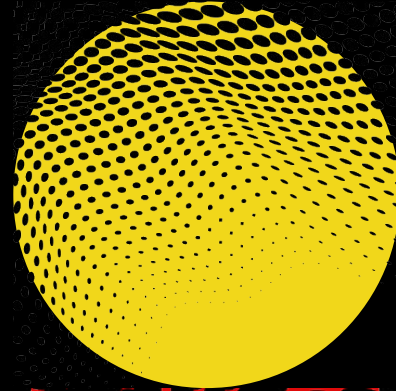
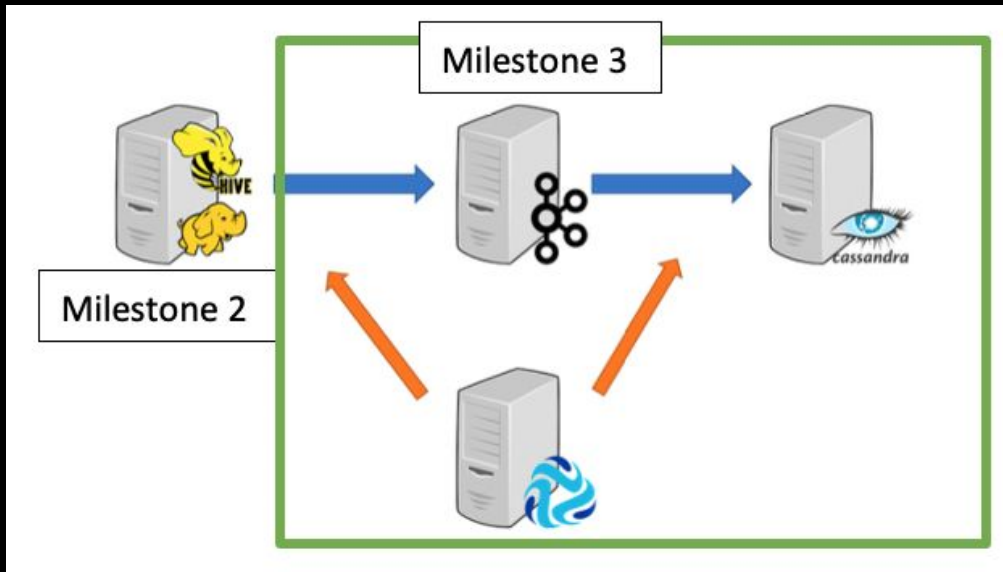
Tablas que reciben los datos desde hdfs



MILESTONE 3

- Construir un modelo en Cassandra que permita replicar el estado actual de las reservaciones para lograr escalabilidad y mejorar tiempos de respuesta.
- Construir un tópico de kafka que gestionará todas las actualizaciones que se envían a Cassandra
- Crear un data pipeline en Streamsets que cargue la información de las reservaciones de la base relacional a Kafka, y un Pipeline que tome la información de Kafka y la envíe a Cassandra.

M3



DOCKER

```
docker-compose.yml • Untitled-1
Users > katherinegarcia > Desktop > Rockstars > docker-compose.yml
1
2  version: '3'
3  services:
4
5      cassandra:
6          image: cassandra:latest
7          ports:
8              - 9042:9042
9          volumes:
10             - ~/apps/casadra:/var/lib/cassandra
11          environment:
12             - CASSANDRA_CLUSTER_NAME = dataRockstars
```


MODELO

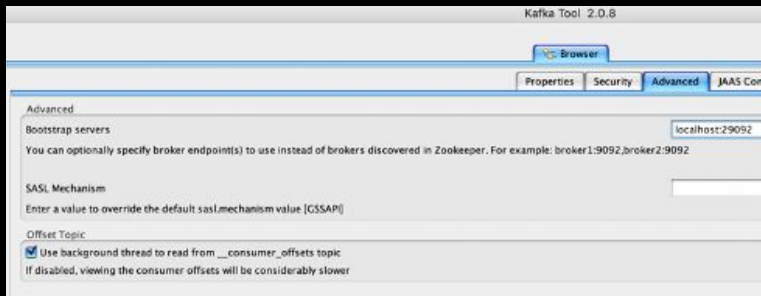
```
Users > katherinegarcia > Desktop > Rockstars > cassandra > ≡ cass.cql
1  -- Construir un modelo en Cassandra que permita replicar
2  -- estado actual de las reservaciones
3  -- para lograr escalabilidad y mejorar tiempos de respuesta.
4
5  -- MODELO
6  CREATE KEYSPACE rocky
7      WITH REPLICATION = {
8          'class' : 'SimpleStrategy'
9          'replication_factor' : 1
10 };
11
12 CREATE TABLE reserva(
13     id_reserva int primary key,
14     id_hotel int,
15     id_habitacion int,
16     id_cliente int,
17     fecha_ingreso date,
18     fecha_egreso date,
19 );
20
21 DROP TABLE rocky.reserva
22
```

