

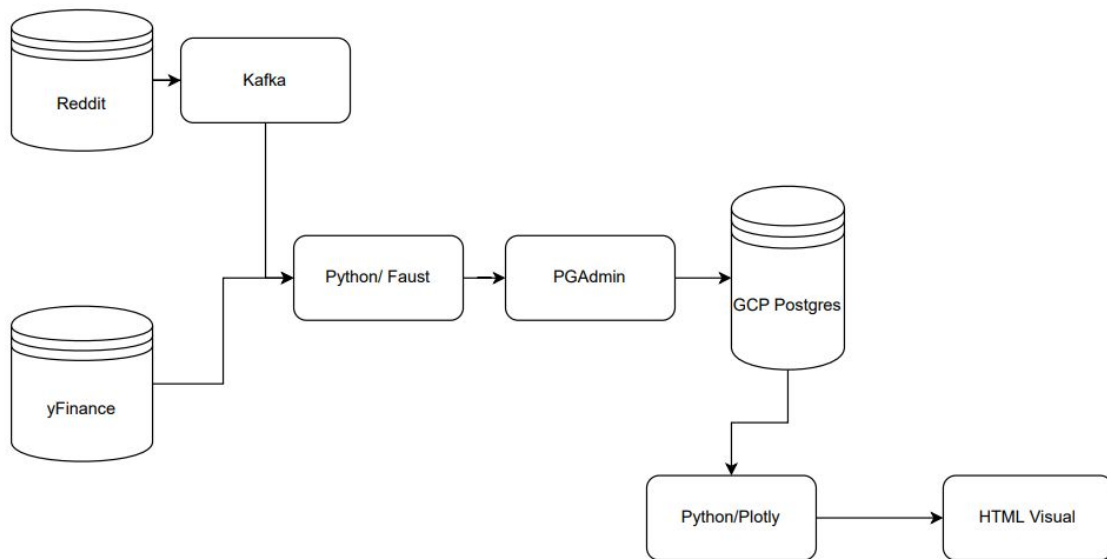
Crypto/Reddit Data Streaming

Michael Adams—Data Engineer

Overview

This project uses kafka, python, Google Cloud Postgres, Faust, confluent plugins, websockets, docker containers, and plotly for visuals. To the right is the data pipeline diagram.

Each component will be explained in the coming slides.



Data Sources

We retrieve reddit info using kafka. We use the landoop kafka docker container. In this container we add a reddit connector plugin and add our desired subreddits to the properties file as well as sending the data to a topic.

We retrieve yFinance data through a python websocket. Match the correct stock tickers with the data we want.

connect-standalone.properties
bootstrap.servers=localhost:9092

offset.storage.file.filename=/tmp/connect.offsets

key.converter=org.apache.kafka.connect.json.JsonConverter
key.converter.schemas.enable=false
value.converter=org.apache.kafka.connect.json.JsonConverter
value.converter.schemas.enable=false

internal.key.converter=org.apache.kafka.connect.json.JsonConverter
internal.key.converter.schemas.enable=false
internal.value.converter=org.apache.kafka.connect.json.JsonConverter
internal.value.converter.schemas.enable=false
Rest API
rest.port=8086
rest.host.name=127.0.0.1
this config is only for standalone workers
#offset.storage.file.filename=/standalone.offsets
offset.flush.interval.ms=10000

plugin.path=/connectors/

kafka-connect-reddit-source.properties

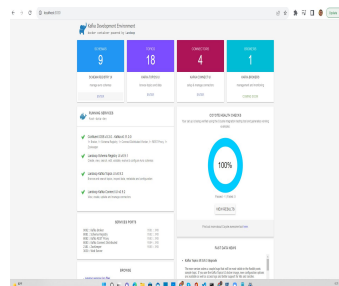
General properties for any connector

connector.class=com.github.c0urante.kafka.connect.reddit.RedditSourceConnector
name=reddit-source
tasks.max=2

Properties specifically for Reddit source connector

Posts and comments can be read from r/all
posts.subreddits=Bitcoin, cardano, solana, XRP, DOGE, ethereum, Tether, Binance
They can also be read from a specific subreddit or list of subreddits
posts.topic=reddit_topic

