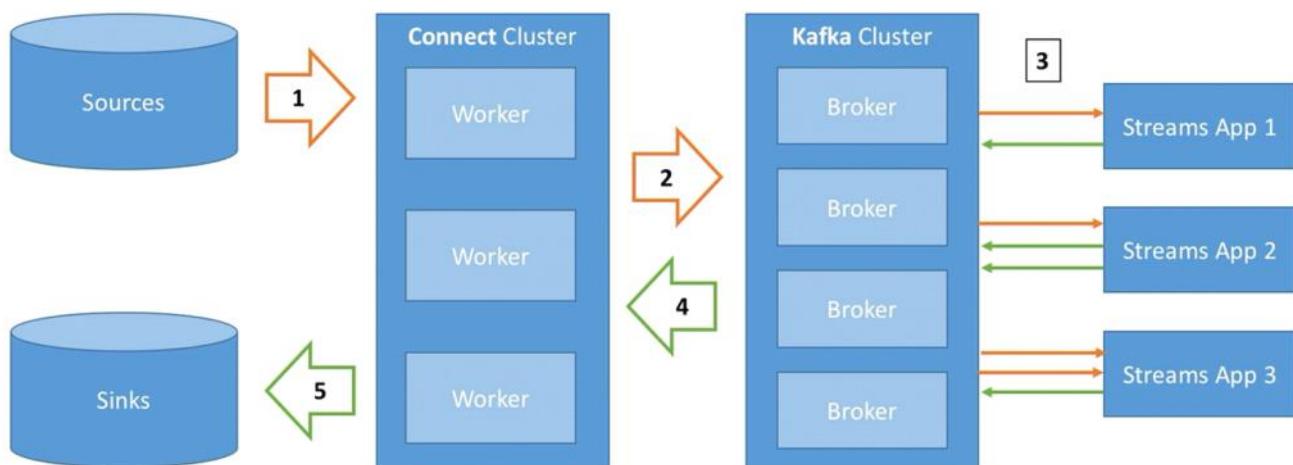


<https://github.com/PacktPublishing/Apache-Kafka-Series---Kafka-Connect-Hands-on-Learning>

Source => Kafka	Producer API	Kafka Connect Source
Kafka => Kafka	Consumer, Producer API	Kafka Streams
Kafka => Sink	Consumer API	Kafka Connect Sink
Kafka => App	Consumer API	

Kafka Connect and Streams Architecture Design



Kafka Connect - Concepts

- Kafka Connect Cluster has multiple loaded **Connectors**
 - Each connector is a re-usable piece of code (java jars)
 - Many connectors exist in the open source world, leverage them!
- Connectors + **User Configuration** => **Tasks**
 - A task is linked to a connector configuration
 - A job configuration may contain multiple tasks

Kafka Connect - Concepts

- Kafka Connect Cluster has multiple loaded **Connectors**
 - Each connector is a re-usable piece of code (java jars)
 - Many connectors exist in the open source world, leverage them!
- Connectors + **User Configuration** => **Tasks**
 - A task is linked to a connector configuration
 - A job configuration may spawn multiple tasks
- Tasks are executed by Kafka Connect **Workers** (servers)
 - A worker is a single java process
 - A worker can be standalone or in a cluster

Standalone vs Distributed Mode

- Standalone:
 - A single process runs your connectors and tasks
 - Configuration is bundled with your process
 - Very easy to get started with, [useful for development and testing](#)
 - Not fault tolerant, no scalability, hard to monitor
- Distributed:
 - Multiple workers run your connectors and tasks
 - Configuration is submitted using a REST API
 - Easy to scale, and fault tolerant (rebalancing in case a worker dies)
 - Useful for [production deployment of connectors](#)

List of available connectors

- There are many source and sinks connectors
- The best place to find them is the Confluent website:

<https://www.confluent.io/product/connectors/>



Secure https://www.confluent.io/product/connectors/

The screenshot shows the Confluent Kafka Connectors page. At the top, there's a navigation bar with links to Product, Solutions, Services, Resources, About, Blog, Download, and a search icon. Below the navigation, a banner states: "purpose is to make it easy to add new systems to your scalable and secure stream data pipelines." A text block explains that users instantiate Kafka Connectors for systems they want to pull data from or push data to. It distinguishes between Source Connectors (import data) and Sink Connectors (export data). A note says this page lists many notable connectors.

CONNECTOR	TAGS	DEVELOPER/SUPPORT	DOWNLOAD
Amazon S3 (Sink)	S3	Confluent	Confluent
Elasticsearch (Sink)	search, Elastic, log, analytics	Confluent	Confluent
HDFS (Sink)	HDFS, Hadoop, Hive	Confluent	Confluent
JDBC (Source)	JDBC, MySQL	Confluent	Confluent
JDBC (Sink)	JDBC, MySQL	Confluent	Confluent
Attunity (Source)	CDC	Attunity	Attunity
Azure IoT Hub (Source)	IoT, messaging	Microsoft	Azure

Want to build a connector?

Interested Open Source Developers and Vendors can get started with the following resources:

- > Kafka Connect Overview
- > Kafka Connect Developers Guide
- > Partner Development Guide for Kafka Connect

Already built a connector?

If you have a connector that you'd like to see added to our list, let us know at: confluent-platform@googlegroups.com

If your connector meets our criteria, we'll post it and send you a free t-shirt. Woohoo!

Contact Us

Hello. Is there anything I can help you with?

For software vendors seeking to build a

Внутренне коннектор/worker это обычный консьюмер/продьюсер

- Internally, Kafka Connect is just a Kafka client using the standard Producer
- Connect is a prebuilt framework which provides the basic client architecture (using the standard Kafka Client API). This framework is written to make Connect highly available and scalable. However, it relies on plugins called Connectors to provide the specific functionality to communicate with your chosen endpoints.
- Connect is restricted to data copying, with light data transformations
- Connect Converters are like wrappers around Serializers and Deserializers. Since a Worker can be both a Producer and a Consumer (depending on the Connector), it was decided that having one object to specify rather than two a more efficient way to configure the Worker.

- **Connectors** are logical jobs responsible for managing the copying of data between Kafka and another system
- Connector **Sources** read data from an external data system into Kafka
 - Internally, source connector is a Kafka **Producer**
- Connector **Sinks** write Kafka data to an external data system
 - Internally, sink connector is a Kafka **Consumer Group**

An instance of Connect (called a **worker**) can function as both a producer and consumer, depending on the type of connector you install. In fact, internally connectors are running the standard client code: Producer (source connectors) and Consumer (sink connectors).

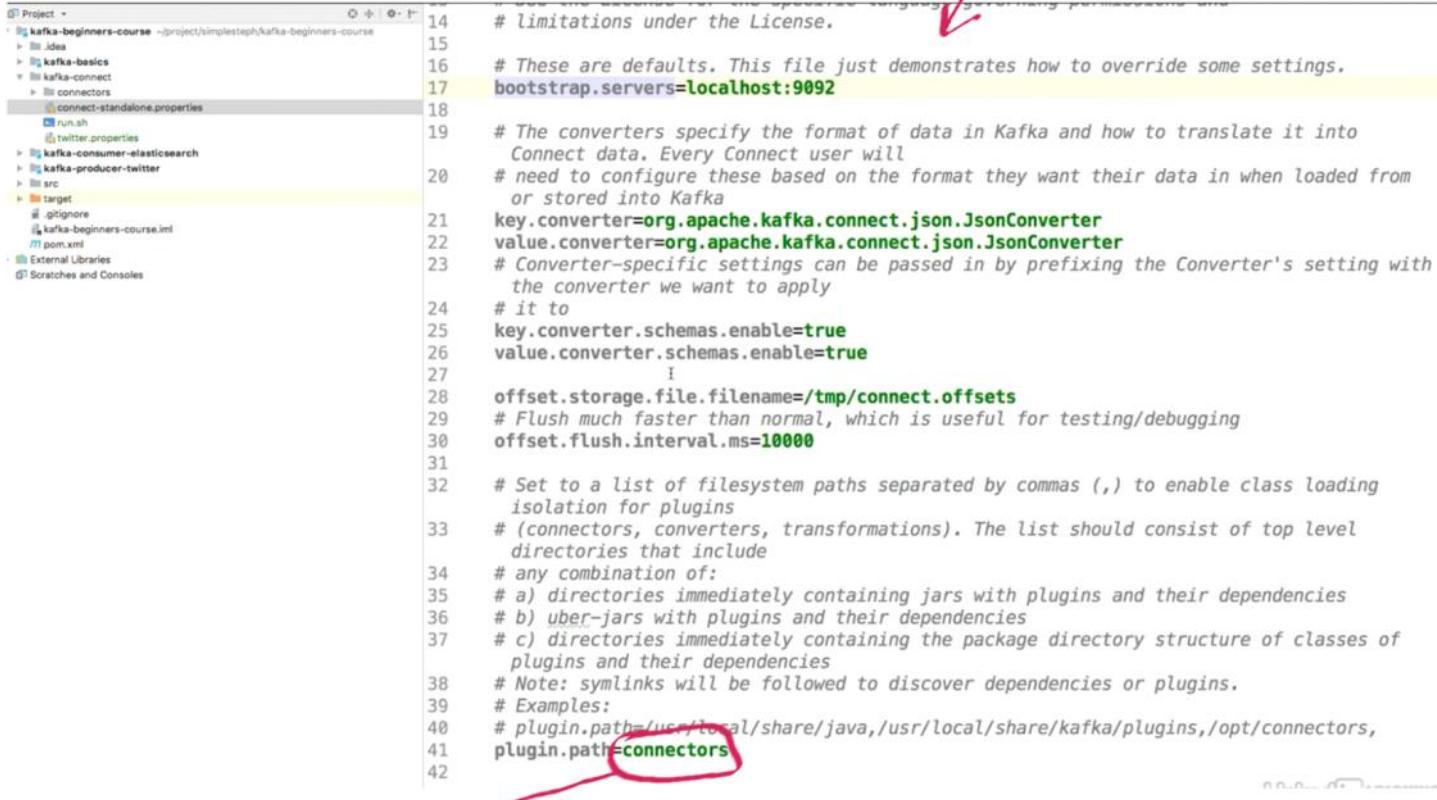
Source connectors read from external data stores (e.g., databases). Sink Connectors pull data from a Kafka topic and write it out to an external application (e.g., HDFS, Elasticsearch).

The term "Connector" can be used for either the type of plugin (e.g., the HDFS Connector) or a configured instance of a connector (e.g., the mycompany_HDFS connector).

