

# Apache Flink para processamento em tempo real de dados IoT na nuvem

1<sup>st</sup> Luiza Martins de Freitas Cintra  
*Instituto de Informática*  
*Universidade Federal de Goiás*  
Goiânia, Goiás  
cintraluiza@discente.ufg.br

2<sup>nd</sup> Matheus Yosimura Lima  
*Instituto de Informática*  
*Universidade Federal de Goiás*  
Goiânia, Goiás  
matheusyosimura@discente.ufg.br

3<sup>rd</sup> Mateus Eduardo Silva Ribeiro  
*Instituto de Informática*  
*Universidade Federal de Goiás*  
Goiânia, Goiás  
mateus.eduardo@discente.ufg.br

4<sup>th</sup> Humberto Vitória Almeida Guadagnini  
*Instituto de Informática*  
*Universidade Federal de Goiás*  
Goiânia, Goiás  
hvittorioguadagnini@discente.ufg.br

5<sup>th</sup> Isabela de Queiroz Rodrigues  
*Instituto de Informática*  
*Universidade Federal de Goiás*  
Goiânia, Goiás  
Isabela.queiroz@discente.ufg.br

**Abstract**—O presente estudo visa promover uma resolução de problema relacionado ao processamento de dados em tempo real coletados por um dispositivo IoT. O paradigma envolve a captação de dados de temperatura em tempo real através da utilização de sensores. Por meio destes, haverá a emissão de um aviso alertando sobre picos de temperatura.

**Index Terms**—IoT, sensores, Sistema Operacional, Cloud

## I. INTRODUÇÃO E REVISÃO BIBLIOGRÁFICA

A “Internet das Coisas” surgiu recentemente como um novo conceito de rede, que abrange comunicações e processamento dos mais diversos equipamentos [1]. IoT, como é popularmente conhecida, é uma nova visão para a internet, em que ela passa a abranger não só computadores, como também, objetos do dia a dia [2].

As aplicações mais vislumbradas pela IoT, atualmente, é a telemetria, coleta de dados, atuação direta com os mais diversos objetos, relacionamentos em rede, entre muitos outros casos. Ela também aparece relacionada ao conceito de “Big Data”, muito difundido e estudado na contemporaneidade, cuja ideia está relacionada a uma expansão inimaginável de dados. Esses dados são gerados e coletados continuamente por objetos e computadores, indicando um enorme volume de memorização e processamento. A possibilidade de existência da IoT ocorre com um avanço específico do protocolo da internet, em que cada equipamento passa a ter o seu “endereço IP” próprio.

O Apache Flink é um mecanismo distribuído de código aberto para processamento com estado em conjunto de dados ilimitados (fluxos) e limitados (lotes) [3]. As aplicações de processamento são projetadas para serem executadas continuamente, tendo um mínimo de tempo de inatividade, e para que haja o processamento de dados à medida em que são ingeridos. O Apache Flink foi pensado para processamento de baixa latência, execução de cálculos na memória, alta

disponibilidade, remoção de pontos únicos de falhas e aumento da escala na horizontal.

Os atributos do Apache Flink incluem gerenciamento avançado de estado com garantias de consistência exatamente uma vez, semântica de processamento de tempo de evento com tratamento sofisticado de dados atrasados e fora de ordem [3]. Ele foi desenvolvido para priorizar o streaming. Justamente por conter todas essas características, ele foi pensado para a utilização de captações de picos de temperatura em ambientes, conectado em sensores.

## II. FUNDAMENTOS TEÓRICOS

O Apache Flink é usado para criar os mais diferentes tipos de aplicações de transmissão e em lotes, devido ao amplo conjunto de atributos. As aplicações mais comuns desenvolvidas com essa ferramenta são:

Aplicações orientadas por eventos: ingerindo eventos de um ou mais fluxos de eventos e executando cálculos, atualizações de estado ou ações externas.

Aplicações de análise de dados: extraindo informações e insights de dados. É executada consultando conjuntos de dados finitos e reexecutando as consultas ou alterando os resultados para que consiga incorporar novos dados. Com ele, é possível que ocorra a atualização de dados em tempo real, de forma constante.

