

# 项目申请书

项目名称：openGauss 向量数据库 AI 数据 ETL 框架最佳实践输出

项目编号：2522d0324

项目导师：mpb159753

申请人：曾繁曦

申请人邮箱：[2628690042@qq.com](mailto:2628690042@qq.com)

## 一、项目背景

### 1.1 项目简述

完成 Apache Kafka、Firecrawl、VectorETL 对接 openGauss 最佳实践教程输出

### 1.2 项目产出要求

1. 基于 openGauss 部署指导文档，完成 openGauss Docker 部署
2. 参考 Milvus weaviate 等文档，完成 Apache Kafka、Firecrawl、VectorETL 对接 openGauss 最佳实践教程输出
3. 将文档提交至 openGauss 社区

### 1.3 项目技术要求

1. 了解 docker 构建部署
2. 了解 openGauss 的基础功能
3. 了解 Python 编程开发
4. 了解 RAG/LLMs 基础知识

### 1.4 项目成果仓库

- <https://gitcode.com/opengauss/examples>

## 二、Apache Kafka 对接 openGauss 实现方案

### 1. 架构设计

Apache Kafka 与 openGauss 的集成主要通过以下两种模式实现：

## 1.1 实时数据管道模式

代码块

```
1 [生产者应用] → [Kafka Topic] → [Kafka Connect] → [openGauss]
```

## 1.2 自定义消费者模式

代码块

```
1 [生产者应用] → [Kafka Topic] → [自定义消费者应用] → [openGauss]
```

## 2. 前置准备

### 2.1 环境要求

- Apache Kafka 2.8+ 集群
- openGauss 3.0+ 数据库（已安装vector扩展）
- JDK 11+ 环境
- 网络互通性验证

### 2.2 依赖组件

- Kafka Connect 框架
- openGauss JDBC 驱动（opengauss-jdbc-5.0.0.jar）
- Confluent JDBC Sink Connector 或 Debezium Connector

## 3. 核心实现方案

### 3.1 方案一：使用Kafka Connect JDBC Sink

#### 3.1.1 安装配置步骤

##### 1. 部署JDBC连接器：

代码块

```
1 # 下载Confluent JDBC连接器
2 wget https://packages.confluent.io/maven/io/confluent/kafka-connect-
   jdbc/10.7.0/kafka-connect-jdbc-10.7.0.tar.gz
3 tar -xzf kafka-connect-jdbc-10.7.0.tar.gz -C /usr/share/java/
```

## 2. 配置openGauss驱动：

代码块

```
1 cp opengauss-jdbc-5.0.0.jar /usr/share/java/kafka-connect-jdbc/
```

## 3. 创建Sink配置 (`jdbc-sink-opengauss.properties`) :

代码块

```
1 name=opengauss-sink
2 connector.class=io.confluent.connect.jdbc.JdbcSinkConnector
3 tasks.max=3
4 topics=source_topic
5 connection.url=jdbc:opengauss://opengauss_host:5432/target_db
6 connection.user=kafka_user
7 connection.password=db_password
8 auto.create=false
9 auto.evolve=true
10 insert.mode=upsert
11 pk.mode=record_key
12 pk.fields=id
13 table.name.format=target_table
14 transforms=unwrap
15 transforms.unwrap.type=io.debezium.transforms.ExtractNewRecordState
```

## 4. 启动连接器：

代码块

```
1 curl -X POST -H "Content-Type: application/json" \
2 --data @jdbc-sink-opengauss.properties \
3 http://kafka-connect:8083/connectors
```