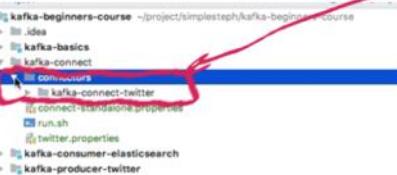
CLI

30 ноября 2020 г. 19:45

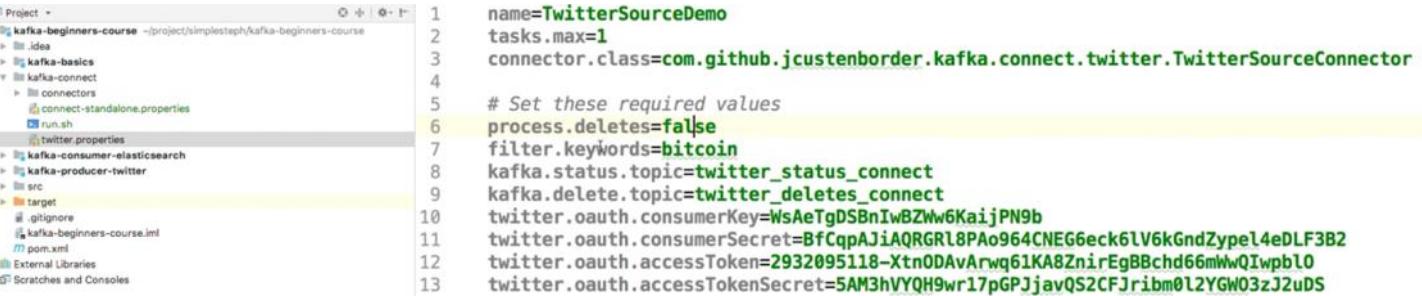
```
~/kafka_2.12-2.0.0 ➤ connect-standalone  
USAGE: /usr/local/Cellar/kafka/2.0.0/libexec/bin/connect-standalone.sh [-daemon] [connect-standalone.properties]
```



```
Project: kafka-beginners-course (~) [File | Open | New | Help]  
  -> kafka-connect  
    -> connectors  
      -> connect-standalone.properties  
        -> run.sh  
        -> twitter.properties  
      -> kafka-connect-twitter  
      -> connect-standalone.properties  
        -> run.sh  
        -> twitter.properties  
      -> kafka-consumer-elasticsearch  
      -> kafka-producer-twitter  
      -> src  
    -> target  
      -> gitignore  
      -> kafka-beginners-course.iml  
    -> pom.xml  
  -> External Libraries  
  -> Scratches and Consoles  
  
14 # limitations under the License.  
15  
16 # These are defaults. This file just demonstrates how to override some settings.  
17 bootstrap.servers=localhost:9092  
18  
19 # The converters specify the format of data in Kafka and how to translate it into  
# Connect data. Every Connect user will  
20 # need to configure these based on the format they want their data in when loaded from  
# or stored into Kafka  
21 key.converter=org.apache.kafka.connect.json.JsonConverter  
22 value.converter=org.apache.kafka.connect.json.JsonConverter  
23 # Converter-specific settings can be passed in by prefixing the Converter's setting with  
# the converter we want to apply  
24 # it to  
25 key.converter.schemas.enable=true  
26 value.converter.schemas.enable=true  
27  
28 offset.storage.file.filename=/tmp/connect.offsets  
29 # Flush much faster than normal, which is useful for testing/debugging  
30 offset.flush.interval.ms=10000  
31  
32 # Set to a list of filesystem paths separated by commas (,) to enable class loading  
# isolation for plugins  
33 # (connectors, converters, transformations). The list should consist of top level  
# directories that include  
34 # any combination of:  
35 # a) directories immediately containing jars with plugins and their dependencies  
36 # b) uber-jars with plugins and their dependencies  
37 # c) directories immediately containing the package directory structure of classes of  
# plugins and their dependencies  
38 # Note: symlinks will be followed to discover dependencies or plugins.  
39 # Examples:  
40 # plugin.path=/usr/local/share/java,/usr/local/share/kafka/plugins,/opt/connectors,  
41 plugin.path=connectors  
42
```



connect-standalone.bat connect-standalone.properties twitter.properties



```
Project: kafka-beginners-course (~) [File | Open | New | Help]  
  -> kafka-connect  
    -> connectors  
      -> connect-standalone.properties  
        -> run.sh  
        -> twitter.properties  
      -> kafka-connect-twitter  
      -> connect-standalone.properties  
        -> run.sh  
        -> twitter.properties  
      -> kafka-consumer-elasticsearch  
      -> kafka-producer-twitter  
      -> src  
    -> target  
      -> gitignore  
      -> kafka-beginners-course.iml  
    -> pom.xml  
  -> External Libraries  
  -> Scratches and Consoles  
  
1 name=TwitterSourceDemo  
2 tasks.max=1  
3 connector.class=com.github.jcugenborder.kafka.connect.twitter.TwitterSourceConnector  
4  
5 # Set these required values  
6 process.delete=false  
7 filter.keywords=bitcoin  
8 kafka.status.topic=twitter_status_connect  
9 kafka.delete.topic=twitter_deletes_connect  
10 twitter.oauth.consumerKey=WsAeTgDSBnIwBZw6KaijPN9b  
11 twitter.oauth.consumerSecret=BfCqpAJiAQGRl8PAo964CNEG6eck6lV6kGndZypel4eDLF3B2  
12 twitter.oauth.accessToken=2932095118-XtnODAvArwg61KA8ZnirEgBBchd66mWQIwpb10  
13 twitter.oauth.accessTokenSecret=5AM3hVYQH9wr17pGPJjavQS2CFJribm0l2YGW03zJ2uDS
```

пример

2. Create a table operators in PostgreSQL:

```
postgres=# CREATE TABLE operators(  
    id int not null primary key,  
    name varchar(50) not null  
)
```

```
$ docker container run --rm -d \
-p 8080:80 \
--net using=connect_confluent \
-e PGADMIN_DEFAULT_EMAIL=student@confluent.io \
-e PGADMIN_DEFAULT_PASSWORD=TopSecret \
dpage/pgadmin4
```

1. Create a new Topic called pg-operators with one partition and one replica:

```
$ kafka-topics \
--bootstrap-server kafka:9092 \
--create \
--topic pg-operators \
--partitions 1 \
--replication-factor 1
```

2. Add a JDBC source connector via command line and REST API of Connect:

```
$ curl -s -X POST \
-H "Content-Type: application/json" \
--data '{
  "name": "Operators-Connector",
  "config": {
    "connector.class": "io.confluent.connect.jdbc.JdbcSourceConnector",
    "connection.url": "jdbc:postgresql://postgres:5432/postgres",
    "connection.user": "postgres",
    "table.whitelist": "operators",
    "mode": "incrementing",
    "incrementing.column.name": "id",
    "table.types": "TABLE",
    "topic.prefix": "pg-",
    "numeric.mapping": "best_fit",
    "transforms": "createKey,extractInt",
    "transforms.createKey.type": "org.apache.kafka.connect.transforms.ValueToKey",
    "transforms.createKey.fields": "id",
    "transforms.extractInt.type": "org.apache.kafka.connect.transforms.ExtractField$Key",
    "transforms.extractInt.field": "id"
  }
}' http://connect:8083/connectors | jq
```

```
$ kafka-console-consumer \
--bootstrap-server kafka:9092 \
--property schema.registry.url=http://schema-registry:8081 \
--topic pg-operators \
--from-beginning \
--property print.key=true
```

Viewing Kafka Connect Logs

- In standalone mode
 - They will directly be in your terminal
- In distributed mode
 - They will be at the URL <http://127.0.0.1:3030/logs/>
 - Take the file 'connect-distributed.log'
- Look for "ERROR" in the log, that's a good indicator of what went wrong (an error message will come with the error)

The screenshot shows a web browser window for 'Landoop Kafka Development' at <https://fast-data-dev.demo.landoop.com>. The main page displays various service status icons and a summary of tests run. A red arrow points from the 'BROWSE' link in the 'running services log files' section to the 'logs' page.

Running Services Log Files:

- running services log files

Logs Page:

Logs

Name	Size	Modified
broker.log	132 kB	30/04/2017, 19:38:21
caddy.log	33 kB	30/04/2017, 19:46:00
connect-distributed.log	109 kB	30/04/2017, 19:46:00
log-to-kafka.log	838 B	30/04/2017, 16:09:00
put.log	557 B	30/04/2017, 19:46:00
rest-proxy.log	88 kB	30/04/2017, 19:46:00
schema-registry.log	21 kB	30/04/2017, 16:11:00
smoke-tests.log	44 kB	30/04/2017, 16:11:00
supervisor.log	2.7 kB	30/04/2017, 16:11:00
zookeeper.log	70 kB	30/04/2017, 16:11:00

Served with Caddy

ETL with KAFKA

3 марта 2021 г. 12:53

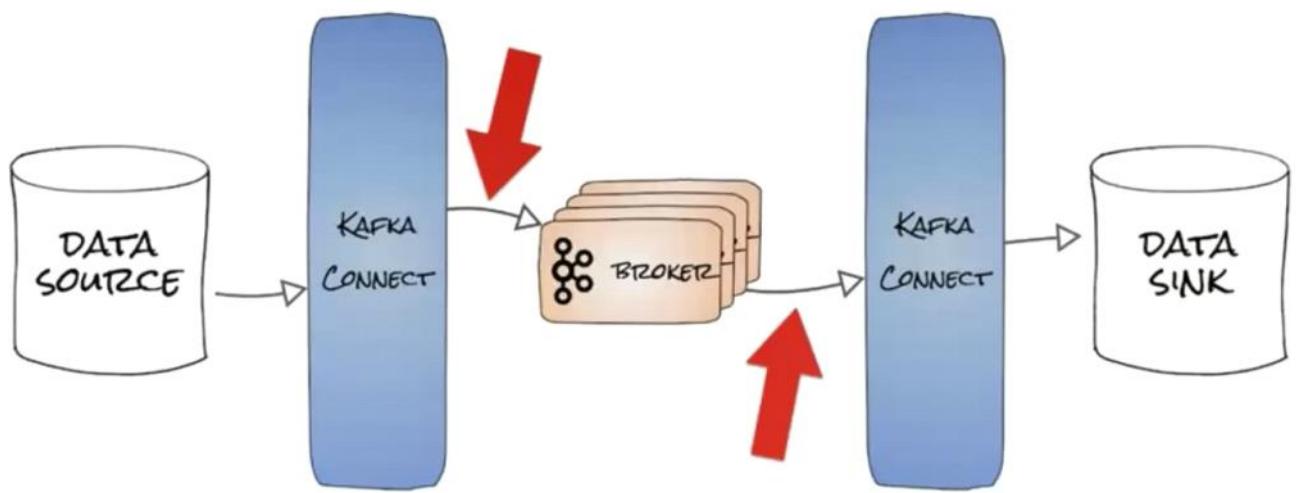
♥ ETL with KAFKA

раньше тулза для этого называлась **copycat**
<https://blog.codecentric.de/en/2018/03/etl-kafka/>



