_____
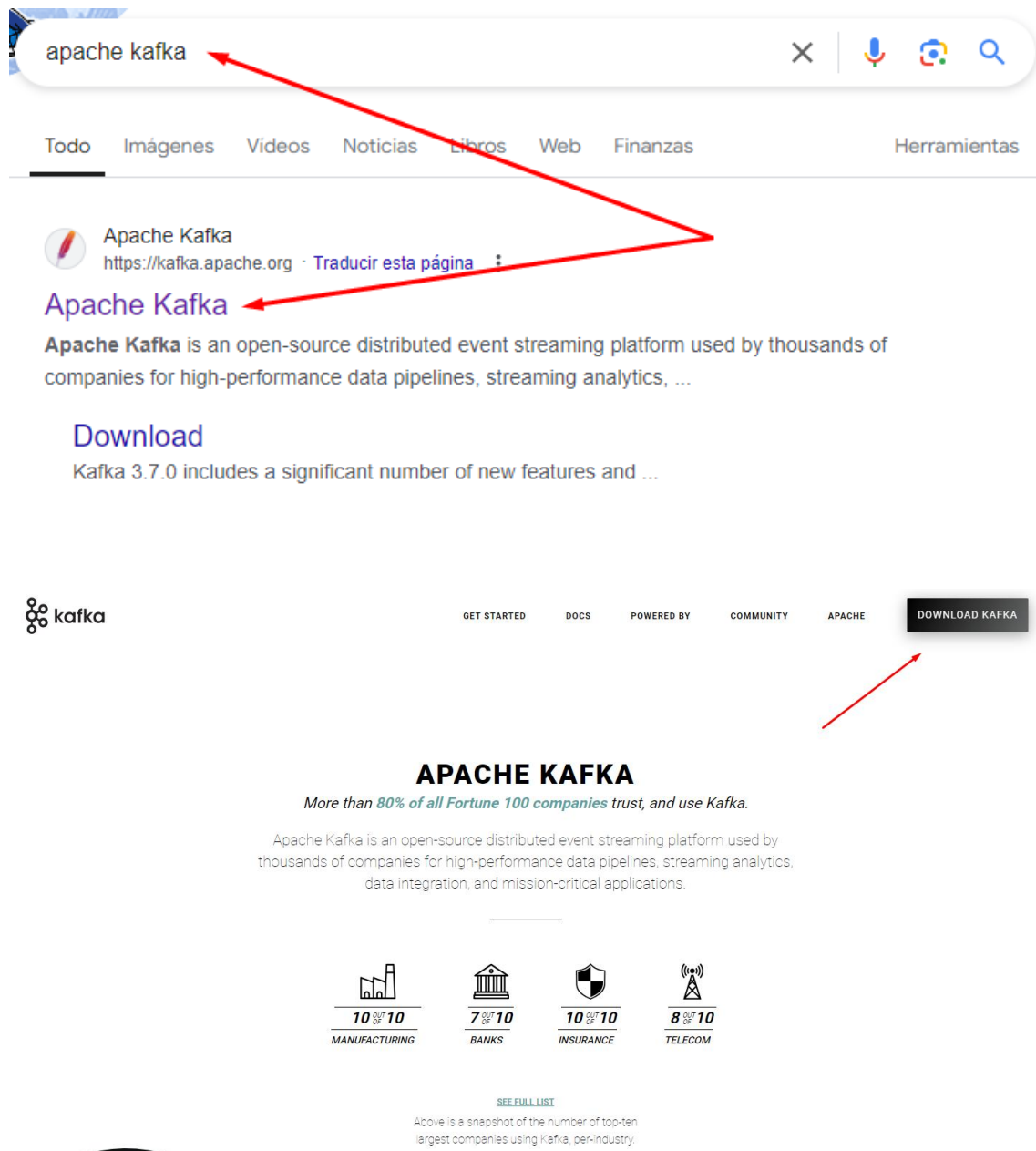
## 5075. Big Data aplicado - 1ª Evaluación (RA 1 – CE b)

## Unidad Didáctica 1. Estrategias de ingestión y almacenamiento de datos en Big Data

Práctica 7. Practicando la ingesta de datos en Streaming.