## 3.2 方案二：自定义Java消费者应用

### 3.2.1 Java实现代码

代码块

```
1 import org.apache.kafka.clients.consumer.*;
2 import org.apache.kafka.common.serialization.StringDeserializer;
3 import java.sql.*;
4 import java.time.Duration;
5 import java.util.Collections;
```

```
6 import java.util.Properties;
7
8 public class OpenGaussKafkaConsumer {
9     private static final String TOPIC = "vector_data";
10    private static final String BOOTSTRAP_SERVERS = "kafka1:9092,kafka2:9092";
11    private static final String GROUP_ID = "opengauss-consumer-group";
12    private static final String OPENGAUSS_URL =
13        "jdbc:opengauss://opengauss:5432/vector_db";
14
15    public static void main(String[] args) {
16        // Kafka消费者配置
17        Properties props = new Properties();
18        props.put(ConsumerConfig.BOOTSTRAP_SERVERS_CONFIG, BOOTSTRAP_SERVERS);
19        props.put(ConsumerConfig.GROUP_ID_CONFIG, GROUP_ID);
20        props.put(ConsumerConfig.KEY_DESERIALIZER_CLASS_CONFIG,
21            StringDeserializer.class.getName());
22        props.put(ConsumerConfig.VALUE_DESERIALIZER_CLASS_CONFIG,
23            StringDeserializer.class.getName());
24        props.put(ConsumerConfig.AUTO_OFFSET_RESET_CONFIG, "earliest");
25        props.put(ConsumerConfig.ENABLE_AUTO_COMMIT_CONFIG, "false");
26
27        // openGauss连接池配置
28        PGSimpleDataSource ds = new PGSimpleDataSource();
29        ds.setURL(OPENGAUSS_URL);
30        ds.setUser("kafka_user");
31        ds.setPassword("db_password");
32
33        try (KafkaConsumer<String, String> consumer = new KafkaConsumer<>
34            (props)) {
35            consumer.subscribe(Collections.singleton(TOPIC));
36
37            while (true) {
38                ConsumerRecords<String, String> records =
39                consumer.poll(Duration.ofMillis(100));
40
41                try (Connection conn = ds.getConnection()) {
42                    conn.setAutoCommit(false);
43
44                    for (ConsumerRecord<String, String> record : records) {
45                        processRecord(conn, record);
46                    }
47
48                    conn.commit();
49                    consumer.commitSync();
50                } catch (SQLException e) {
51                    System.err.println("Database error: " + e.getMessage());
52                }
53            }
54        }
55    }
56}
```

```

48         }
49     }
50 }
51
52     private static void processRecord(Connection conn, ConsumerRecord<String,
53                                         String> record)
54             throws SQLException {
55         // 解析JSON数据 (示例使用简单字符串, 实际可用Jackson/Gson)
56         String[] parts = record.value().split("\\|");
57
58         try (PreparedStatement pstmt = conn.prepareStatement(
59                 "INSERT INTO vector_messages (id, content, vector) " +
60                 "VALUES (?, ?, ?::vector) " +
61                 "ON CONFLICT (id) DO UPDATE SET " +
62                 "content = EXCLUDED.content, " +
63                 "vector = EXCLUDED.vector")) {
64
64             pstmt.setString(1, record.key());
65             pstmt.setString(2, parts[0]);
66             pstmt.setObject(3, conn.createArrayOf("float",
67                     parseVector(parts[1])));
68             pstmt.executeUpdate();
69         }
70     }
71
72     private static Float[] parseVector(String vectorStr) {
73         String[] nums = vectorStr.split(",");
74         Float[] vector = new Float[nums.length];
75         for (int i = 0; i < nums.length; i++) {
76             vector[i] = Float.parseFloat(nums[i].trim());
77         }
78         return vector;
79     }

```

## 4. 高级配置与优化

### 4.1 批量写入优化

#### 代码块

```

1 # 在连接器配置中添加
2 batch.size=5000
3 max.retries=3
4 retry.backoff.ms=1000

```

## 4.2 向量数据处理

对于openGauss的向量类型，需要特殊处理：

代码块

```
1 -- 创建带向量字段的表
2 CREATE TABLE vector_data (
3     id VARCHAR(128) PRIMARY KEY,
4     content TEXT,
5     embedding VECTOR(384),
6     created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
7 );
```