KAFKA

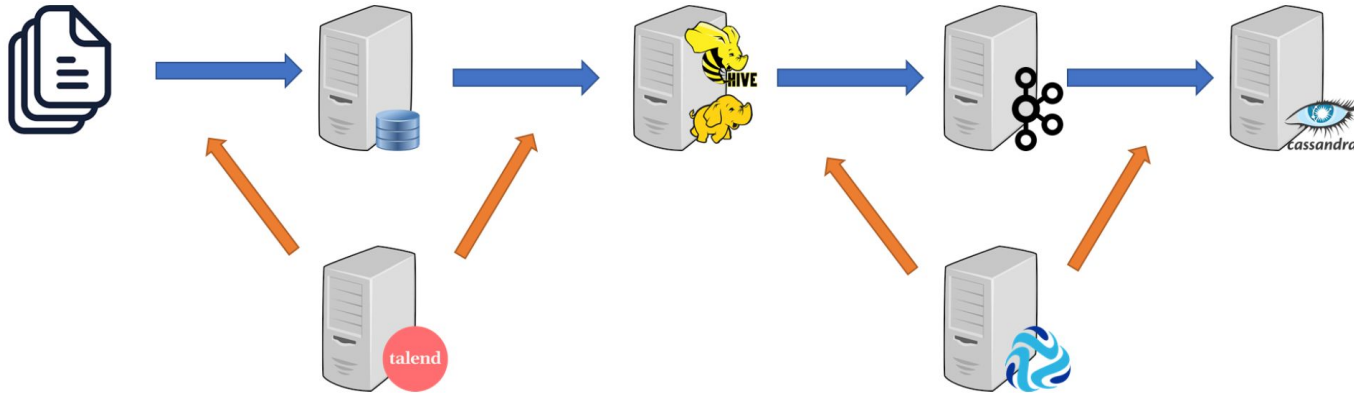
> katherinegarcia > Desktop > Rockstars > 🐳 docker-cor

```
version: '3'
services:
  zookeeper:
    image: confluentinc/cp-zookeeper:latest
    environment:
      ZOOKEEPER_CLIENT_PORT: 2181
      ZOOKEEPER_TICK_TIME: 2000
    ports:
      - 22181:2181

  kafka:
    image: confluentinc/cp-kafka:latest
    depends_on:
      - zookeeper
    ports:
      - 29092:29092
    environment:
      KAFKA_BROKER_ID: 1
      KAFKA_ZOOKEEPER_CONNECT: zookeeper:2181
      KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
```



MODELO COMPLETO



RETOS

CLOUD

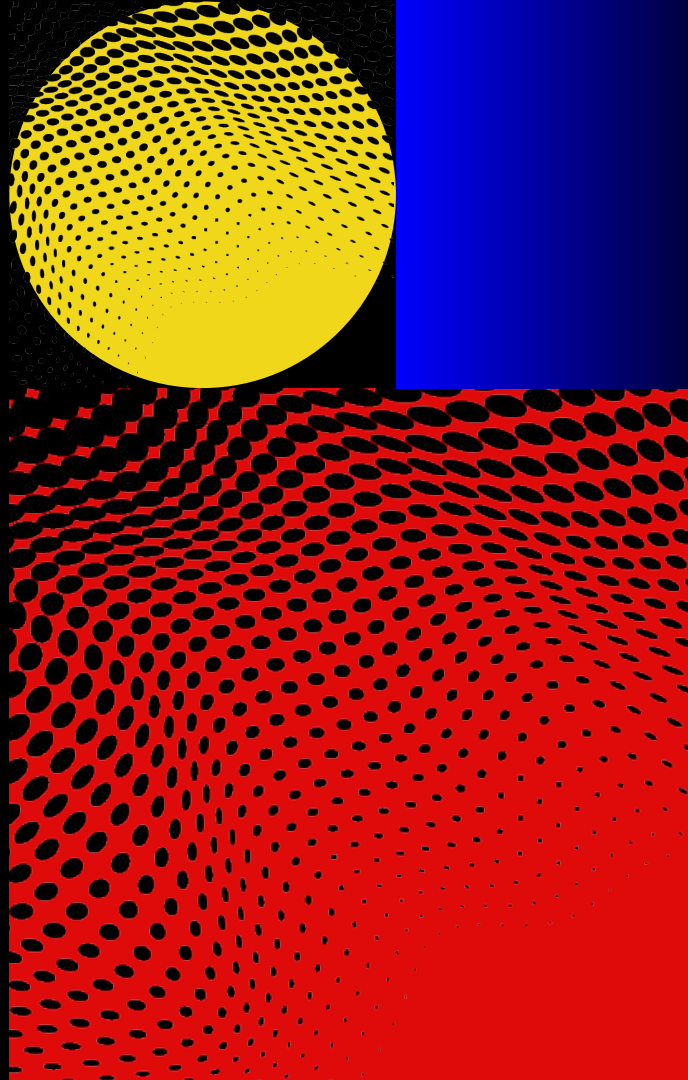
El reto principal fue el tiempo que nos tomó pensar en la sincronización talend-hdfs cloud

Herramientas

En el pasado no habíamos utilizado ninguna de las herramientas necesarias para este proyecto. La curva de aprendizaje y el tiempo fue un reto.

DISTRIBUCIÓN

Al trabajar con recursos locales, la distribución de trabajo fue un reto para nosotros como equipo.



LECCIONES APRENDIDAS



EFFECTIVIDAD DE DOCKER

Poder sincronizar todas las fases localmente por medio de docker.



USO DE NUEVAS HERRAMIENTAS

HIVE
HDFS
CASSANDRA
TALEND
DOCKER



GRACIAS

Do you have any questions?

Andrea Reyes
Katherine García