Aplicações de pipelines de dados: transformando os dados a serem removidos de um armazenamento para outro. Tradicionalmente, extract-transform-load (ETL) é executado periodicamente, em lotes. Vide (Figura 1).

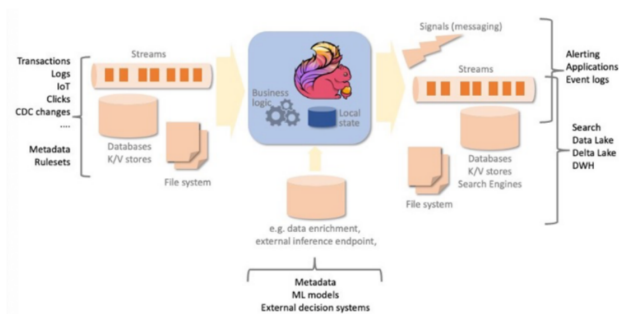


Fig. 1. Execução de aplicações.

O Flink é um mecanismo de processamento de fluxo de alto throughput e baixa latência. Uma aplicação do Flink consiste em um gráfico de fluxo de dados acíclico complexo e arbitrário, composto por fluxos e transformações. Os dados são ingeridos de uma ou mais fontes de dados e enviados para um ou mais destinos. Os sistemas de fonte e destino podem ser fluxos, filas de mensagens ou datastores e incluem arquivos, bancos de dados populares e mecanismos de busca. As transformações podem ser com estado, como agregações ao longo de janelas de tempo ou detecção de padrões complexos [3].

A tolerância a falhas é obtida por dois mecanismos separados: pontos de verificação automática e periódica do estado da aplicação, copiados para um armazenamento persistente, para permitir a recuperação automática em caso de falha; pontos de salvamento sob demanda, salvando uma imagem consistente do estado de execução, para permitir a pausa e retomada, atualizar ou bifurcar sua tarefa do Flink, mantendo o estado da aplicação em todas as pausas e reinicializações. Os mecanismos de ponto de verificação e ponto de salvamento são assíncronos, capturando um snapshot consistente do estado sem “parar o mundo”, enquanto a aplicação continua processando eventos. Vide (Figura 2).

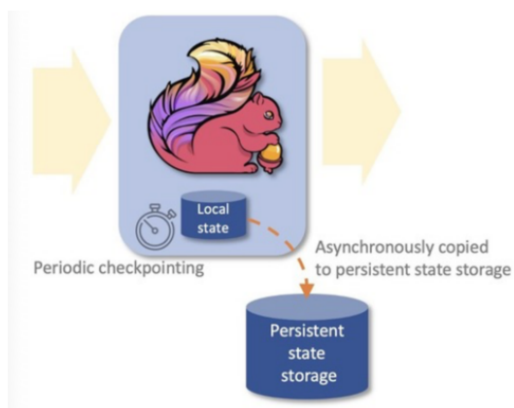


Fig. 2. Mecanismo de funcionamento.

Os benefícios de utilização do Apache Flink incluem:

Processamento de conjunto de dados ilimitados e limitados: o Apache Flink fornece algoritmos e estruturas de

dados para suportar processamento limitado e ilimitado por meio da mesma interface de programação. As aplicações que processam dados ilimitados são executadas continuamente. As aplicações que processam dados limitados encerrarão sua execução ao chegarem ao final dos conjuntos de dados de entrada;

Execução de aplicações em grande escala: o Apache Flink foi projetado para executar aplicações com estado em praticamente qualquer escala. O processamento é paralelo a milhares de tarefas, distribuído em várias máquinas, simultaneamente;

Desempenho na memória: os dados que fluem pela aplicação e pelo estado são particionados em várias máquinas. Portanto, o cálculo pode ser concluído acessando dados locais, geralmente na memória;

Consistência de estado exatamente uma vez: o Apache Flink garante a consistência do estado interno, mesmo em caso de falha e durante a parada e reinicialização da aplicação. O efeito de cada mensagem no estado interno é sempre aplicado exatamente uma vez, independentemente de a aplicação receber duplicatas da fonte de dados na recuperação ou na reinicialização;