在消费者中处理向量数据：

代码块

```
1 // 在processRecord方法中使用向量参数
2 pstmt.setObject(3, conn.createArrayOf("float", vectorArray));
```

## 5. 完整部署示例（Docker Compose）

代码块

```
1 version: '3.8'
2
3 services:
4     zookeeper:
5         image: confluentinc/cp-zookeeper:7.0.0
6         environment:
7             ZOOKEEPER_CLIENT_PORT: 2181
8
9     kafka:
10        image: confluentinc/cp-kafka:7.0.0
11        depends_on:
12            - zookeeper
13        ports:
14            - "9092:9092"
15        environment:
16            KAFKA_ZOOKEEPER_CONNECT: zookeeper:2181
17            KAFKA_ADVERTISED_LISTENERS: PLAINTEXT://kafka:9092
18            KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
19
20     kafka-connect:
21        image: confluentinc/cp-kafka-connect:7.0.0
```

```

22     depends_on:
23         - kafka
24     ports:
25         - "8083:8083"
26     environment:
27         CONNECT_BOOTSTRAP_SERVERS: kafka:9092
28         CONNECT_REST_ADVERTISED_HOST_NAME: kafka-connect
29         CONNECT_GROUP_ID: connect-cluster
30         CONNECT_CONFIG_STORAGE_TOPIC: connect-configs
31         CONNECT_OFFSET_STORAGE_TOPIC: connect-offsets
32         CONNECT_STATUS_STORAGE_TOPIC: connect-status
33         CONNECT_KEY_CONVERTER: org.apache.kafka.connect.json.JsonConverter
34         CONNECT_VALUE_CONVERTER: org.apache.kafka.connect.json.JsonConverter
35         CONNECT_CONFIG_STORAGE_REPLICATION_FACTOR: 1
36         CONNECT_OFFSET_STORAGE_REPLICATION_FACTOR: 1
37         CONNECT_STATUS_STORAGE_REPLICATION_FACTOR: 1
38         CONNECT_PLUGIN_PATH: "/usr/share/java,/usr/share/confluent-hub-
components"
39     volumes:
40         - ./opengauss-jdbc-5.0.0.jar:/usr/share/java/kafka-connect-
jdbc/opengauss-jdbc-5.0.0.jar
41         - ./jdbc-sink-opengauss.properties:/tmp/jdbc-sink-opengauss.properties
42
43     opengauss:
44         image: enmotech/opengauss:3.0.0
45         environment:
46             GS_PASSWORD: "OpenGauss@123"
47         ports:
48             - "5432:5432"
49         volumes:
50             - opengauss_data:/var/lib/opengauss
51
52     volumes:
53         opengauss_data:

```

启动后初始化连接器：

#### 代码块

```

1 docker-compose exec kafka-connect \
2     curl -X POST -H "Content-Type: application/json" \
3     --data @/tmp/jdbc-sink-opengauss.properties \
4     http://localhost:8083/connectors

```

## 三、Firecrawl 对接 openGauss 实现方案

# 1. 整体架构设计

Firecrawl (如 Firecrawl API 或开源爬虫框架) 与 openGauss 的集成主要通过以下流程实现：

代码块

```
1 [Firecrawl 爬取数据] → [数据预处理] → [向量化处理] → [存储到 openGauss]
```

## 2. 前置准备

### 2.1 环境要求

- 已部署的 Firecrawl 服务或 API 访问权限
- openGauss 3.0+ 数据库(已启用向量扩展)
- Python 3.8+ 环境
- 网络连通性(Firecrawl 可访问目标网站，应用服务器可访问 openGauss)

### 2.2 依赖安装

代码块

```
1 pip install opengauss-python psycopg2-binary requests beautifulsoup4 numpy sentence-transformers
```