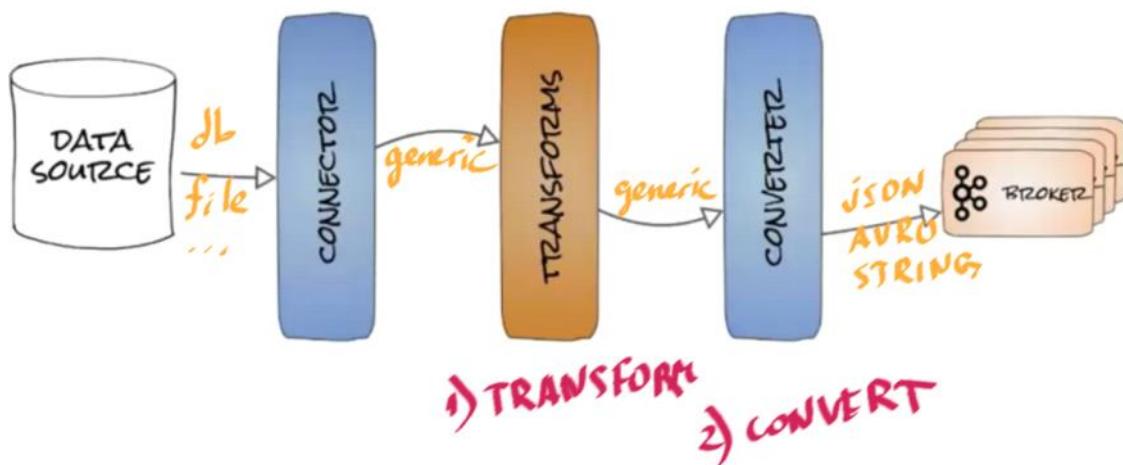
→ Neha Narkhede's talk, Thu 1.15pm

Transforming Data with Connect



- 1) connector преобразует формат источника в общий формат generic format
- 2) generic->generic трансформация
- 3) generic->(avro,JSON,String) преобразует в итоговый формат

Where SMTs Live



Single message transforms are a new stage in the Kafka Connect pipeline. Here we can see where they fit in for a source connector.

Source connectors read data from another system and store the data in Kafka. The connector (or more accurately, the connector's tasks) read data from the data source. At this stage the data will be in whatever format the source system exposes, but the connector will translate that into Kafka Connect's generic data API. Once it is in this generic format, the Kafka Connect single message transformations can be applied. Because they use the generic data format, any transformation can be used with any connector, independent of the original format of the data. The transformation is applied and the resulting data is still in a generic data API format. Finally, the data is passed to a converter, which serializes it into a format like Avro or JSON before it is passed to Kafka for storage. The pipeline for sink connectors looks similar, just reversed.

As you can see, by breaking the processing down into separate steps and leveraging Connect's data API, we can make transformations very general and applicable regardless of which connector you're using, much like we make serialization formats reusable across all connectors.

SMT не годится для серьезных трансформаций, те только для простых преобразований (добавить поле, убрать поле, писать константу в поле для всех сообщений)

- для серьезных трансформаций лучше использовать kafka streams или KSQL

Single Message Transform (SMT)

Modify events before storing in Kafka:

- Mask sensitive information
- Add identifiers
- Tag events
- Lineage/provenance
- Remove unnecessary columns

Modify events going out of Kafka:

- Route high priority events to faster data stores
- Direct events to different Elasticsearch indices
- Cast data types to match destination
- Remove unnecessary columns

Single Message Transform - Details

InsertField	insert a field using attributes from the message metadata or from a configured static value
ReplaceField	rename fields, or apply a blacklist or whitelist to filter
MaskField	replace field with valid <code>null</code> value for the type (0, empty string, etc)
ValueToKey	replace the key with a new key formed from a subset of fields in the value payload
HoistField	wrap the entire event as a single field inside a <code>Struct</code> or a <code>Map</code>
ExtractField	extract a specific field from <code>Struct</code> and <code>Map</code> and include only this field in results
SetSchemaMetadata	modify the schema name or version
TimestampRouter	modify the Topic of a record based on the original Topic name and timestamp
RegexRouter	modify the topic of a record based on original topic, replacement string and a regular expression

Many customers need a way to perform simple filtering and transformations for which building a streams application is too much work. For example, a company might want to take order information and send that to a team for analysis but needed to redact sensitive information such as social security numbers or credit card numbers.

The above properties can be configured at the **connector level** and be applied to either **key** and/or **value**.

TimestampRouter - Modify the topic of a record based on original topic and timestamp. Useful when using a sink that needs to write to different tables or indexes based on timestamps

More info: http://kafka.apache.org/documentation.html#connect_transforms

SMT конфигурируется в настройках коннектора

```
name=local-file-source
connector.class=FileStreamSource
tasks.max=1
file=test.txt
topic=connect-test
transforms=MakeMap,InsertSource, InsertKey, ExtractStoreId, MessageTypeRouter
transforms.MakeMap.type=org.apache.kafka.connect.transforms.HoistField$Value
transforms.MakeMap.field=line
transforms.InsertSource.type=org.apache.kafka.connect.transforms.InsertFields$Value
transforms.InsertSource.static.field=data_source
transforms.InsertSource.static.value=test-file-source
transforms.InsertKey.type=org.apache.kafka.connect.transforms.ValueToKey
transforms.InsertKey.fields=storeId
transforms.ExtractStoreId.type=org.apache.kafka.connect.transforms.ExtractField$Key
transforms.ExtractStoreId.field=storeId
transforms.MessageTypeRouter.type=org.apache.kafka.connect.transforms.RegexRouter
transforms.MessageTypeRouter.regex=(foo|bar|baz)-.*
transforms.MessageTypeRouter.replacement=$1-logs
```

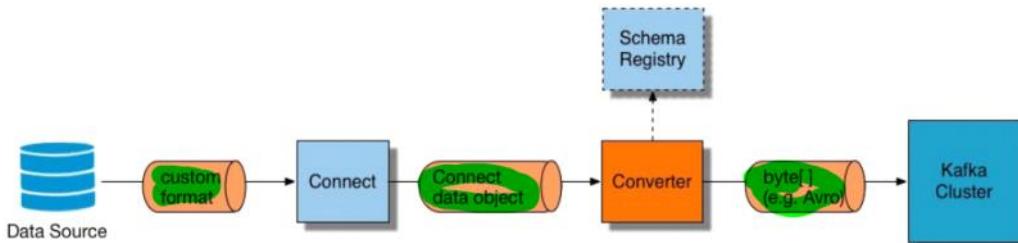
This starts to look kind of unpleasant!

1. Hoist Field: message sent `Foo` → transformed value `{"line": "Foo"}`
2. Insert Source: Add field `data_source` with value `test-file-source` to the imported records
3. Insert Key & ExtractStoreId: Take the value of the field `storeId` and use this as value for the Kafka record `key`
4. MessageTypeRouter: Update the record's topic using the configured regular expression and replacement string

- converter преобразует generic format уже в конкретный формат, который пойдет в кафку (те делает сериализацию и десериализацию)
- в конверторах можно использовать а можно и не использовать схему данных (коннектор интегрируется с реестром схем)

Converting Data

- Converters provide the data format written to or read from Kafka (like Serializers)
- Converters are decoupled from Connectors to enable reuse of converters
 - Allows any connector to work with any serialization format
- Example of format conversion in a source converter (sink converter is the reverse)



Serialization formats may seem like a minor detail, but **not** separating the details of data serialization in Kafka from the details of source or sink systems would result in a lot of inefficiency:

- First, a lot of code for doing simple data conversions are duplicated across a large number of ad hoc connector implementations.
- Second, each connector ultimately contains its own set of serialization options as it is used in more environments—JSON, Avro, Thrift, ProtoBuf, and more.

Much like the serializers in Kafka's producer and consumer, the **Converters** abstract away the details of serialization. Converters are different because they guarantee data is transformed to a common data API defined by Kafka Connect. This API supports both schema and schema-less data, common primitive data types, complex types like structs, and logical type extensions. By sharing this API, connectors write one set of translation code and Converters handle format-specific details. For example, the JDBC connector can easily be used to produce either JSON or Avro to Kafka, without any format-specific code in the connector.

Converter Data Formats

- Converters apply to both the key and value of the message
 - Key and value converters can be set independently
 - `key.converter`
 - `value.converter`
- Pre-defined data formats for Converter
 - Avro: `AvroConverter`
 - JSON: `JsonConverter`
 - String: `StringConverter`
 - Byte Array: `ByteArrayConverter`

The `ByteArrayConverter` (detailed in KIP-128, introduced in Kafka 0.11) provides an option for dealing with raw data, in contrast to structured data. It saves cost (CPU, memory, garbage collection overhead) of deserializing/converting data to Connect format.

For example, this is particularly useful with **Replicator**. This connector doesn't need to convert to/from byte array format and objects - it just copies the raw data as bytes for better performance.

как добавить в конфиге ссылку на реестр схем

Avro Converter as a Best Practice

Use an Avro Converter and Schema Registry!

```
key.converter=io.confluent.connect.avro.AvroConverter
key.converter.schema.registry.url=http://schemaregistry1:8081
value.converter=io.confluent.connect.avro.AvroConverter
value.converter.schema.registry.url=http://schemaregistry1:8081
```

пример

В этом примере мы воспользуемся двумя встроенными преобразованиями: TimestampConverter и ReplaceField

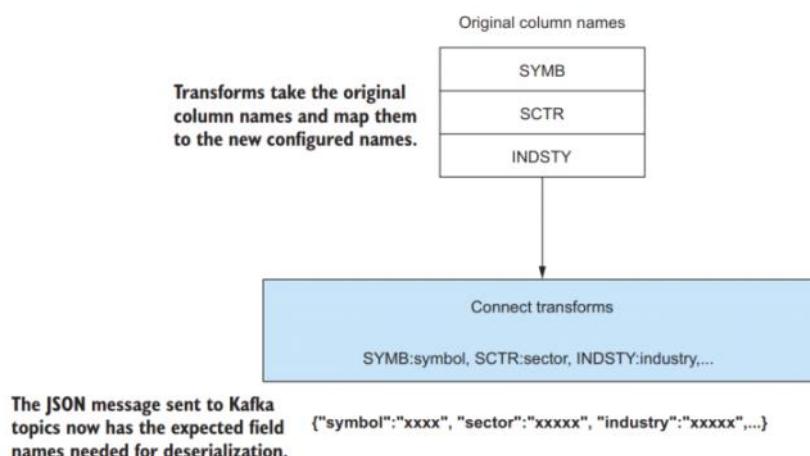


Figure 9.3 Transforming the names of the columns to match expected field names

Listing 9.1 JDBC connector properties

```
Date field to convert
transforms=ConvertDate,Rename
transforms.ConvertDate.type=
  ↳ org.apache.kafka.connect.transforms.TimestampConverter$Value
  ↳ transforms.ConvertDate.field=TXNTS
  ↳ transforms.ConvertDate.target.type=string
  ↳ transforms.ConvertDate.format=yyyy-MM-dd'T'HH:mm:ss.SSS-0400
  ↳ transforms.Rename.type=
  ↳ org.apache.kafka.connect.transforms.ReplaceField$Value
  ↳ transforms.Rename.renames=SMBL:symbol, SCTR:sector, ...
Type for the Rename alias
Aliases for the transformers
Type for the ConvertDate alias
List of column names (truncated for clarity) to replace. The format is Original:Replacement.
Output type of the converted date field
Format of the date
```

преобразования должны быть простыми.