Paso 1. Descargamos e instalamos Apache Kafka.

# DOWNLOAD

The project goal is to have 3 releases a year, which means a release every 4 months. Bugfix releases are made as needed for supported releases only. It is possible to verify every download by following these procedures and using these KEYS.

## SUPPORTED RELEASES

### 3.8.0

- Released July 29, 2024
- Release Notes
- Docker image: apache/kafka:3.8.0.
- Docker Native image: apache/kafka-native:3.8.0.
- Source download: kafka-3.8.0-src.tgz (asc, sha512)
- Binary downloads:

    - Scala 2.12 - kafka_2.12-3.8.0.tgz (asc, sha512)
    - Scala 2.13 - kafka_2.13-3.8.0.tgz (asc, sha512)

    We build for multiple versions of Scala. This only matters if you are using Scala and you want a version built for the same Scala version you use. Otherwise any version should work (2.13 is recommended).

    Kafka 3.8.0 includes a significant number of new features and fixes. For more information, please read our blog post and the detailed Release Notes.

```
hadoop@hadoop-virtualbox:~$ wget https://downloads.apache.org/kafka/3.8.0/kafka_2.13-3.8.0.tgz
--2024-08-31 11:05:14--  https://downloads.apache.org/kafka/3.8.0/kafka_2.13-3.8.0.tgz
Resolviendo downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.208.237
Conectando con downloads.apache.org (downloads.apache.org)[135.181.214.104]:443... conectado.
Petición HTTP enviada, esperando respuesta... 200 OK
Longitud: 120735482 (115M) [application/x-gzip]
Guardando como: 'kafka_2.13-3.8.0.tgz'

kafka_2.13-3.8.0.tgz                    100%[===================================================

2024-08-31 11:05:18 (35,1 MB/s) - 'kafka_2.13-3.8.0.tgz' guardado [120735482/120735482]
```

```
hadoop@hadoop-virtualbox:~$ tar -xvzf kafka_2.13-3.8.0.tgz
kafka_2.13-3.8.0/
kafka_2.13-3.8.0/LICENSE
kafka_2.13-3.8.0/NOTICE
kafka_2.13-3.8.0/bin/
kafka_2.13-3.8.0/bin/connect-distributed.sh
kafka_2.13-3.8.0/bin/connect-mirror-maker.sh
kafka_2.13-3.8.0/bin/connect-plugin-path.sh
kafka_2.13-3.8.0/bin/connect-standalone.sh
kafka_2.13-3.8.0/bin/kafka-acls.sh
kafka_2.13-3.8.0/bin/kafka-broker-api-versions.sh
kafka_2.13-3.8.0/bin/kafka-client-metrics.sh
kafka_2.13-3.8.0/bin/kafka-cluster.sh
kafka_2.13-3.8.0/bin/kafka-configs.sh
kafka_2.13-3.8.0/bin/kafka-console-consumer.sh
kafka_2.13-3.8.0/bin/kafka-console-producer.sh
kafka_2.13-3.8.0/bin/kafka-consumer-groups.sh
kafka_2.13-3.8.0/bin/kafka-consumer-perf-test.sh
kafka_2.13-3.8.0/bin/kafka-delegation-tokens.sh
kafka_2.13-3.8.0/bin/kafka-delete-records.sh
```

_____

```
hadoop@hadoop-virtualbox:~$ sudo mv kafka_2.13-3.8.0 /opt/kafka
hadoop@hadoop-virtualbox:~$
hadoop@hadoop-virtualbox:~$ ls -lart /opt/kafka
total 80
-rw-r--r-- 1 hadoop hadoop 28359 jul 23 10:04 NOTICE
-rw-r--r-- 1 hadoop hadoop 15295 jul 23 10:04 LICENSE
drwxr-xr-x 2 hadoop hadoop  4096 jul 23 10:09 licenses
drwxr-xr-x 3 hadoop hadoop  4096 jul 23 10:09 config
drwxr-xr-x 3 hadoop hadoop  4096 jul 23 10:09 bin
drwxr-xr-x 7 hadoop hadoop  4096 jul 23 10:09 .
drwxr-xr-x 2 hadoop hadoop  4096 jul 23 10:09 site-docs
drwxr-xr-x 2 hadoop hadoop 12288 ago 31 11:07 libs
drwxr-xr-x 5 root   root    4096 ago 31 11:07 ..
hadoop@hadoop-virtualbox:~$
```