## 3. 核心实现方案

### 3.1 方案一：直接对接 Firecrawl API

#### 3.1.1 Python 实现代码

代码块

```
1 import requests
2 import psycopg2
3 from sentence_transformers import SentenceTransformer
4 from bs4 import BeautifulSoup
5
6 # 配置参数
7 FIRE_CRAWL_API_KEY = "your_api_key"
8 OPENGUSS_CONN_STR = "host=opengauss_host port=5432 dbname=vector_db user=user password=pass"
9
10 # 初始化模型
```

```
11 model = SentenceTransformer('paraphrase-multilingual-MiniLM-L12-v2')
12
13 def crawl_and_store(url):
14     # 调用Firecrawl API
15     headers = {
16         "Authorization": f"Bearer {FIRE_CRAWL_API_KEY}",
17         "Content-Type": "application/json"
18     }
19     payload = {
20         "url": url,
21         "pageOptions": {
22             "includeHtml": True
23         }
24     }
25
26     response = requests.post(
27         "https://api.firecrawl.dev/v0/scrape",
28         headers=headers,
29         json=payload
30     )
31     data = response.json()
32
33     # 解析内容
34     soup = BeautifulSoup(data['html'], 'html.parser')
35     text_content = soup.get_text(separator=' ', strip=True)
36     title = data.get('title', 'No Title')
37
38     # 生成向量
39     vector = model.encode(text_content)
40
41     # 存储到openGauss
42     conn = psycopg2.connect(OPENGAUSS_CONN_STR)
43     cursor = conn.cursor()
44
45     try:
46         cursor.execute("""
47             INSERT INTO crawled_data
48                 (url, title, content, vector, created_at)
49             VALUES (%s, %s, %s, %s::vector, NOW())
50             ON CONFLICT (url) DO UPDATE SET
51                 title = EXCLUDED.title,
52                 content = EXCLUDED.content,
53                 vector = EXCLUDED.vector
54             """, (url, title, text_content, list(vector)))
55
56         conn.commit()
57         print(f"Successfully stored: {url}")
```

```
58     except Exception as e:
59         conn.rollback()
60         print(f"Error storing {url}: {str(e)}")
61     finally:
62         cursor.close()
63         conn.close()
64
65 # 示例使用
66 if __name__ == "__main__":
67     crawl_and_store("https://example.com")
```

## 3.2 方案二：使用自托管 Firecrawl 服务

### 3.2.1 部署与集成

#### 1. 部署 Firecrawl 服务:

代码块

```
1 docker run -p 8000:8000 -e API_KEY=your_secret_key firecrawl/firecrawl
```

#### 2. 修改对接代码:

代码块

```
1 # 将API端点改为本地服务
2 FIRE_CRAWL_API = "http://localhost:8000"
3
4 response = requests.post(
5     f"{FIRE_CRAWL_API}/v0/scrape",
6     headers=headers,
7     json=payload
8 )
```

## 4. 数据库准备

### 4.1 创建 openGauss 表结构

代码块

```
1 -- 启用向量扩展
2 CREATE EXTENSION vector;
3
4 -- 创建存储表
5 CREATE TABLE crawled_data (
```

```
6     id SERIAL PRIMARY KEY,
7     url VARCHAR(1024) UNIQUE,
8     title TEXT,
9     content TEXT,
10    vector VECTOR(384), -- 维度根据模型调整
11    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
12    updated_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
13 );
14
15 -- 创建向量索引
16 CREATE INDEX ON crawled_data USING ivfflat (vector vector_l2_ops) WITH (lists =
17 100);
```

## 5. 高级功能实现

### 5.1 批量爬取与存储

代码块

```
1 def batch_crawl(urls, batch_size=10):
2     for i in range(0, len(urls), batch_size):
3         batch = urls[i:i + batch_size]
4         with ThreadPoolExecutor(max_workers=5) as executor:
5             executor.map(crawl_and_store, batch)
```