Ampla variedade de conectores: vários conectores comprovados para sistemas populares de mensagens e streaming, armazenamentos de dados, mecanismos de pesquisa e sistemas de arquivos. Alguns exemplos são Apache Kafka, Amazon Kinesis Data Streams, Amazon SQS, Active MQ, Rabbit MQ, NiFi, OpenSearch e Elasticsearch, DynamoDB, HBase e qualquer banco de dados que forneça o cliente JDBC;

Vários níveis de abstrações: o Apache Flink oferece vários níveis de abstração para a interface de programação. Do SQL de streaming de alto nível e da API de tabela, usando abstrações familiares, como tabela, junções e agrupamento. A API DataStream oferece um nível mais baixo de abstração, mas também mais controle, com a semântica de fluxos, janelas e mapeamento. E, finalmente, a API ProcessFunction oferece controle rigoroso sobre o processamento de cada mensagem e controle direto do estado. Vide (Figura 3).

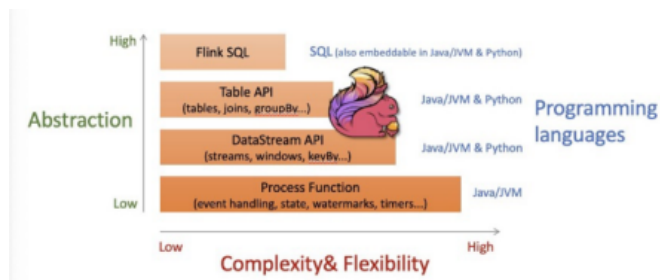


Fig. 3. Abstrações.

Várias linguagens de programação: o Apache Flink pode ser programado com várias linguagens, desde o SQL de streaming de alto nível até Python, Scala, Java, mas também com outras linguagens JVM, como Kotlin.

A IoT é um conceito que está fora do âmbito das tecnologias, pois não deriva delas, e sim as utiliza para cumprir uma série de funcionalidades. Em termos de protocolos diversos

e capazes de expandir a nova “rede de objetos”, há ainda protocolos adicionais do ambiente industrial e predial, pois os objetos existentes e já operando em tais sistemas não só podem como serão envolvidos na grande rede de interligação de objetos.

As funcionalidades do objeto na IoT já estão postas pelo mercado e pelas organizações normativas, e as tecnologias associadas continuam em desenvolvimento. A Intel, que fabrica os componentes para funcionalidade de “processamento” dos objetos, definindo assim suas características, acredita em “inteligente devices to deliver intelligence where needed and to acquire and filter data from the field” [5]. Esses sistemas irão compor sistemas inteligentes, integrando bilhões de dispositivos e provendo soluções e análises para valorizar soluções fim a fim dos clientes.

Para o lado dos desenvolvedores de aplicações, há a multiplicação de novas soluções, e novas expressões podem ser citadas: smart buildings, smart cities, smart transport, smart grid, smart energy, smart health, entre outras [5].

No conjunto das características do objeto na IoT, existem as seguintes atribuições [5]:

Processamento: capacidade de processamento computacional inserida no objeto, ou inteligência capaz de fazê-lo agir e responder às requisições de IoT;

Endereçamento: capacidade do objeto de ser encontrado na IoT, ou seja, ser localizado na rede por roteamento;

Identificação: identidade de cada objeto;

Localização: atributo relacionado ao local físico em que o objeto se encontra.

Ao conectar objetos com diferentes recursos a uma rede, potencializa-se o surgimento de novas aplicações. Neste sentido, conectar esses objetos à Internet significa criar a Internet das Coisas. Na IoT, os objetos podem prover comunicação entre usuários, dispositivos. Com isto emerge uma nova gama de aplicações, tais como coleta de dados de pacientes e monitoramento de idosos, sensoriamento de ambientes de difícil acesso e inóspitos, entre outras.

A arquitetura básica dos objetos inteligentes é composta por quatro unidades: processamento/memória, comunicação, energia e sensores/atuadores. São descritos a seguir a interligação dos componentes de dispositivos:

- **Unidade(s) de processamento/memória:** composta de uma memória interna para armazenamento de dados e programas, um microcontrolador e um conversor analógico-digital para receber sinais dos sensores. As CPUs utilizadas nesses dispositivos são, em geral, as mesmas utilizadas em sistemas embarcados e comumente não apresentam alto poder computacional. Frequentemente existe uma memória externa do tipo flash, que serve como memória secundária, por exemplo, para manter um “log” de dados. As características desejáveis para estas unidades são consumo reduzido de energia e ocupar o menor espaço possível.
- **Unidade(s) de comunicação:** consiste de pelo menos um canal de comunicação com ou sem fio, sendo mais comum o meio sem fio. Neste último caso, a maioria das

plata-formas usam rádio de baixo custo e baixa potência. Como consequência, a comunicação é de curto alcance e apresentam perdas frequentes. Esta unidade básica será objeto de estudo mais detalhado na próxima seção.