- хотя может показаться заманчивым применить преобразования повсеместно, импортируя данные в Kafka через Kafka Connect, помните: преобразования должны быть простыми. Для любых преобразований, кроме самых простых, показанных в наших примерах, лучше извлечь данные в Kafka и воспользоваться Kafka Streams для вычислительно сложных преобразований.

standalone mode (dev)

4 декабря 2020 г. 18:31

Two Modes: Standalone and Distributed

Standalone mode

- 1 worker on 1 machine
- When to use:
 - testing and development
 - non distributable process

Distributed mode

- n workers on m machines
- When to use:
 - production environments
 - for fault tolerance
 - and scalability

Kafka Connect can be run in two modes

- **Standalone**

- 1 worker process on single machine
 - Use cases:
 - testing and development, or
 - when a process should not be distributed (e.g. `tail` a log file)

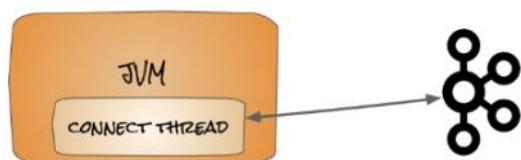
- **Distributed mode**

- Multiple worker processes on one or more machines
 - Use Case: requirements for fault tolerance and scalability

Running in Standalone Mode

To run in standalone mode, start a process by providing arguments

- Standalone config properties file
 - 1...n connector config files
- i** Each connector instance runs in own thread



```
$ connect-standalone connect-standalone.properties \
connector1.properties [connector2.properties connector3.properties ...]
```

A single worker can run as many connectors as you wish. Once the connectors are running, changes to their configuration files will only take effect on restart of the worker. Making changes to the connectors without a restart will be discussed later in the chapter.

Building and running a connector in standalone mode



- Building: `mvn clean package`
- Running (see `./run.sh` – only for linux / mac users):
 1. Export classpath:
 - `export CLASSPATH=$(find target/ -type f -name '*.jar'| grep '\-package' | tr '\n' ':')`
 2. Build the Dockerfile:
 - `docker build . -t simplesteph/kafka-connect-source-github:1.0`
 3. Run the docker image:
 - `docker run -e CLASSPATH=$CLASSPATH \
--net=host --rm -t \
-v $(pwd)/offsets:/kafka-connect-source-github/offsets \
simplesteph/kafka-connect-source-github:1.0`

```
connect-standalone config/worker.properties config/GitHubSourceConnectorExample.properties
```

Configuring Workers: Standalone Mode

Parameter	Description
<code>offset.storage.file.filename</code>	The filename in which to store offset data for the Connectors (Default: ""). This enables a standalone process to be stopped and then resume where it left off.

The provided file `/etc/kafka/connect-standalone.properties` sets this parameter to `/tmp/connect.offsets` which is not a very durable location on disk

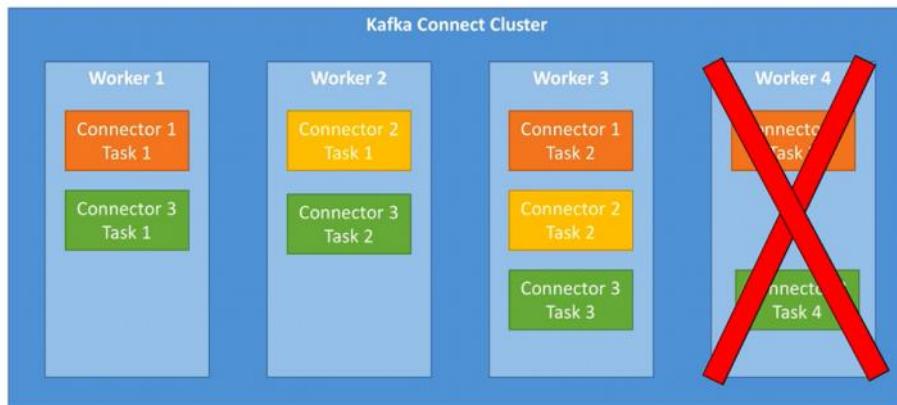
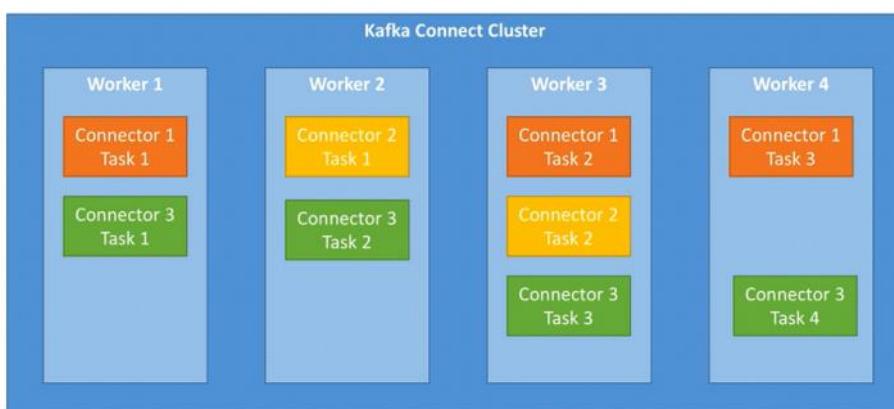
... distributed mode (prod)

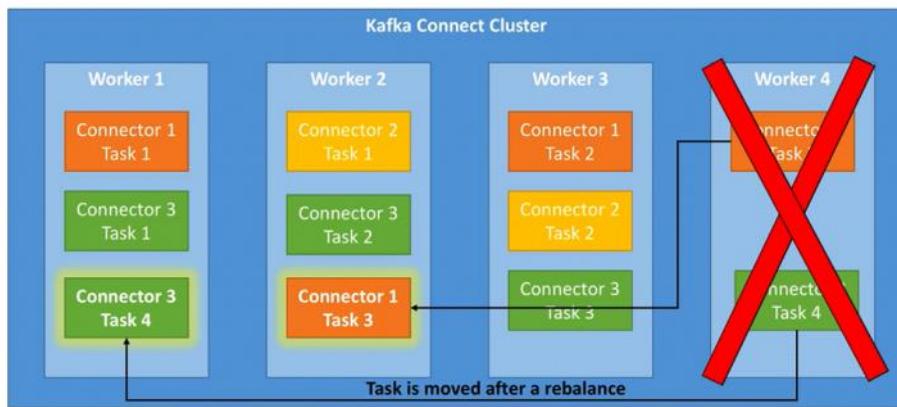
4 декабря 2020 г. 18:31

tasks & workers

- задачи автоматически распределяются между воркерами (и перераспределяются в случае ошибки)
- задачи относящиеся к коннектору не обязательно должны запускаться вместе на одном и том же сервере
- As we can see, Kafka Connect uses a similar technique that we know from consumer groups to reschedule or reallocate workload evenly across the available workers. In our case the connector job B has been moved to worker 2 whilst the connector job C and task C1 have been moved to worker 3.
- Rebalancing/task assignment is performed by the Worker, not by the Connector instance thread. Much like consumer groups, there's a concept of a "leader" Worker which performs task reassessments whenever rebalancing is triggered, which can involve movement of connector threads in addition to tasks. See <https://github.com/apache/kafka/blob/2.2/connect/runtime/src/main/java/org/apache/kafka/connect/runtime/distributed/WorkerCoordinator.java#L220-L258> for more detail

Distributed Architecture in details





В отличие от самописных клиентов, стандартный коннектор обеспечивает load balancing + fault tolerance

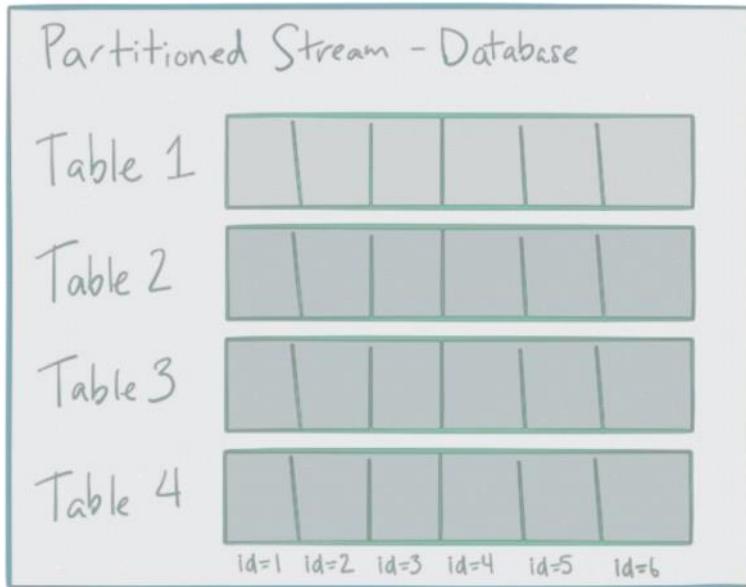
Kafka Connect Reliability and Scalability

- Elasticity and scalability
 - To add resources to Connect simply start more workers
 - Connect Workers discover each other and load-balance the work automatically
 - Allows running more connectors and more tasks
- Reliability and fault-tolerance
 - Failed connector will automatically restart
 - Tasks and connectors on failed worker will automatically (re)start on another worker
- ALL connectors support at-least-once semantics
- **SOME** connectors support exactly-once semantics
- MANY connectors support schema evolution
 - Data format changes in source system will reflect in Kafka and all target systems



To clarify the point "**Failed Connector will automatically restart**": If a worker fails, any connectors running on that worker will automatically be restarted. If a connector thread or task fails (while the worker continues to run), **fault tolerance is not automatic!**

`connector=DB` мы параллельно копируем каждую таблицу в отдельный топик
`connector=FILE` мы параллельно копируем каждый файл в отдельный топик



Kafka Connect also uses the concept of partitioning the input stream of data. Admittedly it is a bit **confusing** to use the same terminology for two different things...