Paso 2. Definimos las variables de entorno para Apache Kafka y lo arrancamos.

```
hadoop@hadoop-virtualbox:~$
hadoop@hadoop-virtualbox:~$ vim .bashrc
hadoop@hadoop-virtualbox:~$
```

```
export KAFKA_HOME=/opt/kafka
export PATH=$KAFKA_HOME/bin:$PATH
```

```
hadoop@hadoop-virtualbox:~$ source .bashrc
hadoop@hadoop-virtualbox:~$
hadoop@hadoop-virtualbox:~$ echo $PATH
/opt/kafka/bin:/opt/nifi/bin:/usr/lib/jvm/java-11-open
sbin:/bin:/usr/games:/usr/local/games:/snap/bin
hadoop@hadoop-virtualbox:~$
```

```
hadoop@hadoop-virtualbox:~$ zookeeper-server-start.sh `echo $KAFKA_HOME`/config/zookeeper.properties
[2024-08-31 11:14:35,837] INFO Reading configuration from: /opt/kafka/config/zookeeper.properties (org.apache.zooke
[2024-08-31 11:14:35,839] INFO clientPortAddress is 0.0.0.0:2181 (org.apache.zookeeper.server.quorum.QuorumPeerConf
[2024-08-31 11:14:35,839] INFO secureClientPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2024-08-31 11:14:35,840] INFO observerMasterPort is not set (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2024-08-31 11:14:35,840] INFO metricsProvider.className is org.apache.zookeeper.metrics.impl.DefaultMetricsProvide
[2024-08-31 11:14:35,841] INFO autopurge.snapRetainCount set to 3 (org.apache.zookeeper.server.DatadirCleanupManage
[2024-08-31 11:14:35,841] INFO autopurge.purgeInterval set to 0 (org.apache.zookeeper.server.DatadirCleanupManager)
[2024-08-31 11:14:35,841] INFO Purge task is not scheduled. (org.apache.zookeeper.server.DatadirCleanupManager)
```

```
zookeeper-server-start.sh `echo $KAFKA_HOME`/config/zookeeper.properties
```

```
hadoop@hadoop-virtualbox:~$ kafka-server-start.sh `echo $KAFKA_HOME`/config/server.properties ←
[2024-08-31 11:16:47,020] INFO Registered kafka:type=kafka.Log4jController MBean (kafka.utils.Log4
[2024-08-31 11:16:47,185] INFO Setting -D jdk.tls.rejectClientInitiatedRenegotiation=true to disab
l)
[2024-08-31 11:16:47,186] INFO RemoteLogManagerConfig values:
        log.local.retention.bytes = -2
        log.local.retention.ms = -2
        remote.fetch.max.wait.ms = 500
        remote.log.index.file.cache.total.size.bytes = 1073741824
        remote.log.manager.copier.thread.pool.size = 10
        remote.log.manager.copy.max.bytes.per.second = 9223372036854775807
```

```
kafka-server-start.sh `echo $KAFKA_HOME`/config/server.properties
```

```
hadoop@hadoop-virtualbox:~$ ls -lart $KAFKA_HOME/config/ ←
total 84
-rw-r--r-- 1 hadoop hadoop 1205 jul 23 10:04 zookeeper.properties
-rw-r--r-- 1 hadoop hadoop 1169 jul 23 10:04 trogdor.conf
-rw-r--r-- 1 hadoop hadoop 1094 jul 23 10:04 tools-log4j.properties
-rw-r--r-- 1 hadoop hadoop 6896 jul 23 10:04 server.properties
-rw-r--r-- 1 hadoop hadoop 2065 jul 23 10:04 producer.properties
-rw-r--r-- 1 hadoop hadoop 4917 jul 23 10:04 log4j.properties
drwxr-xr-x 2 hadoop hadoop 4096 jul 23 10:04 kraft
-rw-r--r-- 1 hadoop hadoop 1221 jul 23 10:04 consumer.properties
-rw-r--r-- 1 hadoop hadoop 2262 jul 23 10:04 connect-standalone.properties
-rw-r--r-- 1 hadoop hadoop 2540 jul 23 10:04 connect-mirror-maker.properties
-rw-r--r-- 1 hadoop hadoop 2063 jul 23 10:04 connect-log4j.properties
-rw-r--r-- 1 hadoop hadoop  881 jul 23 10:04 connect-file-source.properties
-rw-r--r-- 1 hadoop hadoop  883 jul 23 10:04 connect-file-sink.properties
-rw-r--r-- 1 hadoop hadoop 5475 jul 23 10:04 connect-distributed.properties
-rw-r--r-- 1 hadoop hadoop  909 jul 23 10:04 connect-console-source.properties
-rw-r--r-- 1 hadoop hadoop  906 jul 23 10:04 connect-console-sink.properties
drwxr-xr-x 3 hadoop hadoop 4096 jul 23 10:09 .
drwxr-xr-x 8 hadoop hadoop 4096 ago 31 11:14 ..
hadoop@hadoop-virtualbox:~$
```