- **Fonte de energia:** responsável por fornecer energia aos componentes do objeto inteligente. Normalmente, a fonte de energia consiste de uma bateria (recarregável ou não) e um conversor AC-DC e tem a função de alimentar os componentes. Entretanto, existem outras fontes de alimentação como energia elétrica, solar e mesmo a captura de energia do ambiente através de técnicas de conversão (e.g., energia mecânica em energia elétrica), conhecidas como energy harvesting.
- **Unidade(s) de sensor(es)/atuador(es):** realizam o monitoramento do ambiente no qual o objeto se encontra. Os sensores capturam valores de grandezas físicas como temperatura, umidade, pressão e presença. Atualmente, existem literalmente centenas de sensores diferentes que são capazes de capturar essas grandezas. Atuadores, como o nome indica, são dispositivos que produzem alguma ação, atendendo a comandos que podem ser manuais, elétricos ou mecânicos.

A seção subsequente discute uma das aplicações práticas dessas ferramentas, destacando um caso de uso comum e implementação exemplar.

### III. METODOLOGIA

O objetivo central da pesquisa construir em tempo real um meio de comunicação para que analisar e gerar métricas pautadas nos dados de temperatura residenciais e periódicos coletados por sensores, com seu sistemas individuais de processamento, com a finalidade de detectar picos inesperados. Portanto, fundamentalmente a solução definida foi a construção de um pipeline na biblioteca de Python, PyFlink, que seja capaz de receber os dados coletados, analisá-los sistematicamente, gerar métricas em cima das coletas e ser capaz de apresentar dashboards de simples entendimento para informar o usuário da interface.

Inicialmente, foram necessárias a instalação de tecnologias localmente com finalidades diversas como apresentado na referência bibliográfica [5], como será descrito a seguir em ordem de prioridade:

- **Python:** A instalação do Python, primordialmente é indispensável por ser a linguagem trabalhada no core do desenvolvimento.
- **Docker:** O docker atua como uma camada de virtualização do sistema físico de um computador, capaz de carregar consigo inúmeras pré-configurações de bibliotecas que serão acrescentadas no projeto e sem a preocupação de configuração milimetricamente do ambiente de desenvolvimento.
- **Apache Flink:** O Apache é projetado para processar grandes volumes de dados em tempo real e também oferece capacidades de processamento de dados em lote.

- **Apache Kafka:** Será responsável por armazenar os dados coletados e transferir através de APIs para as tecnologias que solicitarem.
- **Kibana:** Uma interface de usuário (UI) que funciona em cima do Elasticsearch, proporcionando recursos avançados de visualização e análise de dados.
- **Elastic Search:** É um mecanismo de busca e análise de dados em tempo real.
- **Zookeeper:** Uma tecnologia responsável por centralizar a manutenção das configurações, sincronização e serviços.
- **Amazon EC2:** A máquina virtual da Aws servirá como ferramenta poderosa de nuvem neste projeto, pois não será necessário realizar as instalações e configurações localmente, todo o procedimento poderá ter um ambiente dedicado para sua realização, economizando espaço local, ganhando velocidade de processamento e garantindo a confiabilidade de um serviço em nuvem.

Há utilização do Kafka para armazenar dados de entrada sobre temperatura. Um script atuando como gerador de dados simples para simular o recebimento de dados de uma IoT, chamado `generate_source_data.py` é fornecido para gravar em tempo real novos registros no `temperature_msg` tópico Kafka. Cada registro foi estruturado da seguinte forma: `{"createTime": "2020-08-12 06:29:02", "orderId": 1597213797, "tempAmount": 29.159}`

- **createTime:** A hora da criação.
- **orderId:** O id da temperatura atual.
- **tempAmount:** O valor da temperatura.

