

DATU MASIBOEN PROZESAMENDURAKO AZPIEGITURAK

3.LANA

Andoni Sudupe eta Oier Ijurko

Aurkibidea

1.	Sarrera.....	3
2.	Datuak.....	4
3.	Erabili ditugun tresnak.....	4
3.1.	Datuen jatorria:MQTT.....	4
3.2.	Datuen biltegitratzea: Kafka.....	5
3.3.	Datuen prozesamendua: Flink.....	5
3.4.	Datuak esportatu: ElasticSearch.....	5
3.5.	Datuak bistaratu: Kibana.....	5
4.	Lana garatzeko prozesua.....	6
4.1.	Compose eta dockerfile.....	6
4.2.	Python script-a.....	7
4.3.	Kafka-connect eta datuak lortu.....	8
4.4.	Flink Shell-eko deiak.....	9
4.5.	Kibanako Dashboarda.....	10
5.	Hobekuntza.....	10
5.1.	Grafana erabiltzen saiatu.....	10
6.	Ondorioak.....	11
7.	Erreferentziak.....	12

1. Sarrera

Lan honen helburu nagusia datu fluxuak prozesatzen dituen soluzio global bat lortzea da. Horretarako, ondorengo osagaiak elkarlanean jarriko ditugu:

1. MQTT, bai broker-ak eta baita bezeroak (igorlea eta harpideduna) ere.
2. Topic-etan antolatutako datu fluxuak biltegitatzeko Kafka erabili.
3. Datu fluxuak prozesatzeko Flink erabiltzea
4. Superset Impala-ra konektatu, bi taulen gainean kontsultak egin ahal izateko.
5. Elasticsearch (+Kibana) erabili datuak bistaratzeko grafikoki

2. Datuak

Erabili ditugun datuak, smartphone eta smartwatch aktibitatea eta biometrikak jasotzen dituen dataset batetik hartu ditugu. 51 pertsonen 3 minutuko 18 aktibitate egiteko eskatu zitzaion, eta mugikorra eta smartwatch-arekin giroskopia eta azelerometroa erabiliz, hauen mugimenduak kontrolatu ziren. Datuak 20 Hz-ko maiztasunean gordeta daude eta dena 15,630,426 datu gordin.

Ezkerreko irudian aktibitate desberdinak agertzen dira eta eskuineko argazkian datu ilara bakoitzak izango duen formatua.

Activity	Code
Walking	A
Jogging	B
Stairs	C
Sitting	D
Standing	E
Typing	F
Brushing Teeth	G
Eating Soup	H
Eating Chips	I
Eating Pasta	J
Drinking from Cup	K
Eating Sandwich	L
Kicking (Soccer Ball)	M
Playing Catch w/Tennis Ball	O
Dribbling (Basketball)	P
Writing	Q
Clapping	R
Folding Clothes	S

Field name	Description
Subject-id	Type: Symbolic numeric identifier. Uniquely identifies the subject. Range: 1600-1650.
Activity code	Type: Symbolic single letter. Identifies a specific activity as listed in Table 2. Range: A-S (no "N" value)
Timestamp	Type: Integer. Linux time
x	Type Numeric: real. Sensor value for x axis. May be positive or negative.
y	Same as x but for y axis
z	Same as x but for z axis

3. Erabili ditugun tresnak

Lan hau garatzeko hainbat tresna erabili ditugu. Hona hemen sarreraren aipaturako helburua gauzatzeko beharrezkoak izango ditugun baliabideak:

3.1. Datuen jatorria: MQTT

MQTT, subscriber-publisher motako protokoloa da. Normalean, baliabide mugatuak (memoria edo banda-zabalera) dituzten gailuen artean mezuak trukatzeko erabiltzen da protokolo hau. Bere abantailen artean, erabiltzeko oso erraza duen API bat duela dago.