Paso 3. Creamos un tópico para poder publicar y consumir del mismo.

```
hadoop@hadoop-virtualbox:~$ kafka-topics.sh --create --topic pruebaBDA --bootstrap-server localhost:9092
Created topic pruebaBDA.
hadoop@hadoop-virtualbox:~$
```

```
kafka-topics.sh --create --topic pruebaBDA --bootstrap-server localhost:9092
```

```
hadoop@hadoop-virtualbox:~$ kafka-console-producer.sh --topic pruebaBDA --bootstrap-server localhost:9092
>Linea 1
>Linea 2
>█
```

_____

```
kafka-console-producer.sh --topic pruebaBDA --bootstrap-server localhost:9092
```

```
hadoop@hadoop-virtualbox:~$ kafka-console-consumer.sh --topic pruebaBDA --from-beginning --bootstrap-server localhost:9092
Linea 1
Linea 2
```

```
kafka-console-consumer.sh --topic pruebaBDA --from-beginning --bootstrap-server localhost:9092
```

Paso 4. Añadimos los procesos para publicar en Kafka desde NiFi.

Displaying 1 of 360                                          GenerateFlo

| Type ▲ | Version | Tags |
|---|---|---|
| GenerateFlowFile | 1.27.0 | random, test, load, generate |

**Configure Processor** | GenerateFlowFile 1.27.0

⚠ Invalid

| SETTINGS | SCHEDULING | PROPERTIES | RELATIONSHIPS | COMMENTS |

Scheduling Strategy ❷

Timer driven ⌄

Run Duration ❷

0ms  25ms  50ms  100ms  250ms  500ms  1s  2s

Lower latency                                    Higher throughput

Concurrent Tasks ❷

1

Run Schedule ❷

10 sec

Execution ❷

All nodes ⌄

**Configure Processor** | GenerateFlowFile 1.27.0

⚠ Invalid

| SETTINGS | SCHEDULING | PROPERTIES | RELATIONSHIPS | COMMENTS |

Required field                                                    ⊙  +

| Property | | Value | |
|---|---|---|---|
| File Size | ❷ | 0B | |
| Batch Size | ❷ | 1 | |
| Data Format | ❷ | Text | |
| Unique FlowFiles | ❷ | false | |
| Custom Text | ❷ | current time: ${now():format("yyyy-MM-dd HH:mm:ss")} | |
| Character Set | ❷ | UTF-8 | |
| Mime Type | ❷ | No value set | |