In case of a relational DB as source, the **partitions** of the input stream of data would correspond to the **tables** of the DB

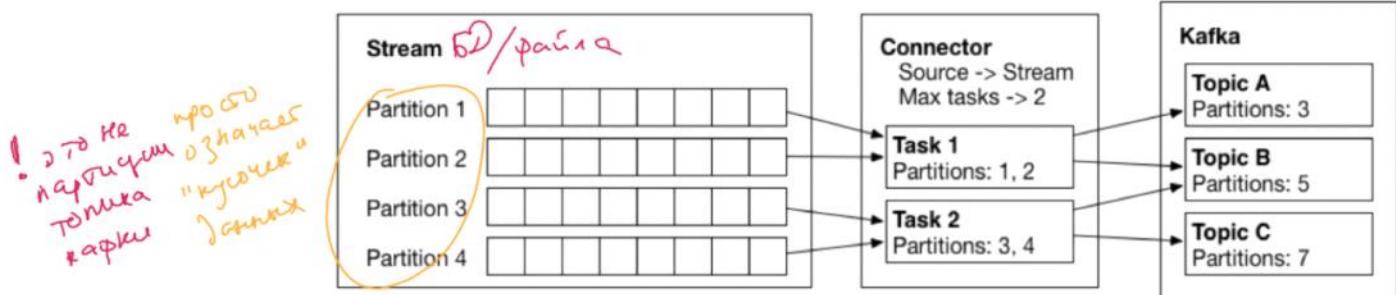
In case of importing CSV files, the **partitions** would correspond to **files**

в kafka коннект термин партиция обозначает не партицию_кафки а просто кусочек данных исходной системы (но идея похожа: например партиция-файл не может быть на значена на два таска)

- см выше: в БД-коннекторе партиция это таблица, а НЕ ее части ,
- в FILE-коннекторе партиция это файл. а НЕ куски одного и того же файла

Providing Parallelism and Scalability

- Partitioning → parallelism & scalability
- 1 Connector job → 1..n **tasks**
- 1 Worker → 1..n **tasks**
- 1 Task per Thread



- Splitting the workload into smaller pieces provides the parallelism and scalability
- Connector jobs are broken down into **tasks** that do the actual copying of the data
- **Workers** are processes running one or more tasks, each in a different thread

A Connector is a Connector class and a configuration. Each connector defines and updates a set of tasks that actually copy the data. **Connect** distributes these tasks across **workers**.

In the case of Connect, the term **partition** can mean any subset of data in the source or the sink. How a partition is represented depends on the type of Connector (e.g., tables are partitions for the JDBC connector, files are the partitions for the FileStream connector).

Worker processes are not managed by Kafka. Other tools can do this: Kubernetes, Mesos, YARN, Chef, Puppet, custom scripts, etc.



not all connectors support multiple tasks and parallelism. For example, the **syslog** source connector only supports one task.

коннектор отслеживает (и записывает себе в специальный топик) позицию прочетных данных файла источника или таблицы БД

- для обеспечения EOS в sink-коннекторах требуется хранить позицию прочетных данных не в кафке (например если sink=файл то и позиция в отдельном файле)

Source and Sink Offsets

- Kafka Connect tracks the produced and consumed offsets so it can restart at the correct place, in case of failure
- What the offset corresponds to depends on the Connector, for example:
 - File input: offset → position in file
 - Database input: offset → timestamp or sequence id
- The method of tracking the offset depends on the specific Connector
 - Each Connector can determine its own way of doing this
- Examples of source and sink offset tracking methods:
 - JDBC source in distributed mode: a special Kafka Topic
 - HDFS sink: an HDFS file
 - FileStream source: a separate local file

As with Kafka topics, the position within the partitions used by Connect must be tracked as well to prevent the unexpected replay of data. As with the partitions, the object used as the offset varies from connector to connector, as does the method of tracking the offset. The most common way to track the offset is through the use of a Kafka topic, though the Connector developer can use whatever he or she wants.

The variation is primarily with source connectors. With sink connectors, the offsets are Kafka topic partition offsets, and they're usually stored in the external system where the data is being written.

Configuring the Connector

Parameter	Description
<code>name</code>	Connector's unique name
<code>connector.class</code>	Connector's Java class
<code>tasks.max</code>	Maximum tasks to create. The Connector may create fewer if it cannot achieve this level of parallelism (Default: 1)
<code>key.converter</code>	(optional) Override the worker key converter
<code>value.converter</code>	(optional) Override the worker value converter
<code>topics</code> (Sink connectors only)	List of input Topics (to consume from)

используем landoop UI

```

version: '2'

services:
  # this is our kafka cluster.
  kafka-cluster:
    image: landoop/fast-data-dev:cp3.3.0
    environment:
      ADV_HOST: 127.0.0.1      # Change to 192.168.99.100 if using Docker Toolbox
      RUNTESTS: 0              # Disable Running tests so the cluster starts faster
    ports:
      - 2181:2181              # Zookeeper
      - 3030:3030              # Landoop UI
      - 8081-8083:8081-8083   # REST Proxy, Schema Registry, Kafka Connect ports
      - 9581-9585:9581-9585   # JMX Ports
      - 9092:9092              # Kafka Broker

  # we will use elasticsearch as one of our sinks.
  # This configuration allows you to start elasticsearch
  elasticsearch:
    image: itzg/elasticsearch:2.4.3
    environment:
      PLUGINS: appbaseio/dejavu
      OPTS: -Dindex.number_of_shards=1 -Dindex.number_of_replicas=0
    ports:
      - "9200:9200"

  # we will use postgres as one of our sinks.
  # This configuration allows you to start postgres
  postgres:
    image: postgres:9.5-alpine
    environment:
      POSTGRES_USER: postgres  # define credentials
      POSTGRES_PASSWORD: postgres # define credentials
      POSTGRES_DB: postgres    # define database
    ports:
      - 5432:5432              # Postgres port

```

как создать distributed connector через UI

Kafka Development Environment
docker container powered by Landoop

SCHEMAS 5
SCHEMA REGISTRY UI manage avro schemas
ENTER

TOPICS 12
KAFKA TOPICS UI browse topics and data
ENTER

CONNECTORS 5
KAFKA CONNECT UI setup, manage connectors
ENTER

BROKERS 1
KAFKA BROKERS management and monitoring
COMING SOON

RUNNING SERVICES
fast-data-dev

KAFKA CONNECT

5 Connectors **NEW**

Search connectors

- logs-broker
- logs-connect-distributed
- logs-rest-proxy
- logs-schema-registry
- logs-zookeeper

Kafka Connect : /api/kafka-connect
Kafka Connect Version : 0.10.1.1
Kafka Connect UI Version : 0.9.0

New Connector

Search

Sources

- twitter** Subscribe to feeds using the Twitter API and stream data into a kafka topic
- ftp** Tail directories in remote FTP location and bring messages in Kafka
- file** Read files and stream data into Kafka Topics
- bloomberg** A connector to subscribe to Bloomberg feeds via the Bloomberg Ibis open API and write data to Kafka

KAFKA CONNECT

5 Connectors

Search connectors

- logs-broker
- logs-connect-distributed
- logs-rest-proxy
- logs-schema-registry
- logs-zookeeper

Kafka Connect : /api/kafka-connect
Kafka Connect Version : 0.10.1.1
Kafka Connect UI Version : 0.9.0

New Connector (Sink): FileStreamSinkConnector

C FileStreamSinkConnector
class: org.apache.kafka.connect.file.FileStreamSinkConnector

Show cURL command

```

1 These are standard kafka connect parameters, need for ALL connector
2 name=file-stream-demo-distributed
3 connector.class=org.apache.kafka.connect.file.FileStreamSourceConnector
4 tasks.max=1
5 # Parameters can be found here: https://github.com/apache/kafka/blob/trunk/connector/file/src/main/java/org/apache
6 file=demo-file.txt
7 topic=demo-2-distributed
8 # Added configuration for the distributed mode:
9 key.converter=org.apache.kafka.connect.json.JsonConverter
10 key.converter.schemas.enable=true
11 value.converter=org.apache.kafka.connect.json.JsonConverter
12 value.converter.schemas.enable=true

```

connector properties

Working task

... Connectors with Confluent Control Center

18 декабря 2020 г. 11:40

Configuring Connectors with Confluent Control Center

The left screenshot shows the 'Browse' page of the Confluent Control Center, displaying a list of available connectors: ActiveMQSource Connector, JdbcSource Connector, and FileStreamSink Connector. The right screenshot shows the 'Add Connector' page for a JDBC source connector, with fields for connector class (JdbcSourceConnector), name (My-JDBC-Source), tasks max (3), key converter class, value converter class, and header converter class.

Connectors can also be configured in Confluent Control Center. We can define source and sink connectors there.



Not all connectors may be fully configurable in the Control Center. In this case one has to use the Connect REST API to create and configure the connector. In the hands-on lab we give an example of this.

я могу динамически менять конфигурацию работающего коннектора

connect REST API

4 декабря 2020 г. 18:31

Using the REST API

Some important REST endpoints

Method	Path	Description
GET	/connectors	Get a list of active connectors
POST	/connectors	Create a new Connector
GET	/connectors/(string: name)/config	Get configuration information for a Connector
PUT	/connectors/(string: name)/config	Create a new Connector, or update the configuration of an existing Connector

More information on the REST API can be found at <https://docs.confluent.io/current/connect/references/restapi.html>

- Add, modify delete connectors
- Distributed Mode:
 - Config **only** via REST API
 - Config stored in Kafka topic
 - REST call to **any** worker
- Standalone Mode:
 - **Config also via REST API**
 - **Changes not persisted!**
- Control Center uses REST API

```
1  "name": "Operators-Connector",
2  "config": {
3    "connector.class": "io.confluent.connect.jdbc.JdbcSourceConnector",
4    "connection.url": "jdbc:postgresql://postgres:5432/postgres",
5    "connection.user": "postgres",
6    "table.whitelist": "operators",
7    "mode": "incrementing",
8    "incrementing.column.name": "id",
9    "table.type": "TABLE",
10   "topic.prefix": "pg_",
11   "numeric.mapping": "best_fit",
12   "transforms": "createKey,extractInt",
13   "transforms.createKey.type": "org.apache.kafka.connect.transforms.ValueMapper$ValueMapperBuilder",
14   "transforms.createKey.fields": "id",
15   "transforms.extractInt.type": "org.apache.kafka.connect.transforms.Extractor$ExtractorBuilder",
16   "transforms.extractInt.field": "id"
17 }
18
19 }
```

- Connectors can be added, modified, and deleted via a REST API on port 8083
- In distributed mode, configuration can be done only via this REST API
 - Changes made this way will persist after a worker process restart
 - Connector configuration data is stored in a special Kafka Topic
 - The REST requests can be made to any worker
- In standalone mode, configuration can also be done via a REST API
 - However, typically configuration is done via a properties file
 - Changes made via the REST API when running in standalone mode will not persist after worker restart
- Confluent Control Center leverages this REST API to let users configure and manage connectors through the GUI

To make changes to connectors, the Worker will also run a REST API as part of its code. This allows HTTP calls to be sent to any of the Workers that will then configure the connectors. Changes made via the REST API take effect immediately and without a reboot.