3.2. Datuen biltegitratzea: Kafka



Apache Kafka gertaera sakabanatuko biltegia eta streaming bidez prozesatzeko open source plataforma bat da. Honen ideia, denbora errealeko datuekin lan egiteko sendoa, errendimendu handikoa eta latentzi baxukoa den plataforma bat eskaintzea da.

3.3. Datuen prozesamendua: Flink

Apache Flink fluxu-prozesamendurako erabiltzen den framework-a da. Eskala handiko datu-fluxuak prozesatzeko eta streaming aplikazioarekin prozesatutako datuei buruzko informazio analitikoa denbora errealean emateko erabiltzen da.

Lan honetan, Flink-eko SQL API-a erabiliko dugu taulak sortu eta Kafkan gordetako datuak Elasticsearch-era esportatzeko.



3.4. Datuak esportatu: Elasticsearch

Elasticsearch kode irekiko bilaketa- eta analisi-motor bat da, eta datuen azterketa errazten duena. Honen bitartez, kafkan gordeta ditugun datuak Kibana bezalako datu-bistaratzeko ahalbidetuko duen aplikazioetara bidaltzeko ahalmena lortuko dugu.

3.5. Datuak bistaratu: Kibana

Elastic Stack-en gainean dagoen frontend-en aplikazioa da Kibana, eta datuak bistaratzeko eta Elasticsearch-en indexatutako datuak bilatzeko gaitasuna ematen du.

Honen bitartez, hasieran aipatutako datuak grafikoki bistaratzeko gai izango gara, dashboard bat sortuz lan honen amaieran.



Kibana

4. Lana garatzeko prozesua

4.1. Compose eta Dockerfile

Egindako beste lanetan bezala, .yml fitxategietan adierazten dira lanean erabiliko diren kontainerrak. Kasu honetan, ondorengo kontainerrak izango ditugu: Mosquitto, Zookeeper, Kafka, Kafka-connect, ElasticSearch, Kibana, sql-client, jobmanager, taskmanager eta python kodea exekutatu duen kontainer bat.

Mosquitto eta pythoneko kontainerra erabiliz, mqtt-client bat sortu eta broker eta portu bat erabiliz, deskargatutako datuak lortu ahal izango ditugu. Ondorengo atalean azalduko dugu hau zehaztasun handiagorekin.

Kafka eta Kafka-connect kontainerrekin, datuak gorde eta source eta sink-ak adierazteko gai izango gara. Modu honetan, MQTT-tik lortutako datuak gorde eta Flink-era eramateko gai izango gara. Horretarako, Kafka-connect-ek environment atalean dituen “connect” atalak erabiliko ditugu, baita deskargatutako jar-ak erabili ere.

Mosquitto eta Kafka arteko konexioa funtzionatu ahal izateko, json formatuan dagoen ondorengo connector-a erabili behar izan dugu:

```
{
  "name": "mqtt-source",
  "config": {
    "connector.class": "io.confluent.connect.mqtt.MqttSourceConnector",
    "tasks.max": 1,
    "mqtt.server.uri": "tcp://mosquitto:1883",
    "mqtt.topics": "smart",
    "kafka.topic": "smart",
    "value.converter": "org.apache.kafka.connect.converters.ByteArrayConverter",
    "confluent.topic.bootstrap.servers": "kafka:9092",
    "confluent.topic.replication.factor": 1
  }
}
```

Bertan, mosquitto-ren portua, Kafka-ren portua, topic-a eta bestelako informazioa adierazi dugu funtzionatu ahal izateko.

Zookeeper kontainerren helburua zerbitzuen konfigurazioa eta sinkronizazioa ematea izango da. Normalean, Kafka, Zookeeper-rekin batera erabiltzen da.

Flink erabili ahal izateko, lehen aipatutako sql-client, jobmanager eta taskmanager kontainerrak erabiliko ditugu. Hauen bitartez, SQL API-a erabiliko dugu Kafkatik datuak hartu eta ondoren azalduko dugun elasticsearch-era esportatzeko. Datuak esportatzerakoan, Job berri bat sortuko da, dagokien Job ID-arekin.