### 5.2 增量爬取策略

代码块

```
1 def get_last_crawl_time(url):
2     conn = psycopg2.connect(OPENGAUSS_CONN_STR)
3     cursor = conn.cursor()
4     cursor.execute("SELECT updated_at FROM crawled_data WHERE url = %s",
5 (url,))
6     result = cursor.fetchone()
7     cursor.close()
8     conn.close()
9     return result[0] if result else None
10
11 def smart_crawl(url):
12     last_crawl = get_last_crawl_time(url)
13     if last_crawl and (datetime.now() - last_crawl).days < 1:
14         print(f"Skipping recently crawled URL: {url}")
15         return
```

## 5.3 向量相似度搜索集成

代码块

```

1  def search_similar_content(query, top_k=5):
2      conn = psycopg2.connect(OPENGAUSS_CONN_STR)
3      cursor = conn.cursor()
4
5      query_vector = model.encode(query)
6
7      cursor.execute("""
8          SELECT url, title, content,
9              vector <-> %s::vector AS similarity
10         FROM crawled_data
11        ORDER BY similarity ASC
12        LIMIT %s
13      """, (list(query_vector), top_k))
14
15      results = cursor.fetchall()
16      cursor.close()
17      conn.close()
18
19      return [
20          {
21              'url': r[0],
22              'title': r[1],
23              'content': r[2][:200] + '...', # 摘要
24              'score': float(r[3])
25          } for r in results]

```

## 6. 完整部署方案

### 6.1 Docker Compose 部署

代码块

```

1  version: '3.8'
2
3  services:
4      firecrawl:
5          image: firecrawl/firecrawl
6          environment:
7              API_KEY: your_firecrawl_key
8          ports:

```

```

9      - "8000:8000"
10
11 crawler:
12   build: .
13   environment:
14     FIRE_CRAWL_API: "http://firecrawl:8000"
15     OPENGAUSS_HOST: "opengauss"
16   depends_on:
17     - firecrawl
18     - opengauss
19
20 opengauss:
21   image: opengauss:3.0
22   environment:
23     GS_PASSWORD: "OpenGauss@123"
24   ports:
25     - "5432:5432"
26   volumes:
27     - opengauss_data:/var/lib/opengauss
28
29 volumes:
30   opengauss_data:

```

## 6.2 自动化爬取调度

使用 Celery 实现定时任务：

### 代码块

```

1  from celery import Celery
2
3  app = Celery('crawler', broker='redis://redis:6379/0')
4
5  @app.task
6  def scheduled_crawl(url):
7      crawl_and_store(url)
8
9  # 配置定时任务
10 app.conf.beat_schedule = {
11     'daily-crawl': {
12         'task': 'scheduled_crawl',
13         'schedule': 86400.0,  # 每天
14         'args': ('https://example.com',)
15     },
16 }

```

# 四、VectorETL 对接 openGauss 实现方案

## 1. 整体架构设计

VectorETL 与 openGauss 的集成主要通过以下数据处理流程实现：

代码块

```
1 [原始数据源] → [VectorETL 处理] → [向量化转换] → [openGauss 向量存储] → [向量检索服务]
```

## 2. 前置准备

### 2.1 环境要求

- VectorETL 运行环境 (Python 3.8+ 或 Docker)
- openGauss 3.0+ 数据库 (已安装 vector 扩展)
- 向量模型服务 (如 Sentence-Transformers、HuggingFace 或自定义模型)
- 网络互通性验证

### 2.2 依赖安装

代码块

```
1 # Python 环境
2 pip install vector-etl opengauss-python psycopg2-binary numpy sentence-
   transformers pandas
3
4 # 或使用 Docker
5 docker pull vectordotdev/vector-etl
```