_____

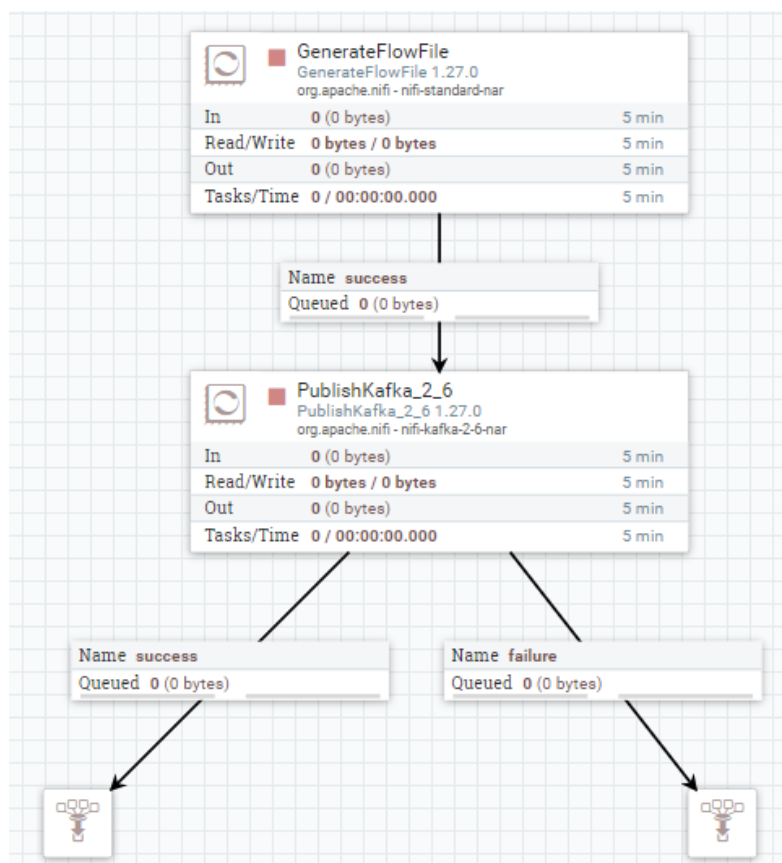| Displaying 3 of 360 | | PublishKafka_ |
|---|---|---|
| **Type ▲** | **Version** | **Tags** |
| PublishKafka_1_0 | 1.27.0 | PubSub, 1.0, Message, Kafka, A... |
| PublishKafka_2_0 | 1.27.0 | PubSub, Message, 2.0, Kafka, A... |
| PublishKafka_2_6 | 1.27.0 | PubSub, Message, Kafka, 2.6, A... |

**Configure Processor** | PublishKafka_2_6 1.27.0

⚠ Invalid

| SETTINGS | SCHEDULING | PROPERTIES | RELATIONSHIPS | COMMENTS |
|---|---|---|---|---|

**Required field**                                                        ⊘  +

| Property | | Value | |
|---|---|---|---|
| Kafka Brokers | ❓ | localhost:9092 | |
| Topic Name | ❓ | pruebaBDA | |
| Use Transactions | ❓ | true | |
| Transactional Id Prefix | ❓ | No value set | |
| Message Demarcator | ❓ | No value set | |
| Failure Strategy | ❓ | Route to Failure | |

Paso 5. Comprobamos el flujo.



success

Displaying 5 of 5 (165.00 bytes)                                                                           The source of this queue is currently running. This listing may no longer be accurate.

| | Position | UUID | Filename | File Size | Queued Duration | Lineage Duration | Penalized | |
|---|---|---|---|---|---|---|---|---|
| ⓘ | 1 | 10543252-19d7-41c3-aadb-1c9551c71067 | 10543252-19d7-41c3-aadb-1c9551c71067 | 33.00 bytes | 00:00:45.135 | 00:00:45.281 | No | ⬇ 👁 ❡ |
| ⓘ | 2 | bb469923-c28c-42e2-9a9f-6e5bec6cb6eb | bb469923-c28c-42e2-9a9f-6e5bec6cb6eb | 33.00 bytes | 00:00:36.249 | 00:00:36.265 | No | ⬇ 👁 ❡ |
| ⓘ | 3 | 21cc3c48-2d5b-47c3-9e37-c8841902fb30 | 21cc3c48-2d5b-47c3-9e37-c8841902fb30 | 33.00 bytes | 00:00:26.247 | 00:00:26.262 | No | ⬇ 👁 ❡ |
| ⓘ | 4 | b221e537-be64-4dc9-98a6-ec4cc1d08549 | b221e537-be64-4dc9-98a6-ec4cc1d08549 | 33.00 bytes | 00:00:16.250 | 00:00:16.261 | No | ⬇ 👁 ❡ |
| ⓘ | 5 | b0ebee57-80bb-4cb0-ac36-e5f11ed93f92 | b0ebee57-80bb-4cb0-ac36-e5f11ed93f92 | 33.00 bytes | 00:00:06.235 | 00:00:06.257 | No | ⬇ 👁 ❡ |