Elasticsearch kontainerren zeregina, Flink-en izango ditugun datuak esportatzea izango da, portu bat esleituz eta ondoren, Kibana bezalako aplikazio batean bistaratzu (honek ere bere kontainerra izango du, bere portuarekin).

Dockerfile-ari dagokionez, python script-ean client bat sortu ahal izateko erabiliko dugun paho.mqtt paketea instalatuko dugu, baita ondoren azalduko dugun .py-a exekutatu ere.

4.2. Python script-a

Datuak irakurri eta publikatzeko, python fitxategi bat sortuko dugu enuntziatuan emandako kodea abiapuntu hartuta. Hasieran, WISDM dataseta dagoen direktorio desberdinak definitzen dira, eta parametro batzuk irakurketa ongi egiteko.

```
import os, time
import random
from paho.mqtt import client as mqtt_client

rootpath = "/data/"
paths = ["raw/phone/accel/", "raw/phone/gyro/", "raw/watch/accel/", "raw/watch/gyro/"]
WAIT_TIME = 0.5
MAX_LINES = 1000000000 #connect mosquitto
i = 0
```

Ondoren, mqtt eta Kafka konektatzeko konektoreak martxan jartzeko komandoa eta sleep bat ditugu. Hau eskuz egin daiteke terminal batetik, baina horrela automatikoki egingo du. MQTT-kin publikatzeko beharrezko hasieraketak ere jarri ditugu, broker bat eta portu bat definitu ditugu bezeroa konektatzeko.

```
time.sleep(20)
#os.system('curl -d @/home/oier/Documents/kafka/config/connect-mqtt-source.json -H "Con
broker = 'mosquitto'
port = 1883
client = mqtt_client.Client()
client.connect(broker,port)
```

Honekin datuak irakurtzen hasi gaitzke. Fitxategi guztiak zeharkatuko ditugu eta fitxategiko ilara bakoitza hartu, json formatuko string bat sortu eta “smart” topikoan publikatuko dugu 0.5 segunduro . Honela, sorgailu programa bat sortu dugu, gure datuak stream prozesamendua egiteko moduan bidaliko ditu datuak broker batera smart gaiari lotuta.

```
for path in paths:
    data_files = sorted(os.listdir(rootpath+path))
    for data_file in data_files:
        with open(rootpath+path+data_file, "r") as f:
            maxl = 0
            while True:
                try:
                    line = f.readline()
                except UnicodeDecodeError:
                    continue
                if not line or len(line) == 1:
                    break
                line = f.readline().split(",")
                try:
                    if line[-1][-1] == ";":
                        line[-1] = line[-1][:-1]
                    elif line[-1][-1] == "\n":
                        line[-1] = line[-1][:-2]
                except:
                    continue
                line = '{"usid": ' + line[0] + ', "action": "' + line[1]
                #publish egin
                time.sleep(WAIT_TIME)
                result = client.publish("smart",line)
                maxl += 1
                if maxl > MAX_LINES:
                    break
            i += 1
```

Ondo publikatzen direla ikusteko, mosquitto bezero bat definitu dezakegu “smart” topikotik mezuak jasotzeko subscriber moduan, eta mezuak jasotzen diren ikusi.

```
{ "usid":1600,"action": "D","ts": 241602851953437,"x":3.9605713,"y":4.569107,"z":
6.9317017}
{"usid":1600,"action": "D","ts": 241602952661445,"x":3.9846191,"y":4.529175,"z":
6.9095764}
{"usid":1600,"action": "D","ts": 241603053369453,"x":3.9538574,"y":4.555298,"z":
6.9383545}
{"usid":1600,"action": "D","ts": 241603154077461,"x":3.9732666,"y":4.572586,"z":
6.887512}
{"usid":1600,"action": "D","ts": 241603254785469,"x":3.9729614,"y":4.568268,"z":
6.9502716}
^Coier@oier-R0G-Zephyrus-G14-GA401QE-GA401QE:~$ docker run -it --rm --name mqtt-
publisher --network kafka_default efrecon/mqtt-client sub -h mosquitto -t "smart
"
```