Os dados de temperatura serão processados com PyFlink usando o script Python `temperature_msg.py`. Este script primeiro mapeará os `orderId` registros e calculará a soma dos valores das transações para cada id. O ambiente é baseado no Docker Compose. Ele usa uma imagem personalizada para ativar o Apache Flink (JobManager + TaskManager), Kafka+Zookeeper, o gerador de dados e contêineres Elastic-Search+Kibana.

Ademais, foi construída uma máquina virtual através do Amazon EC2, onde configurou-se uma imagem Docker e sendo baixadas todas as imagens das tecnologias mencionadas anteriormente. Após a ambientação estar completa os serviços disponibilizaram no url `http://localhost:8081` as métricas necessárias para análise. Durante todo o processo para o fortalecimento da metodologia foram realizados testes por grandes quantidades de tempo com temperaturas fictícias variadas.

#### IV. RESULTADOS E CONCLUSÕES

Por fim, obteve-se uma profunda quantidade amostral de dados capaz da geração de diversas métricas importantes para a análise. Sendo capaz, portanto, de aprofundar os saberes no tópico Apache Flink para processamento em tempo real de dados IoT na nuvem, sendo capaz futuramente de impulsionar trabalhos científicos interessados em utilizar dados oriundo de IoT e processados em tempo real na nuvem para automatizar ambientes residenciais ou empresariais.

Os resultados obtidos foram muito satisfatórios para o esperado do projeto, as métricas apresentadas estavam de acordo com o esperado do espaço amostral com uma eficácia de 100%. Portanto, conclui-se a oportunidade tecnológica para um cenário que apresenta como core este processamento, sem a necessidade da utilização de recursos locais.

A figura a seguir não estará apresentando um resultado oriundo dos resultados obtidos neste artigo, porém servirão como exemplo prático do produto final de entrega semelhante:

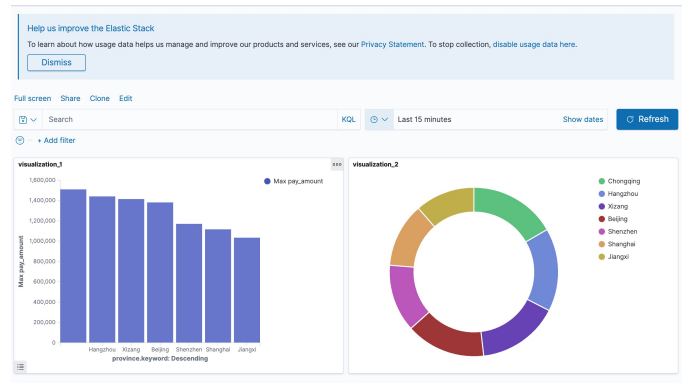


Fig. 4. Dashboard das métricas no Kibana.

#### REFERENCES

- [1] Faccioni Filho, Mauro. "Internet das coisas." Unisul Virtual(2016)
- [2] <https://aws.amazon.com/pt/what-is/apache-flink/>
- [3] Santos, Bruno P., et al. "Internet das coisas: da teoria à prática." Minicursos SBRC-Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos 31 (2016): 16.
- [4] Santaella, Lucia, et al. "Desvelando a Internet das coisas." Revista GEMInIS 4.2 (2013): 19-32.
- [5] <https://github.com/apache/flink-playgrounds/tree/master/pyflink-walkthrough>
- [6] <https://thingspeak.com/channels/public>
- [7] <https://medium.com/@daeynasvistas/a-iot-internet-das-coisas-surgiu-como-a-nova-gera>
- [8] <https://www.confluent.io/blog/apache-flink-stream-processing-use-cases-with-examples/>

## V. APÊNDICE

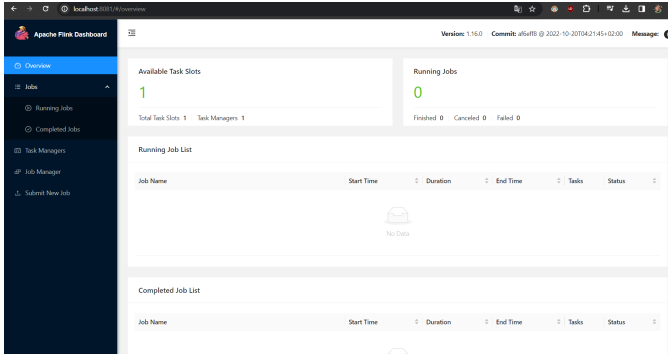


Fig. 5. Interface Web Apache Flink.