View as:  original

```
1   current time: 2024-08-31 11:51:31
```

_____

Paso 6. Añadimos los procesos para consumir en Kafka desde NiFi.

| Displaying 3 of 360 | | ConsumeKafka_ |
|---|---|---|

| Type ▲ | Version | Tags |
|---|---|---|
| ConsumeKafka_1_0 | 1.27.0 | PubSub, Consume, 1.0, Ingest, … |
| ConsumeKafka_2_0 | 1.27.0 | PubSub, Consume, Ingest, 2.0, … |
| ConsumeKafka_2_6 | 1.27.0 | PubSub, Consume, Ingest, Get, … |

**Configure Processor** | ConsumeKafka_2_6 1.27.0

⚠ Invalid

| SETTINGS | SCHEDULING | PROPERTIES | RELATIONSHIPS | COMMENTS |
|---|---|---|---|---|

Required field

| Property | | Value | |
|---|---|---|---|
| Kafka Brokers | ❓ | localhost:9092 | |
| Topic Name(s) | ❓ | pruebaBDA | |
| Topic Name Format | ❓ | names | |
| Group ID | ❓ | grupo1 | |

| Displaying 4 of 360 | | Merge |
|---|---|---|

| Type ▲ | Version | Tags |
|---|---|---|
| JoinEnrichment | 1.27.0 | enrichment, fork, record, merge,… |
| MergeContent | 1.27.0 | zip, flowfile-stream-v3, flowfile-… |
| MergeRecord | 1.27.0 | correlation, stream, merge, reco… |
| UnpackContent | 1.27.0 | zip, flowfile-stream-v3, flowfile-… |

**Configure Processor** | MergeContent 1.27.0

⚠ Invalid

| SETTINGS | SCHEDULING | PROPERTIES | RELATIONSHIPS | COMMENTS |
|---|---|---|---|---|

Required field

| Property | | Value |
|---|---|---|
| Merge Strategy | ❓ | Bin-Packing Algorithm |
| Merge Format | ❓ | Binary Concatenation |
| Attribute Strategy | ❓ | Keep Only Common Attributes |
| Correlation Attribute Name | ❓ | No value set |
| Minimum Number of Entries | ❓ | 50 |
| Maximum Number of Entries | ❓ | 100 |

_____

**Configure Processor** | UpdateAttribute 1.27.0

⚠ Invalid

| SETTINGS | SCHEDULING | PROPERTIES | RELATIONSHIPS | COMMENTS |

Required field

| Property | | Value | |
|---|---|---|---|
| Delete Attributes Expression | ❓ | No value set | |
| Store State | ❓ | Do not store state | |
| Stateful Variables Initial Value | ❓ | No value set | |
| Cache Value Lookup Cache Size | ❓ | 100 | |
| filename | ❓ | Kafka_${now():format("yyyyMMdd_HHmmss")}.txt | 🗑 |

**Configure Processor** | PutFile 1.27.0

■ Stopped

| SETTINGS | SCHEDULING | PROPERTIES | RELATIONSHIPS | COMMENTS |

Required field

| Property | | Value | |
|---|---|---|---|
| Directory | ❓ | /home/hadoop/workspace/salidaNifi | |
| Conflict Resolution Strategy | ❓ | replace | |

Paso 7. Comprobamos el flujo.