4.3. Kafka connect eta datuak lortu

MQTT sorgailuak publikatutako datu fluxuak, Kafka kluster batean biltegitratuko dira. Hau egiteko, Kafka eta MQTT konektatzen dituen konektore bat definitu behar dugu, hau Kafka-connect izango da. Kontainer hau confluentinc/cp-Kafka-connect iruditik jaurtiko dugu, eta datu sarrera eta irteera egoki egiteko driverrak dituzten jar-ak dauden direktorioa definituko dugu. Honela, gure MQTT-Kafka konektorea gorde ahalko dugu eta Kafkan datuak jaso.

```
{ "usid":1600,"action": "A","ts": 252233498416973,"x":-2.034546,"y":14.139221,"z":
-0.7138214}
{"usid":1600,"action": "A","ts": 252233599124981,"x":-1.1264801,"y":9.675171,"z":
3.4104156}
{"usid":1600,"action": "A","ts": 252233699832989,"x":-5.0158386,"y":9.436523,"z":
-0.64216614}
{"usid":1600,"action": "A","ts": 252233800540996,"x":0.77900696,"y":5.7157593,"z":
-0.046676636}
{"usid":1600,"action": "A","ts": 252231383548809,"x":-1.4826813,"y":12.160522,"z":
4.1444244}
{"usid":1600,"action": "A","ts": 252231484256817,"x":-2.9437866,"y":12.0946045,"z":
0.52001953}
{"usid":1600,"action": "A","ts": 252231584964825,"x":-0.02571106,"y":8.554916,"z":
2.3144226}
^CProcessed a total of 10932 messages
oier@oier-R0G-Zephyrus-G14-GA401QE-GA401QE:~$ docker run --rm --network kafka_de
fault confluentinc/cp-kafka:5.1.0 kafka-console-consumer --bootstrap-server kafk
a:9092 --topic connect-custom --from-beginning
```

```
7 - zookeeper
8 kafka-connect:
9 image: confluentinc/cp-kafka-connect:5.1.0
10 hostname: kafka-connect
11 container_name: kafka-connect
12 ports:
13 - "8083:8083"
14 environment:
15 CONNECT_BOOTSTRAP_SERVERS: "kafka:9092"
16 CONNECT_REST_ADVERTISED_HOST_NAME: connect
17 CONNECT_REST_PORT: 8083
18 CONNECT_GROUP_ID: compose-connect-group
19 CONNECT_CONFIG_STORAGE_TOPIC: docker-connect-configs
20 CONNECT_OFFSET_STORAGE_TOPIC: docker-connect-offsets
21 CONNECT_STATUS_STORAGE_TOPIC: docker-connect-status
22 CONNECT_KEY_CONVERTER: org.apache.kafka.connect.json.Json
23 CONNECT_VALUE_CONVERTER: org.apache.kafka.connect.json.Json
24 CONNECT_INTERNAL_KEY_CONVERTER: "org.apache.kafka.conne
25 CONNECT_INTERNAL_VALUE_CONVERTER: "org.apache.kafka.con
26 CONNECT_CONFIG_STORAGE_REPLICATION_FACTOR: "1"
27 CONNECT_OFFSET_STORAGE_REPLICATION_FACTOR: "1"
28 CONNECT_STATUS_STORAGE_REPLICATION_FACTOR: "1"
29 CONNECT_PLUGIN_PATH: "/usr/share/java,/etc/kafka-conne
30 CONNECT_CONFLUENT_TOPIC_REPLICATION_FACTOR: 1
31 volumes:
32 - ./tmp/custom/jars:/etc/kafka-connect/jars
33 depends_on:
34 - zookeeper
35 - kafka
36 - mosquitto
```

Komando honekin jar dezakegu martxan Kafka-mqtt, eta dena ondo badoa, lehenago azaldu den subscriberrean bezala, datuak agertzen dira. Honek esan nahi du, MQTT-tik Kafkara datuak ongi pasatzen direla.