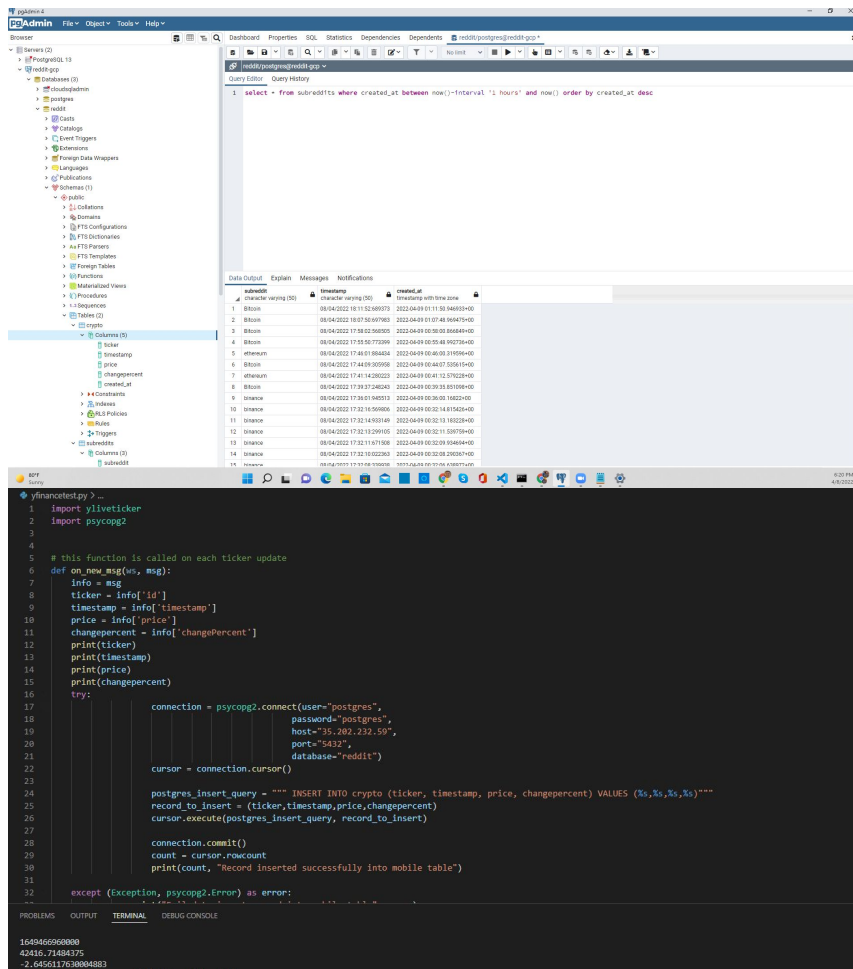
Enable this for debugging
reddit.log.http.requests=false



Data Transformation

We use Faust to take streaming data from kafka, add a timestamp and send it to the PGAdmin connected to our Google Cloud Postgres Server (we'll discuss this part in the next)

The yFinance python script directly puts the data into PGAdmin



Google Cloud Postgres

We create a postgres server on google cloud, noting IP addresses usernames and passwords to allow PGAdmin/Python to access database.

This is where all our data will be stored.

SQL

Instances

[+ CREATE INSTANCE](#)

[MIGRATE DATA](#)

[SHOW INFO PANEL](#)

Filter

Enter property name or value

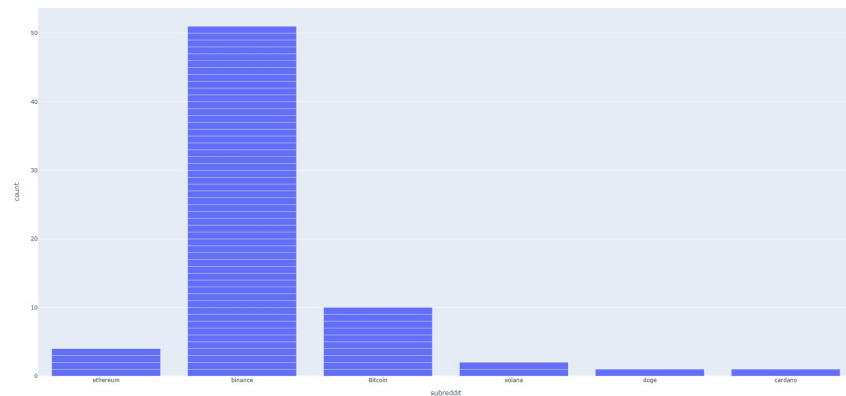
<input type="checkbox"/>	Instance ID ↑	Type	Public IP address	Private IP address	Instance connection name	High availability	Location	Storage used	Labels	Actions
<input type="checkbox"/>	reddit	PostgreSQL 13	35.202.232.59		second-capsule-342720-us-central1:reddit	ADD	us-central1-f	<div><div></div></div> 313 MB of 100 GB		

Visualize

We use python and plotly.express to visualize the data in HTML



```
1 # venv/bin/python
2 import plotly.express as px
3 import psycopg2
4 import pandas as pd
5
6 try:
7     conn = psycopg2.connect(user="postgres",
8                             password="postgres",
9                             host="ps-202.232.59",
10                             port="5432",
11                             database="reddit")
12
13     cur = conn.cursor()
14     query1 = "select * from crypto where created_at between now()-interval '1 hour' and now() and ticker = 'BTC-USD' order by created_at desc"
15     query2 = "select * from subreddit where created_at between now()-interval '1 hour' and now() order by created_at desc"
16     cur.execute(query1)
17     cryptoDF = pd.DataFrame(cur.fetchall(), columns = ['ticker', 'timestamp', 'price', 'changepercent', 'created_at'])
18     cur.execute(query2)
19     subredditDF = pd.DataFrame(cur.fetchall(), columns = ['subreddit', 'timestamp', 'created_at'])
20
21 except (Exception, psycopg2.Error) as error:
22     print("Error while fetching data from PostgreSQL", error)
23
24 finally:
25     # closing database connection.
26     if conn:
27         cur.close()
28         conn.close()
29         print("PostgreSQL connection is closed")
30
31
32 fig = px.line(cryptoDF, x="created_at", y="price", title="Bitcoin Price")
33 fig.show()
34 figure = px.bar(subredditDF, x="subreddit")
35 figure.show()
36
37
```



Conclusion

This project is easily adaptable. Pick whichever crypto you want to follow and add them to your kafka properties files and yFinance websocket.

Please take this project and make build to it!