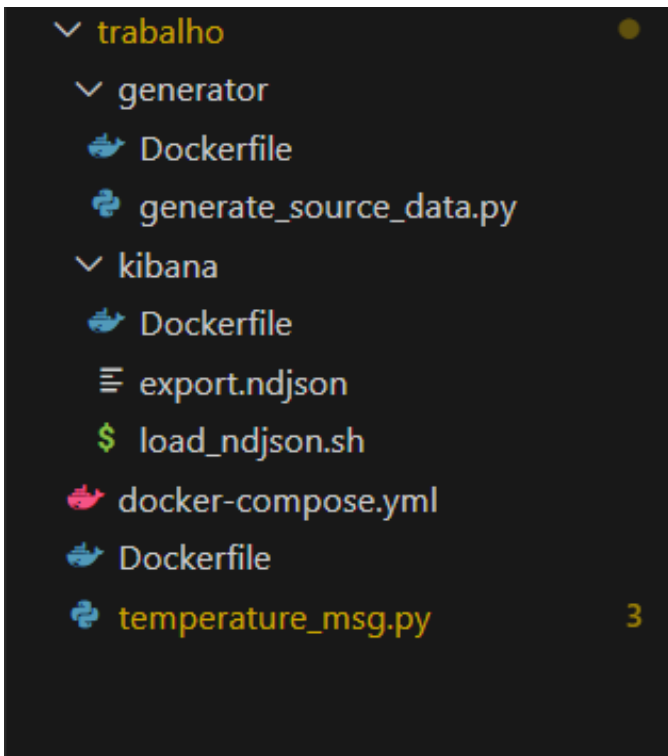


Fig. 6. Hierarquia da Árvore do Projeto.

```
1 version: '2.1'
2 services:
3   jobmanager:
4     build: .
5     image: pyflink/pyflink:1.16.0-scala_2.12
6     volumes:
7       - ./opt/trabalho
8     hostname: "jobmanager"
9     expose:
10       - "6123"
11     ports:
12       - "8081:8081"
13     command: jobmanager
14     environment:
15       - JOB_MANAGER_RPC_ADDRESS=jobmanager
16   taskmanager:
17     image: pyflink/pyflink:1.16.0-scala_2.12
18     volumes:
19       - ./opt/trabalho
20     expose:
21       - "6121"
22       - "6122"
23     depends_on:
24       - jobmanager
25     command: taskmanager
26     links:
27       - jobmanager:jobmanager
28     environment:
29       - JOB_MANAGER_RPC_ADDRESS=jobmanager
30   zookeeper:
31     image: wurstmeister/zookeeper:3.4.6
32     ulimits:
33       nofile:
34         soft: 65536
35         hard: 65536
36     ports:
37       - "2181:2181"
38   kafka:
39     image: wurstmeister/kafka:2.13-2.8.1
40     ports:
41       - "9092:9092"
42     depends_on:
43       - zookeeper
44     environment:
45       HOSTNAME_COMMAND: "route -n | awk '/UG[ \t]/{print \$2}'"
46       KAFKA_ZOOKEEPER_CONNECT: zookeeper:2181
47       KAFKA_CREATE_TOPICS: "temperature_msg:1:1"
48   generator:
49     build: generator
50     image: generator:1.0
51     depends_on:
52       - kafka
53   elasticsearch:
54     image: docker.elastic.co/elasticsearch/elasticsearch:7.8.0
55     environment:
56       - cluster.name=docker-cluster
57       - bootstrap.memory_lock=true
58       - "ES_JAVA_OPTS=-Xms512m -Xmx512m"
59       - discovery.type=single-node
60     ports:
61       - "9200:9200"
62       - "9300:9300"
63     ulimits:
64       memlock:
65         soft: -1
66         hard: -1
67       nofile:
68         soft: 65536
69         hard: 65536
70   kibana:
71     image: docker.elastic.co/kibana/kibana:7.8.0
72     ports:
73       - "5601:5601"
74     depends_on:
75       - elasticsearch
76   load-kibana-dashboard:
77     build: ./kibana
78     command: ['/bin/bash', '-c', 'cat /tmp/load/load_ndjson.sh | tr -d "\n" | bash']
79     depends_on:
80       - kibana
81
```