## 3. 核心实现方案

### 3.1 方案一：使用 Python SDK 直接对接

#### 3.1.1 基础实现代码

代码块

```
1 from vector_etl import VectorETL
2 import psycopg2
3 from sentence_transformers import SentenceTransformer
4 import numpy as np
```

```
5
6 # 配置参数
7 OPENGAUSS_CONN_STR = "host=openauss_host port=5432 dbname=vector_db user=user
password=pass"
8 MODEL_NAME = "paraphrase-multilingual-MiniLM-L12-v2"
9
10 # 初始化组件
11 model = SentenceTransformer(MODEL_NAME)
12 etl = VectorETL()
13
14 def process_and_store(data_source):
15     # 1. 数据抽取
16     raw_data = etl.extract(data_source)
17
18     # 2. 数据转换
19     processed_data = []
20     for item in raw_data:
21         # 文本清洗和预处理
22         clean_text = etl.clean_text(item['content'])
23
24         # 向量化
25         vector = model.encode(clean_text)
26
27         processed_data.append({
28             'id': item['id'],
29             'content': clean_text,
30             'metadata': item.get('metadata', {}),
31             'vector': vector
32         })
33
34     # 3. 加载到openGauss
35     conn = psycopg2.connect(OPENGAUSS_CONN_STR)
36     cursor = conn.cursor()
37
38     try:
39         # 批量插入
40         records = [(d['id'], d['content'], d['metadata'], list(d['vector']))]
41         for d in processed_data:
42             cursor.executemany("""
43                 INSERT INTO vector_docs (id, content, metadata, vector)
44                 VALUES (%s, %s, %s::jsonb, %s::vector)
45                 ON CONFLICT (id) DO UPDATE SET
46                     content = EXCLUDED.content,
47                     metadata = EXCLUDED.metadata,
48                     vector = EXCLUDED.vector
49                 """, records)
```

```

50         conn.commit()
51         print(f"成功插入/更新 {len(processed_data)} 条记录")
52     except Exception as e:
53         conn.rollback()
54         print(f"数据库操作失败: {str(e)}")
55     finally:
56         cursor.close()
57         conn.close()
58
59 # 示例使用
60 if __name__ == "__main__":
61     data_source = {
62         "type": "csv",
63         "path": "data/documents.csv",
64         "options": {"delimiter": ","}
65     }
66     process_and_store(data_source)

```

## 3.2 方案二：使用 VectorETL 配置文件

### 3.2.1 配置 YAML 文件 (`vector-etl-config.yaml`)

代码块

```

1  version: "1"
2  sources:
3      - name: "document_source"
4          type: "csv"
5          path: "data/documents.csv"
6          options:
7              delimiter: ","
8
9  transformations:
10     - name: "text_clean"
11         type: "text_processing"
12         operations:
13             - remove_special_chars
14             - normalize_whitespace
15             - lowercase
16
17     - name: "vectorize"
18         type: "embedding"
19         model: "sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2"
20         input_field: "content"
21         output_field: "vector"
22

```

```
23 sinks:
24   - name: "opengauss_sink"
25     type: "opengauss"
26     connection:
27       host: "opengauss_host"
28       port: 5432
29       database: "vector_db"
30       user: "user"
31       password: "pass"
32     table: "vector_docs"
33     mapping:
34       id: "id"
35       content: "content"
36       metadata: "metadata"
37       vector: "vector"
38     batch_size: 1000
```

### 3.2.2 运行命令

代码块

```
1 vector-etl --config vector-etl-config.yaml
```

## 4. 数据库准备

### 4.1 创建 openGauss 表结构

代码块

```
1 -- 启用向量扩展
2 CREATE EXTENSION vector;
3
4 -- 创建存储表
5 CREATE TABLE vector_docs (
6   id VARCHAR(128) PRIMARY KEY,
7   content TEXT,
8   metadata JSONB,
9   vector VECTOR(384), -- 维度根据模型调整
10  created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
11  updated_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
12 );
13
14 -- 创建向量索引
15 CREATE INDEX ON vector_docs USING ivfflat (vector vector_l2_ops) WITH (lists =
16 100);
```

```
16  
17 -- 创建全文检索索引 (可选)  
18 CREATE EXTENSION pg_trgm;  
19 CREATE INDEX ON vector_docs USING gin (content gin_trgm_ops);
```