The REST API included with a Connect Worker is different from the Confluent REST Proxy.

kafka connect полностью работает только через REST API

Kafka Connect REST API



- All the actions performed by Landoop Kafka Connect UI are actually triggering REST API calls to Kafka Connect.
- Let's learn about those
(<http://docs.confluent.io/3.2.0/connect/managing.html#common-rest-examples>)

- | | |
|---|-----------------------------------|
| 1. Get Worker information | 5. Get Connector Status |
| 2. List Connectors available on a Worker | 6. Pause / Resume a Connector |
| 3. Ask about Active Connectors | 7. Delete our Connector |
| 4. Get information about a Connector Tasks and Config | 8. Create a new Connector |
| | 9. Update Connector configuration |
| | 10. Get Connector Configuration |

The screenshot shows a browser window with the URL docs.confluent.io/3.2.0/connect/managing.html#common-rest-examples. The page title is "Configuration Options". The left sidebar contains navigation links for the Confluent Platform, including "What is the Confluent Platform?", "Confluent Platform 3.1.1 Release Notes", "Data Serialization and Evolution", "Installation", "Upgrade", "Confluent Platform Quickstart", "Application Development", "Operations", "Docker", "Confluent Control Center", "Kafka Security", "Kafka Streams", and "Kafka Connect". The main content area is titled "Configuration Options" and lists configuration parameters for the "Confluent Elasticsearch Connector". The parameters shown are:

- connection.url**: The URL to connect to Elasticsearch. Type: string, Default: "", Importance: high.
- type.name**: The type to use for each index. Type: string, Default: "", Importance: high.
- key.ignore**: Whether to ignore the key during indexing. When this is set to true, only the value from the message will be written to Elasticsearch. Note that this is a global config that applies to all topics. If this is set to true, Use `topic.key.ignore` to config for different topics. This value will be overridden by the per topic configuration. Type: boolean, Default: false, Importance: high.
- batch.size**: The number of requests to process as a batch when writing to Elasticsearch.

пример

```
#!/bin/bash

# let's start a command line to all have linux commands
docker run --rm -it --net=host landoop/fast-data-dev bash
# Install jq to pretty print json
apk update && apk add jq

# Examples are taken from here: http://docs.confluent.io/3.2.0/connect/managing.html#common-rest-examples
# Replace 127.0.0.1 by 192.168.99.100 if you're using docker toolbox
```

1) Get Worker information
`curl -s 127.0.0.1:8083/ | jq`

2) List Connectors available on a Worker
`curl -s 127.0.0.1:8083/connector-plugins | jq`

```
root@fast-data-dev / $ curl -s 127.0.0.1:8083/ | jq
{
  "version": "0.10.1.1",
  "commit": "f10ef2720b63b247"
}
```

3) Ask about Active Connectors
`curl -s 127.0.0.1:8083/connectors | jq`

4) Get information about a Connector Tasks and Config
`curl -s 127.0.0.1:8083/connectors/source-twitter-distributed/tasks | jq`

```

root@fast-data-dev / $ curl -s 127.0.0.1:8083/connectors/source-twitter-distributed/tasks | jq
[{"id": {"connector": "source-twitter-distributed", "task": 0}, "config": {"connector.class": "com.eneco.trading.kafka.connect.twitter.TwitterSourceConnector", "twitter.token": "2932095118-2451tfsj2H01M0rRq0MFNlPHTIxER4e3Cj2ho0", "tasks.max": "1", "twitter.secret": "7kAn2qXA7wJ8k46ySj1lr19Reey9h2NLAxLPtFUKjVgb", "track.terms": "programming,java,kafka,scala", "language": "en", "key.converter.schemas.enable": "true", "task.class": "com.eneco.trading.kafka.connect.twitter.TwitterSourceTask", "value.converter.schemas.enable": "true", "name": "source-twitter-distributed", "topic": "demo-3-twitter", "twitter.consumersecret": "ggQgZ1pZLcexZ0PjTX1zRtG0bhHuWFHmx4DxGfRAVD2xwrHR4", "value.converter": "org.apache.kafka.connect.json.JsonConverter", "key.converter": "org.apache.kafka.connect.json.JsonConverter", "twitter.consumerkey": "9jT7rYU09ELMSibxgVK7tRpRs"}]

```

The screenshot shows the Kafka Connect UI interface. On the left, there's a sidebar with '8 Connectors' listed: file-stream-demo-distributed, logs-broker, logs-connect-distributed, logs-rest-proxy, logs-schema-registry, logs-zookeeper, sink-elastic-twitter-distributed, and source-twitter-distributed. The 'source-twitter-distributed' connector is selected. The main panel displays its configuration details:

- SOURCE : source-twitter-distributed**
- Connector status:** RUNNING
- TOPICS:** (empty)
- CONFIGURATION:** (empty)
- CODE:** (Shows the JSON configuration code from the previous terminal output.)
- TASKS:** (empty)

5) Get Connector Status

```
curl -s 127.0.0.1:8083/connectors/file-stream-demo-distributed/status | jq
```

6) Pause / Resume a Connector (no response if the call is successful)

```
curl -s -X PUT 127.0.0.1:8083/connectors/file-stream-demo-distributed/pause
curl -s -X PUT 127.0.0.1:8083/connectors/file-stream-demo-distributed/resume
```

7) Get Connector Configuration

```
curl -s 127.0.0.1:8083/connectors/file-stream-demo-distributed | jq
```

8) Delete our Connector

```
curl -s -X DELETE 127.0.0.1:8083/connectors/file-stream-demo-distributed
```

9) Create a new Connector

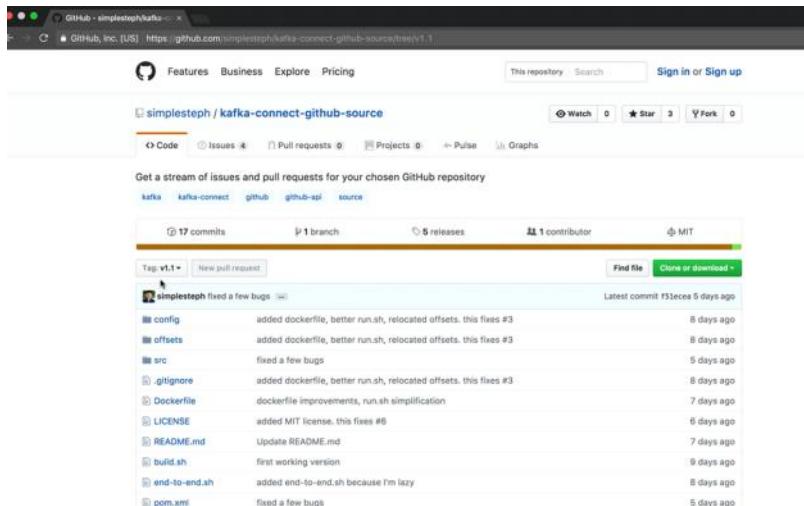
```
curl -s -X POST -H "Content-Type: application/json" --data '{"name": "file-stream-demo-distributed", "config": {"connector.class": "org.apache.kafka.connect.file.FileStreamSourceConnector", "key.converter.schemas.enable": "true", "file": "demo-file.txt", "tasks.max": "1", "value.converter.schemas.enable": "true", "name": "file-stream-demo-distributed", "topic": "demo-2-distributed", "value.converter": "org.apache.kafka.connect.json.JsonConverter", "key.converter": "org.apache.kafka.connect.json.JsonConverter", "key.converter.json.JsonConverter": "http://127.0.0.1:8083/connectors"} }' http://127.0.0.1:8083/connectors | jq
```

10) Update Connector configuration

```
curl -s -X PUT -H "Content-Type: application/json" --data
'{"connector.class":"org.apache.kafka.connect.file.FileStreamSourceConnector","key.converter.schemas.enable":"true"
,"file":"demo-file.txt","tasks.max":"2","value.converter.schemas.enable":"true","name":"file-stream-demo-
distributed","topic":"demo-2-
distributed","value.converter":"org.apache.kafka.connect.json.JsonConverter","key.converter":"org.apache.kafka.conn
ect.json.JsonConverter"}' 127.0.0.1:8083/connectors/file-stream-demo-distributed/config | jq
```

Section goals

- Understand how connectors are made, using the [GitHubSourceConnector](#) as an example
 - Dependencies
 - ConfigDef
 - Connector
 - Schema & Struct
 - Source Partition & Source Offsets
 - Task
- Learn how to deploy your connectors (or any connector)
 - Package Jars
 - Run jars in standalone mode
 - Deploy Jars
- The code can be found at <https://github.com/simplesteph/kafka-connect-github-source/tree/v1.1>
- It's my own code that I put together specially for this class!
- **Knowledge of Java is necessary for this section**
- I recommend you install IntelliJ Community Editions on your laptop to have the same environment as me:
<https://www.jetbrains.com/idea/download/>
- You need to install Java JDK 8 (Java SE Development Kit) as well:
<http://www.oracle.com/technetwork/java/javase/downloads/jdk8-downloads-2133151.html>
- <https://www.confluent.io/product/connectors/>
 - List of open source connectors
 - Read the code!!!
- <https://www.confluent.io/wp-content/uploads/Partner-Dev-Guide-for-Kafka-Connect.pdf?x18424>
 - General guidelines
- <http://docs.confluent.io/3.2.0/connect/devguide.html>