4.4. Flink shell-eko deiak

Kafkan jasotako datua SQL bezala jaso eta tratatzeko Flink erabili dezakegu. “`docker-compose exec sql-client ./sql-client.sh`” eginez, Flinkeko Shell-era sar gaitezke. Bertan, gure datuen gainean SQL deiak egin ahal izango ditugu.

```
pi@pi-hole-R0G-Zephyrus-G14-GA401QE-GA401QE:~/Documents/kafka$ docker compose exec sql-client ./sql-client.sh
Reading default environment from: file:/opt/flink/conf/sql-client-conf.yaml
No session environment specified.

Command history file path: /root/.flink-sql-history
```

Kafkako datuak hemen sartu nahi baditugu, taula bat sortu beharko dugu, gure json moduko datuen zutabeak errespetatuz. Gure kasuan, taula honi LANA izena jarri diogu eta 6 zelai definituko dizkiogu, json-ak dituen bezala. Bestalde, WITH zatian, Kafkatik “smart” topikoan dauden datuak hartu nahi ditugula esango diogu, eta hauek json formatuan daudela.

```
Flink SQL> CREATE TABLE LANA (
>   usid BIGINT,
>   action STRING,
>   ts BIGINT,
>   x BIGINT,
>   y BIGINT,
>   z BIGINT
> ) WITH (
>   'connector' = 'kafka', -- using kafka connector
>   'topic' = 'smart', -- kafka topic
>   'scan.startup.mode' = 'earliest-offset', -- reading from the beginning
>   'properties.bootstrap.servers' = 'kafka:9092', -- kafka broker address
>   'format' = 'json' -- the data format is json
> );
[INFO] Table has been created.
```

Lan honen azken zatia, datu hauek bistaratzea da, hau elasticsearchera esportatuta eta kibana erabilia lor dezakegu. Horretarako, interesgarriak zaizkigun datuak gordeko dituen taula bat sortu behar dugu eta hau elasticsearchera esportatu behar dugu.

```
Flink SQL> CREATE TABLE datuak (x BIGINT, y BIGINT, z BIGINT) WITH ('connector' = 'elasticsearch-7',
>   'hosts' = 'http://elasticsearch:9200',
>   'index' = 'datuak');
[INFO] Table has been created.

Flink SQL> INSERT INTO datuak
> SELECT x,y,z
> FROM LANA;
[INFO] Submitting SQL update statement to the cluster...
[INFO] Table update statement has been successfully submitted to the cluster:
Job ID: b8f1a7e289fdeedbea9a1ba3161ec54f

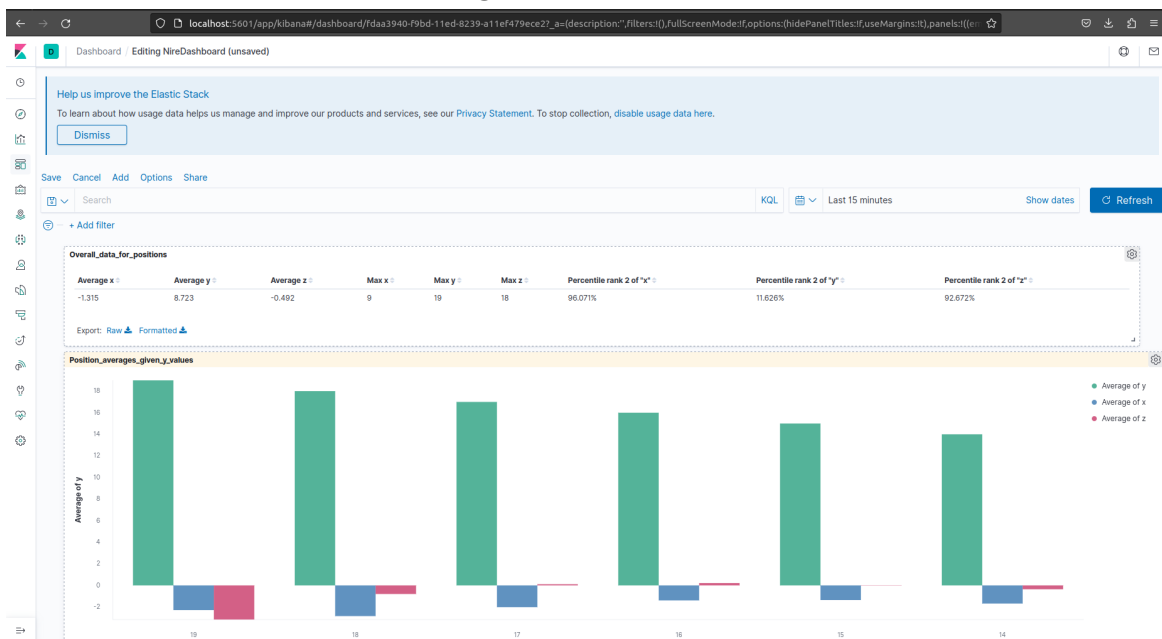
Flink SQL>
```