## 5. 高级功能实现

### 5.1 增量数据处理

代码块

```
1 def get_last_processed_id():  
2     conn = psycopg2.connect(OPENGAUSS_CONN_STR)  
3     cursor = conn.cursor()  
4     cursor.execute("SELECT MAX(id) FROM vector_docs")  
5     last_id = cursor.fetchone()[0]  
6     cursor.close()  
7     conn.close()  
8     return last_id  
9  
10 def incremental_update(data_source):  
11     last_id = get_last_processed_id()  
12  
13     # 修改VectorETL提取逻辑，只获取新数据  
14     raw_data = etl.extract(  
15         data_source,  
16         filter_logic=f"id > '{last_id}'" if last_id else None  
17     )  
18  
19     if raw_data:  
20         process_and_store(raw_data)
```

### 5.2 多模型向量支持

代码块

```
1 class MultiVectorETL:  
2     def __init__(self):  
3         self.text_model = SentenceTransformer('paraphrase-MiniLM-L6-v2')  
4         self.image_model = load_image_model() # 自定义图像模型  
5  
6     def process_item(self, item):  
7         vectors = []  
8  
9         if item['type'] == 'text':
```

```
10         vectors['text_vector'] = self.text_model.encode(item['content'])
11     elif item['type'] == 'image':
12         vectors['image_vector'] =
13             self.image_model.encode(item['image_data'])
14
15
16 # 在存储时需要调整表结构和SQL
```

## 5.3 分布式处理

代码块

```
1  from multiprocessing import Pool
2
3  def parallel_process(data_chunk):
4      # 每个进程独立的数据库连接
5      local_etl = VectorETL()
6      local_model = SentenceTransformer(MODEL_NAME)
7      # 处理逻辑...
8      return processed_chunk
9
10 if __name__ == "__main__":
11     raw_data = etl.extract(data_source)
12     chunks = np.array_split(raw_data, 4)    # 分为4份
13
14     with Pool(processes=4) as pool:
15         results = pool.map(parallel_process, chunks)
16
17     # 合并结果并存储
18     all_processed = [item for chunk in results for item in chunk]
19     batch_insert(all_processed)
```

## 6. 完整部署方案

### 6.1 Docker Compose 部署

代码块

```
1  version: '3.8'
2
3  services:
4      vector-etl:
5          image: vectordotdev/vector-etl:latest
6          volumes:
```

```
7     - ./config:/app/config
8     - ./data:/app/data
9   environment:
10     CONFIG_PATH: "/app/config/vector-etl-config.yaml"
11   depends_on:
12     - opengauss
13
14   opengauss:
15     image: opengauss:3.0
16     environment:
17       GS_PASSWORD: "OpenGauss@123"
18     ports:
19       - "5432:5432"
20     volumes:
21       - opengauss_data:/var/lib/opengauss
22
23   volumes:
24     opengauss_data:
```

## 6.2 Kubernetes 部署示例

代码块

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: vector-etl
5  spec:
6    replicas: 3
7    selector:
8      matchLabels:
9        app: vector-etl
10   template:
11     metadata:
12       labels:
13         app: vector-etl
14     spec:
15       containers:
16         - name: vector-etl
17           image: vectordotdev/vector-etl:latest
18           volumeMounts:
19             - name: config
20               mountPath: /app/config
21             env:
22               - name: CONFIG_PATH
23                 value: "/app/config/vector-etl-config.yaml"
```