GitHub - simplesteph/kafka-connect-github-source

simplesteph / kafka-connect-github-source

Code Issues Pull requests Projects Pulses Graphs

Get a stream of issues and pull requests for your chosen GitHub repository

kafka kafka-connect github github-api source

17 commits 1 branch 5 releases 1 contributor MIT

Tag: v1.1 New pull request

simplesteph fixed a few bugs

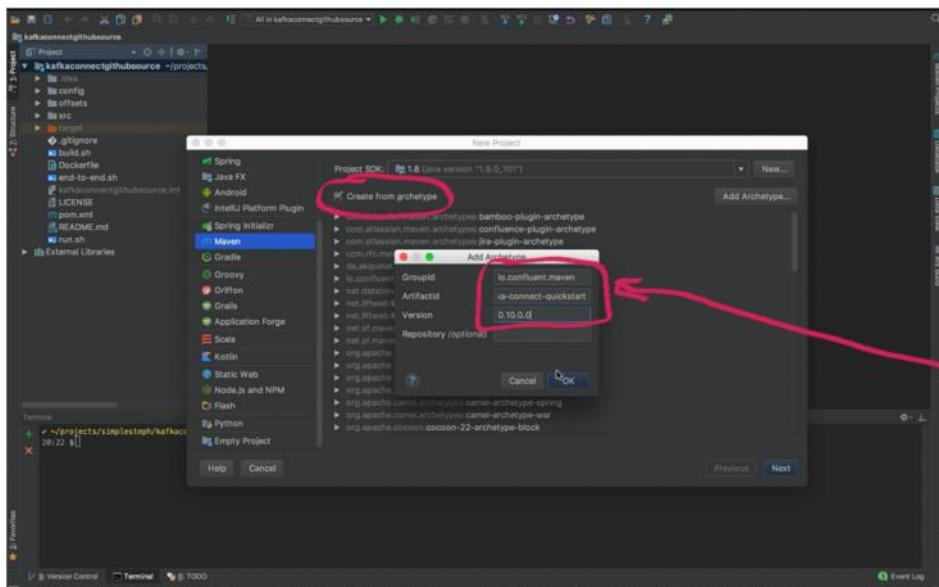
File	Commit Message	Time Ago
config	added dockerfile, better run.sh, relocated offsets. this fixes #3	8 days ago
offsets	added dockerfile, better run.sh, relocated offsets. this fixes #3	8 days ago
src	fixed a few bugs	5 days ago
.gitignore	added dockerfile, better run.sh, relocated offsets. this fixes #3	8 days ago
Dockerfile	dockerfile improvements, run.sh simplification	7 days ago
LICENSE	added MIT license, this fixes #6	8 days ago
README.md	Update README.md	7 days ago
build.sh	first working version	8 days ago
end-to-end.sh	added end-to-end.sh because I'm lazy	8 days ago
pom.xml	fixed a few bugs	5 days ago

intelliJ IDEA проект будет создаваться из архетипа

Maven - Archetype



- The easiest way to get started on a Kafka Connect project is to use an archetype provided here:
<https://github.com/jcugenborder/kafka-connect-archetype>
- We will do a walkthrough on IntelliJ to install and use that archetype
- We will change the Kafka dependency to be 0.10.2.0 (use the one matching the version from Landoop)



Issues | GitHub Developer Guide | kubernetes/kafka-connect-archetype | https://api.github.com/repos/kubernetes/kafka-connect-archetype/branches/master | 10 commits | 2 branches | 1 release

Branch: master | New pull request

joustenborder Added parent. Changed the package to my personal namespace. (#1)

.gitignore Initial commit for Confluent Platform 3.0.0

Jenkinsfile Changed the package to my personal namespace. (#1)

LICENSE Changed versioning to follow Kafka versions. Added test cases.

README.md Changed versioning to follow Kafka versions. Added test cases.

pom.xml Added parent.

READEME.md

This maven quickstart is used to generate a skeleton plugin for Kafka Connect.

```
mvn archetype:generate -DarchetypeGroupId=io.confluent.maven -DarchetypeArtifactId=kafka-connect-quicke...
```

Connector



- That's the class of your Kafka Connector that is referenced from the configs. It should load the config, and create a few tasks.

```
public String version()
public void start(Map<String, String> map)
public Class<? extends Task> taskClass()
public List<Map<String, String>> taskConfigs(int i)
public void stop()
public ConfigDef config()
```

```
GitHubSourceConnector
1 package com.simplesteph.kafka;
2
3 import ...
4
5 public class GitHubSourceConnector extends SourceConnector {
6     private static Logger log = LoggerFactory.getLogger(GitHubSourceConnector.class);
7     private GitHubSourceConnectorConfig config;
8
9     @Override
10    public String version() { return VersionUtil.getVersion(); }
11
12    @Override
13    public void start(Map<String, String> map) { config = new GitHubSourceConnectorConfig(map); }
14
15    @Override
16    public Class<? extends Task> taskClass() { return GitHubSourceTask.class; }
17
18    @Override
19    public List<Map<String, String>> taskConfigs(int i) {
20        // Define the individual task configurations that will be executed.
21        ArrayList<Map<String, String>> configs = new ArrayList<>(initialCapacity: 1);
22        configs.add(config.originalsStrings());
23        return configs;
24    }
25
26    @Override
27    public void stop() {
28        // Do things that are necessary to stop your connector.
29        // nothing is necessary to stop for this connector
30    }
31}
```

... как задеплоить коннектор (distributed mode)

5 декабря 2020 г. 0:25

confluent-hub install как задеплоить стандартный коннектор, который есть на confluent hub

From Confluent Hub

- Use `confluent-hub` client included with Confluent Platform

```
$ confluent-hub install debezium/debezium-connector-mysql:latest
```

как задеплоить коннектор из обычного джарника

plugin.path в конфиге нужно указать путь к джарнику коннектора на сервере

From other sources

- Package Connector in JAR file
- **Install JAR file on all Kafka Connect worker machines**
 - Connectors are installed as plugins
 - There is **library isolation** between plugins

```
plugin.path=/path/to/my/plugins
```

Prior to Kafka 0.11.0, JARs were placed in a path specified by `CLASSPATH` and there was no library isolation between connectors

```
$ export CLASSPATH="$CLASSPATH:/path/to/my/connectors/*"
```

QUESTION: Why is library isolation important?

Answer: Connectors built by different developers might use different versions of the same library.

У всех коннекторов должен быть единий способ общая очередь для хранения общего конфига

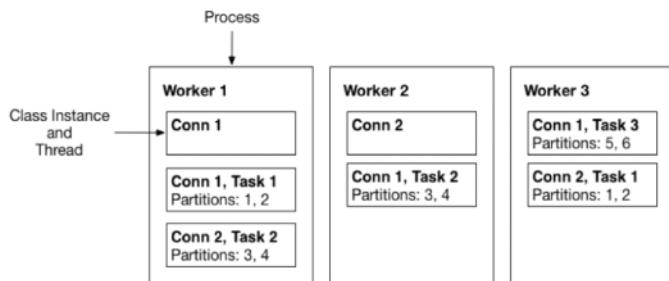
Running in Distributed Mode

Start Kafka Connect on **each** worker node

```
$ connect-distributed connect-distributed.properties
```

Group coordination

- Connect leverages Kafka's group membership protocol
 - Configure workers with the same `group.id`
- Workers distribute load within this Kafka Connect "cluster"



In distributed mode, the connector configurations **cannot** be kept on the Worker systems; a failure of a worker should not make the configurations unavailable. Instead, **distributed workers** keep their connector configurations in a special Kafka topic which is specified in the **worker configuration file**.

Workers coordinate their tasks to distribute the workload using the same mechanisms as Consumer Groups.

как коннектор задеплоить в distributed через GUI

Deploying your Connectors

1. Extract the jars
2. Mount the jars into your Landoop Image
3. Use the Connect UI to launch the connector!

положить все jar моего коннектора в папку /connectors на сервере с кафкой

```
13:18 $ docker run -it --rm -p 2181:2181 -p 3030:3030 -p 8081:8081 -p 8082:8082 -p 8083:8083 -p 9092:9092 -e A
DV_HOST=127.0.0.1 -e RUNTESTS=0 -v ~/projects/simplesteph/kafkaconnectgithubsource/target/kafka-connect-github
-source-1.1-package/share/java/kafka-connect-github-source:/connectors landoop/fast-data-dev
```

коннектор автоматически подхватится и отобразится в UI

The screenshot shows the Kafka Development Environment dashboard at localhost:3030. The interface has four main cards: Schemas (3), Topics (7), Connectors (1), and Brokers (1). Below each card are links to their respective UIs. A context menu is open over the CONNECTORS card, listing options like 'Open Link in New Tab', 'Save Link As...', and 'Inspect'. The 'CONNECTORS' card also displays the text 'Your set up'.

The screenshot shows the Kafka Connect UI at localhost:3030/kafka-connect-ui/#/cluster/fast-data-dev/select-connector. It lists various connectors: Bloomberg, Redis, JMS, Kudu, MQTT, JMS, RethinkDB, HDFS, CoAP, File, Schemas, and GitHubSourceConnector. Each connector has a brief description and a small icon.

как задеплоить в продакшен

Setting up Kafka Connect in Production (2/2)



To setup Kafka Connect in production

1. Download Kafka Binaries
2. Set up the connect-distributed.properties as needed
3. Add your jars where needed (plugins.path or classpath)
4. And launch Kafka Connect!
5. (optional) Setup the Landoop Kafka Connect UI

You will be able to interact with Kafka Connect using the REST API, or the UI if you set it up

```
#!/bin/bash

# 0. Make sure Java is installed on your machine https://java.com/en/download/
# 1. Download Kafka at https://kafka.apache.org/downloads (>= 0.11.0.1)
# 2. Unzip Kafka to a directory of your choice (for example ~/kafka)

# 3. Open a terminal or text editor on the Kafka directory
# 4. Edit the file at config/connect-distributed.config
nano config/connect-distributed.properties
vi config/connect-distributed.properties
atom config/connect-distributed.properties

# 5. change bootstrap.servers to the correct ones
# 6. change rest.port=8083 to rest.port=8084 (or any available port)
# 7. Optional (if you want a separate cluster) -
# 7a. Set offset.storage.topic=connect-offsets-2
# 7b. Set config.storage.topic=connect-configs-2
# 7c. Set status.storage.topic=connect-status-2
# 7d. Set group.id=connect-cluster-2
# 8. Set plugin.path=/directory/of/your/choice
# 9. Adjust any other settings that matters to your setup (see doc)
# 10. Place any connector (jars) you need in the plugin.path sub-directory

# 11. In order to add connect workers, duplicate the connect-distributed file and change the rest.port if
running on the same host. The rest of the settings are the same

# 12. Optionally, setup Kafka Connect UI from Landoop
# Instructions are at: https://github.com/landoop/kafka-connect-ui
# You can build from source or use docker for that (see their github)
```

/etc/kafka/ файл конфига всегда лежит в одной и той же папке

- You can modify Connect configuration settings
 - Distributed mode in /etc/kafka/connect-distributed.properties
 - Standalone mode in /etc/kafka/connect-standalone.properties
 - <http://docs.confluent.io/current/connect/>

config/connect-distributed.properties подправим настройки конфигов на всех коннектор-нодах кафки

```
~/kafka_2.11-0.11.0.1> ll config
total 128
-rw-r--r--@ 1 stephanemaarek staff 906B 6 Sep 04:18 connect-console-sink.properties
-rw-r--r--@ 1 stephanemaarek staff 909B 6 Sep 04:18 connect-console-source.properties
-rw-r--r--@ 1 stephanemaarek staff 5.7K 6 Sep 04:18 connect-distributed.properties
-rw-r--r--@ 1 stephanemaarek staff 883B 6 Sep 04:18 connect-file-sink.properties
-rw-r--r--@ 1 stephanemaarek staff 881B 6 Sep 04:18 connect-file-source.properties
-rw-r--r--@ 1 stephanemaarek staff 1.1K 6 Sep 04:18 connect-log4j.properties
-rw-r--r--@ 1 stephanemaarek staff 2.7K 6 Sep 04:18 connect-standalone.properties
-rw-r--r--@ 1 stephanemaarek staff 1.2K 6 Sep 04:18 consumer.properties
-rw-r--r--@ 1 stephanemaarek staff 4.6K 6 Sep 04:18 log4j.properties
-rw-r--r--@ 1 stephanemaarek staff 1.9K 6 Sep 04:18 producer.properties
-rw-r--r--@ 1 stephanemaarek staff 6.8K 6 Sep 04:18 server.properties
-rw-r--r--@ 1 stephanemaarek staff 1.0K 6 Sep 04:18 tools-log4j.properties
-rw-r--r--@ 1 stephanemaarek staff 1.0K 6 Sep 04:18 zookeeper.properties
```