Ondoren, taula honetara beste taulatik interesgarriak zaizkigun datuak txertatuko ditugu. Gure kasua, x, y eta z zatiak gordeko dituen taula bat dugu. Horrela, kibana bezalako aplikazio bat erabiliz, hauek bistaratu ahalko ditugu.

4.5. Kibanako Dashboard-a

Behin elasticsearch-era datuak esportatuta, datu horiek Kibanan bistartzeko gai izango gara. Flink-eko shellean sortutako indizea zein den adieraziz Kibanan, datuen gainean deiak egingo ditugu.

Hona hemen sortutako dashboard-a gure datuekin:



Bertan, x, y eta z balioentzako batazbestekoak, balio maximoak eta bakoitzaren rank 2 pertzentilak atera ditugu. Bestalde, y balio desberdinetarako x eta z balioen fluktuazioa adierazten duen grafiko bat atera dugu ere bai.

5. Hobekuntza

5.1. Grafana erabiltzen saiatu

Elasticsearch-ekin erabili daitezkeen beste bistaratze-aplikazioak ikusiz, Grafanarekin aurkitu gara.

Kibana bezala, Grafanarentzat beste kontainer bat eginez docker compose fitxategian, eta bertan portu bat esleituko diogu sartu ahal izateko.

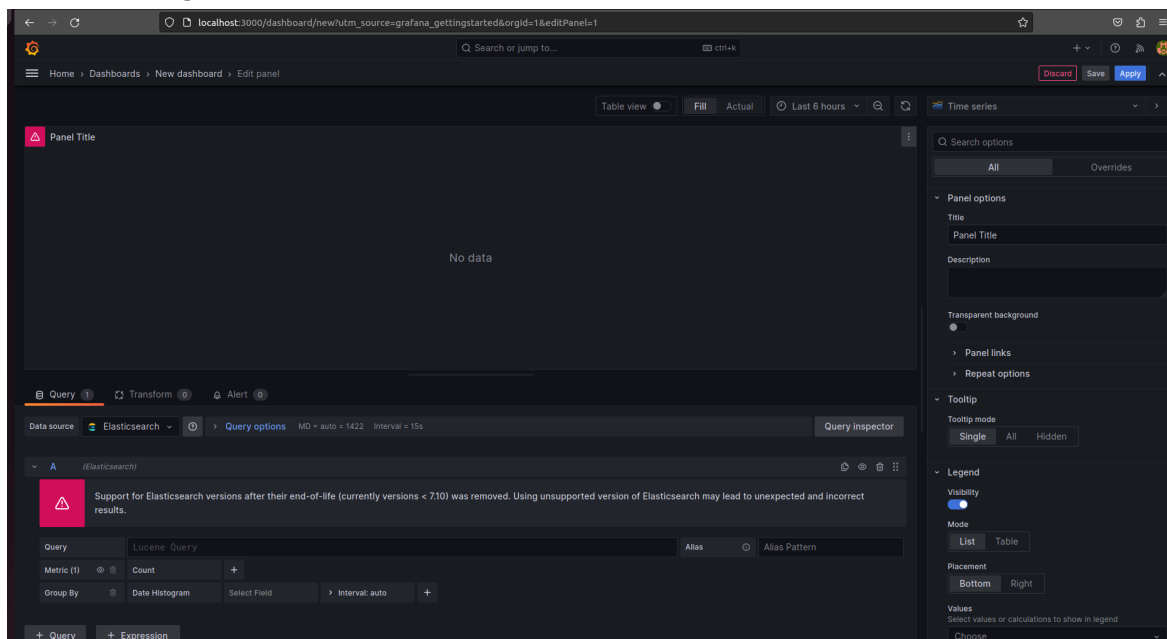
Kibanan ez bezala, Grafanan TimeStamp-a markatzen duen atributu bat izatea beharrezkoa zela ematen zuenez, lehenago egindako taulak (Flink-en) aldatu behar izan ditugu, ts atributua gehituz. Hona hemen erabilitako deia:

```

Flink SQL> INSERT INTO datuak
> SELECT x,y,z, HOUR(ts)
> FROM LANA;
[INFO] Submitting SQL update statement to the cluster...
[INFO] Table update statement has been successfully submitted to the cluster:
Job ID: 9ab8b50f12e9e842b9af537ac7a92c40

```