Fig. 7. Código para Configurar o Docker-Compose.

```

19 import random
20 import time, calendar
21 from random import randint
22 from kafka import KafkaProducer
23 from kafka import errors
24 from json import dumps
25 from time import sleep
26
27
28 def write_data(producer):
29     data_cnt = 20000
30     order_id = calendar.timegm(time.gmtime())
31     max_temp = 100
32     topic = "temperature_msg"
33
34     for i in range(data_cnt):
35         ts = time.strftime("%Y-%m-%d %H:%M:%S", time.localtime())
36         rd = random.random()
37         order_id += 1
38         temp_amount = max_temp * rd
39         cur_data = {"createTime": ts, "orderId": order_id, "tempAmount": temp_amount}
40         producer.send(topic, value=cur_data)
41         sleep(0.5)
42
43 def create_producer():
44     print("Connecting to Kafka brokers")
45     for i in range(0, 6):
46         try:
47             producer = KafkaProducer(bootstrap_servers=['kafka:9092'],
48                                     value_serializer=lambda x: dumps(x).encode('utf-8'))
49             print("Connected to Kafka")
50             return producer
51         except errors.NoBrokersAvailable:
52             print("Waiting for brokers to become available")
53             sleep(10)
54
55     raise RuntimeError("Failed to connect to brokers within 60 seconds")
56
57 if __name__ == '__main__':
58     producer = create_producer()
59     write_data(producer)
60

```

```

trabalho > temperature_msg.py > process_data
1 from pyflink.datastream import StreamExecutionEnvironment, TimeCharacteristic
2 from pyflink.table import StreamTableEnvironment, DataTypes
3 from pyflink.table.udf import udf
4
5 def process_data(t_env):
6     try:
7         print('2-')
8
9         create_kafka_source_ddl = """
10         CREATE TABLE temperature_msg(
11             createTime VARCHAR,
12             orderId BIGINT,
13             tempAmount DOUBLE
14         ) WITH (
15             'connector' = 'kafka',
16             'topic' = 'temperature_msg',
17             'properties.bootstrap.servers' = 'kafka:9092',
18             'properties.group.id' = 'test_3',
19             'scan.startup.mode' = 'latest-offset',
20             'format' = 'json'
21         )
22         """
23
24         print('4-')
25         create_es_sink_ddl = """
26         CREATE TABLE es_sink(
27             orderId BIGINT PRIMARY KEY,
28             tempAmount DOUBLE
29         ) WITH (
30             'connector' = 'elasticsearch-7',
31             'hosts' = 'http://elasticsearch:9200',
32             'index' = 'platform_temp_amount_1',
33             'document-id.key-delimiter' = '$',
34             'sink.bulk-flush.max-size' = '42mb',
35             'sink.bulk-flush.max-actions' = '32',
36             'sink.bulk-flush.interval' = '1000',
37             'sink.bulk-flush.backoff.delay' = '1000',
38             'format' = 'json'
39         )
40         """
41         print('5-')
42
43         t_env.execute_sql(create_kafka_source_ddl)
44         t_env.execute_sql(create_es_sink_ddl)
45         print('6-')
46
47         t_env.from_path("temperature_msg") \
48             .select("tempAmount") \
49             .group_by() \
50             .select("avg(tempAmount) as avg_temperature") \
51             .execute_insert("es_sink")
52
53         print('7-')
54     except Exception as e:
55         print(f'Error: {e}')
56
57 if __name__ == '__main__':
58     print('1-')
59
60     env = StreamExecutionEnvironment.get_execution_environment()
61     t_env = StreamTableEnvironment.create(stream_execution_environment=env)
62     t_env.get_config().get_configuration().set_boolean("python.fn-execution.memory.managed", True)
63
64     # Chamando a função de processamento
65     process_data(t_env)
66
67     print("Execution completed.")
68

```