```

24     volumes:
25       - name: config
26         configMap:
27           name: vector-etl-config
28   ---
29   apiVersion: v1
30   kind: ConfigMap
31   metadata:
32     name: vector-etl-config
33   data:
34     vector-etl-config.yaml: |
35       # 这里放置之前的YAML配置内容

```

## 五、项目开发时间规划

### 1. 项目总体时间安排

**总周期：**2025年7月1日 - 2025年9月30日（共13周）

### 2. 详细时间计划表

**第一阶段：环境准备与基础部署 (2025.07.01-2025.07.14)**

任务	时间	交付物
1.1 openGauss Docker化部署研究	07.01-07.03	部署方案文档
1.2 Docker环境配置与测试	07.04-07.07	可运行的Docker环境
1.3 openGauss向量扩展验证	07.08-07.10	测试报告
1.4 Kafka集群搭建	07.11-07.12	Kafka测试环境
1.5 Firecrawl环境准备	07.13-07.14	爬虫测试环境

**阶段目标：**完成所有基础环境搭建并通过验证

**第二阶段：核心功能开发 (2025.07.15-2025.08.25)**

**模块一：Kafka集成 (07.15-07.28)**

任务	时间	交付物

2.1.1 Kafka Connect配置开发	07.15-07.18	连接器配置文件
2.1.2 自定义消费者实现	07.19-07.22	Java消费者代码
2.1.3 向量数据特殊处理	07.23-07.25	向量转换模块
2.1.4 性能测试与优化	07.26-07.28	性能测试报告

## 模块二：Firecrawl集成 (07.29-08.11)

任务	时间	交付物
2.2.1 网页抓取管道开发	07.29-08.01	Python爬取脚本
2.2.2 内容向量化处理	08.02-08.05	向量化服务
2.2.3 增量爬取机制	08.06-08.08	增量处理模块
2.2.4 分布式爬取测试	08.09-08.11	分布式方案文档

## 模块三：VectorETL集成 (08.12-08.25)

任务	时间	交付物
2.3.1 ETL流程设计	08.12-08.14	ETL流程图
2.3.2 多模型向量支持	08.15-08.18	多模型适配器
2.3.3 批量处理优化	08.19-08.22	批量处理模块
2.3.4 端到端测试	08.23-08.25	集成测试报告

**阶段目标：**完成所有核心功能模块开发并通过集成测试

## 第三阶段：文档与优化 (2025.08.26-2025.09.15)

任务	时间	交付物
3.1 部署指南编写	08.26-08.28	部署文档
3.2 最佳实践手册	08.29-09.01	实践指南
3.3 性能调优白皮书	09.02-09.05	优化方案

3.4 常见问题整理	09.06-09.08	FAQ文档
3.5 安全配置指南	09.09-09.11	安全手册
3.6 代码审查与重构	09.12-09.15	优化后的代码

阶段目标：完成所有技术文档并优化代码质量

## 第四阶段：社区提交与收尾 (2025.09.16-2025.09.30)

任务	时间	交付物
4.1 文档格式转换	09.16-09.18	社区标准文档
4.2 示例代码整理	09.19-09.21	可运行示例包
4.3 社区PR提交	09.22-09.24	提交记录
4.4 根据反馈修改	09.25-09.27	修订版本
4.5 项目总结报告	09.28-09.30	总结报告

阶段目标：完成社区贡献并通过验收

## 3. 关键里程碑节点

1. 07.14 - 基础环境验收
2. 07.28 - Kafka集成完成
3. 08.11 - Firecrawl集成完成
4. 08.25 - VectorETL集成完成
5. 09.15 - 文档初稿完成
6. 09.30 - 社区提交完成

## 4. 每周工作时间安排

- 工作日：每天6小时（9:00-12:00, 14:00-17:00）
- 总开发时间：约400小时（13周×30小时）

## 5. 风险管理计划

风险	应对措施	负责人

组件兼容性问题	提前建立版本矩阵	本人
性能不达标	预留2周优化时间	本人
文档不符合要求	提前与社区沟通格式	本人
个人时间冲突	关键阶段增加周末工时	本人

## 6. 开发环境配置

- 硬件：笔记本(i7/16GB/512GB SSD) + 云测试服务器(4C8G)
- 软件：
  - Docker Desktop
  - IntelliJ IDEA/PyCharm
  - PostgreSQL/OpenGauss客户端工具
  - Kafka/Zookeeper套件

## 7. 质量保证措施

1. 每周日进行代码审查
2. 每个模块完成后进行单元测试
3. 关键功能进行性能压测
4. 文档经过三次校审