Taula egitea lortu badugu ere, ezin izan ditugu datu hauekin grafikoak lortu. Ondorengo irudian ikusten den bezala, erabiltzen ari garen elasticsearch-aren bertsioa ez dago sostendua Grafanarekin funtzionatzeko:



Bertsioarekin errore bat izan denez, txostenean sartzea erabaki dugu halere. Ez dugu bertsio desberdinekin probak egiteko denbora izan, baino Grafanarekin funtzionatu ahal izateko egin beharrezko pauso guztiak egin dira (Flink-eko taulak aldatu TimeStamp-a izan dadin, Compose fitxategian Grafana kontainera gehitu...).

6. Ondorioak

Lan honetan, MQTT, Kafka, Flink, Elasticsearch, Kibana eta Grafana elkarlanean nola jartzen diren ikasi dugu, datu fluxuak prozesatzen dituen soluzio global bat lortuz. Broker-ak zer diren eta nola funtzionatzen duten ere ikasi dugu, teknologia desberdinen arteko konexioa nola egin ikasiz aldi berean.

Aipatzekoa da, compose.yml fitxategietatik bolumenak kenduta ditugunez, gure proiektuan docker up egiten den bakoitzean, Flink atalean sortutako taulak berriz ere sortu behar direla, Kafkarekin eta Elasticsearch-ekin berririo konektatuz.

7. Erreferentziak

- Gure proiektuaren github biltegia:
https://github.com/oierIM/DMPA_MQTT-Kafka-Flink-ElasticSearch
- Kafka connect egiteko argibideak:
<https://www.baeldung.com/kafka-connect-mqtt-mongodb>
- Kafka konektoreak lortzeko:
<https://www.confluent.io/hub/>
- WISDM datasetaren informazioa:
<https://archive.ics.uci.edu/ml/datasets/WISDM+Smartphone+and+Smartwatch+Activity+and+Biometrics+Dataset+>
- Flink erabiltzeko argibideak:
<https://flink.apache.org/2020/07/28/flink-sql-demo-building-an-end-to-end-streaming-application/>