A list of host/port pairs to use for establishing the initial connection to the Kafka cluster. **указать бутстрап сервер кафки**
bootstrap.servers=127.0.0.1:9092

unique name for the cluster, used in forming the Connect cluster group. Note that this must not conflict with consumer group IDs **лучше создать свой кластер коннекторов (у всех коннекторов единица групппы)**
group.id=connect-cluster-2

The converters specify the format of data in Kafka and how to translate it into Connect data. Every Connect user will
need to configure these based on the format they want their data in when loaded from or stored into Kafka
key.converter=org.apache.kafka.connect.json.JsonConverter
value.converter=org.apache.kafka.connect.json.JsonConverter
Converter-specific settings can be passed in by prefixing the Converter's setting with the converter we want to apply
it to
key.converter.schemas.enable=true
value.converter.schemas.enable=true

The internal converter used for offsets, config, and status data is configurable and must be specified, but most users will
always want to use the built-in default. Offset, config, and status data is never visible outside of Kafka Connect in this format.
internal.key.converter=org.apache.kafka.connect.json.JsonConverter
internal.value.converter=org.apache.kafka.connect.json.JsonConverter
internal.key.converter.schemas.enable=false
internal.value.converter.schemas.enable=false

Topic to use for storing offsets. This topic should have many partitions and be replicated and compacted.
Kafka Connect will attempt to create the topic automatically when needed, but you can always manually create
the topic before starting Kafka Connect if a specific topic configuration is needed.
Most users will want to use the built-in default replication factor of 3 or in some cases even specify a larger value.
Since this means there must be at least as many brokers as the maximum replication factor used, we'd like to be able
to run this example on a single-broker cluster and so here we instead set the replication factor to 1. **лучше создать свой топик для хранения оффсетов**
offset.storage.topic=connect-offsets-2
offset.storage.replication.factor=1
#offset.storage.partitions=25

Topic to use for storing connector and task configurations; note that this should be a single partition, highly replicated,
and compacted topic. Kafka Connect will attempt to create the topic automatically when needed, but you can always manually create
the topic before starting Kafka Connect if a specific topic configuration is needed.
Most users will want to use the built-in default replication factor of 3 or in some cases even specify a larger value.
Since this means there must be at least as many brokers as the maximum replication factor used, we'd like to be able
to run this example on a single-broker cluster and so here we instead set the replication factor to 1. **лучше создать свой топик для хранения конфигов**
config.storage.topic=connect-configs-2
config.storage.replication.factor=1

Topic to use for storing statuses. This topic can have multiple partitions and should be replicated and compacted.
Kafka Connect will attempt to create the topic automatically when needed, but you can always manually create
the topic before starting Kafka Connect if a specific topic configuration is needed.
Most users will want to use the built-in default replication factor of 3 or in some cases even specify a larger value.
Since this means there must be at least as many brokers as the maximum replication factor used, we'd like to be able
to run this example on a single-broker cluster and so here we instead set the replication factor to 1. **лучше создать свой топик для хранения статусов**
status.storage.topic=connect-status-2
status.storage.replication.factor=1
#status.storage.partitions=5

Flush much faster than normal, which is useful for testing/debugging
offset.flush.interval.ms=10000

These are provided to inform the user about the presence of the REST host and port configs
Hostname & Port for the REST API to listen on. If this is set, it will bind to the interface used to listen to requests. **на второй коннектор ноде нужно указать свой порт, например 8085**
#rest.host.name=
rest.port=8084

The Hostname & Port that will be given out to other workers to connect to i.e. URLs that are routable from other servers.
#rest.advertised.host.name=
#rest.advertised.port=

Set to a list of filesystem paths separated by commas (,) to enable class loading isolation for plugins
(connectors, converters, transformations). The list should consist of top level directories that include
any combination of:
a) directories immediately containing jars with plugins and their dependencies
b) uber-jars with plugins and their dependencies
c) directories immediately containing the package directory structure of classes of plugins and their dependencies
Examples:
plugin.path=/usr/local/share/java,/usr/local/share/kafka/plugins,/opt/connectors, **путь к моему коннектору на сервере**
plugin.path=/Users/stephanemaarek/kafka_2.11-0.11.0.1/connectors

поместим на коннектор-нода в папку /connectors с коннекторами яг-ник коннекторов



запустим на каждой коннек-ноде

```
~/kafka_2.11-0.11.0.1 ➤ bin/connect-distributed.sh config/connect-distributed.properties
```

```
[2017-10-05 18:57:39,867] INFO Successfully joined group connect-cluster-2 with generation 1 (org.apache.kafka.clients.consumer.internals.AbstractCoordinator:409)
[2017-10-05 18:57:39,867] INFO Joined group and got assignment: Assignment{error=0, leader='connect-1-4acb32d5-5f35-4106-ad84-1d29bdb84f11', leaderUrl='http://192.168.1.5:8084', offset=-1, connectorIds=[], taskIds[]} (org.apache.kafka.connect.runtime.distributed.DistributedHerder:1166)
[2017-10-05 18:57:39,868] INFO Starting connectors and tasks using config offset -1 (org.apache.kafka.connect.runtime.distributed.DistributedHerder:815)
[2017-10-05 18:57:39,868] INFO Finished starting connectors and tasks (org.apache.kafka.connect.runtime.distributed.DistributedHerder:825)

consumer.internals.AbstractCoordinator:409)
[2017-10-05 19:06:52,433] INFO Joined group and got assignment: Assignment{error=0, leader='connect-1-71ac447d-878a-4d3a-a660-44e741a4bef7', leaderUrl='http://192.168.1.5:8084', offset=-1, connectorIds=[], taskIds[]} (org.apache.kafka.connect.runtime.distributed.DistributedHerder:1166)
```

проверим что обе ноды с коннекторами работают

```
~/kafka_2.11-0.11.0.1 ➤ curl localhost:8085
{"version":"0.11.0.1","commit":"c2a0d5f9bf45bf5"} ➤
~/kafka_2.11-0.11.0.1 ➤ curl localhost:8084
{"version":"0.11.0.1","commit":"c2a0d5f9bf45bf5"} ➤
```

проверим что появились топики (connect-offset-2) для кластера-коннекторов

TOTAL TOPICS	BROKERS
13	1

KAFKA REST /api/kafka-rest-proxy
KAFKA-TOPICS-UI v0.9.2

установим landoop connect UI (можно даже на отдельном докере) и указать URL к коннект-нодам

```
✗ ➤ ~/kafka_2.11-0.11.0.1 ➤ docker run --net=host --rm -it \
-e "CONNECT_URL=http://127.0.0.1:8084" \
landoop/kafka-connect-ui
```

Apache Download Mirrors | Landoop/kafka-connect-ui | Landoop Kafka Development | Stéphane Perrin

Code Issues Pull requests Projects Wiki Insights

Web tool for Kafka Connect | http://landoop.com/demos

kafka kafka-connect elasticsearch cassandra s3 documentdb redis jms mqtt hdfs influxdb twitter

122 commits 2 branches 9 releases 9 contributors

Branch: master New pull request Create new file Upload file Find file Clone or download

- andmarlos Docker: option to set custom cluster names (e.g. -e "CONNECT_URL=http://...") 1 Latest commit at7idle on 10 Aug
- docker Docker: option to set custom cluster names (e.g. -e "CONNECT_URL=http://...") 2 months ago
- src Update to exclude JMS Topic In JMSSource 5 months ago
- .gitignore initial commit 10 months ago
- Gruntfile.js pagination fix 7 months ago
- LICENSE.md initial commit 10 months ago
- README.md correct changing Ic 5 months ago
- bower.json cleanup dependencies 10 months ago
- env.js Different header background color for each cluster, show version at f... 10 months ago
- index.html box shadow - 3.2 fix - ui re-arrangement 6 months ago
- package.json version 0.9.2 5 months ago
- README.md

kafka-connect-ui

Apache Download Mirrors | Landoop/kafka-connect-ui | Landoop Kafka Development | Stéphane Perrin

Code Issues Pull requests Projects Wiki Insights

Kafka-Connect-UI

release v0.9.2 docker pulls 79 11 chat on gitter

This is a web tool for Kafka Connect for setting up and managing connectors for multiple connect clusters.

Live Demo

kafka-connect-ui.landoop.com

Run standalone with docker

```
docker run --rm -it -p 8000:8000 \
-e "CONNECT_URL=http://connect.distributed.url" \
landoop/kafka-connect-ui
```

The CONNECT_URL can be a comma separated array of Connect worker endpoints. E.g:
CONNECT_URL=http://connect.1.url,http://connect.2.url"

Web UI will be available at <http://localhost:8000>

Build from source

```
git clone https://github.com/Landoop/kafka-connect-ui.git
cd kafka-connect-ui
npm install
http-server .
```

Web UI will be available at <http://localhost:8000>

Apache Download Mirrors | Landoop/kafka-connect-ui | Landoop Kafka Development | Stéphane Perrin

Code Issues Pull requests Projects Wiki Insights

KAFKA CONNECT

0 Connectors NEW

No connectors found

Kafka Connect : /api/kafka-connect-1
Kafka Connect Version : 0.11.0.1
Kafka Connect UI Version : 0.9.2

Dashboard EXPORT CONFIG

SINK CONNECTORS 0 SOURCE CONNECTORS 0 TOPICS USED BY CONNECTORS 0

Connect topology

Powered by [LANDOOP](#)

.... ДОПОЛНИТЕЛЬНЫЕ НАСТРОЙКИ КОНЕКТОРА

18 декабря 2020 г. 11:29

Configuring Workers: Distributed Mode

Parameter	Description
group.id	A unique string that identifies the Kafka Connect cluster group the worker belongs to
session.timeout.ms	Timeout used to detect failures when using Kafka's group management facilities
heartbeat.interval.ms	Expected time between heartbeats to the group coordinator when using Kafka's group management facilities. Must be smaller than session.timeout.ms
config.storage.topic	Topic in which to store Connector & task config. data
offset.storage.topic	Topic in which to store offset data for Connectors
status.storage.topic	Topic in which to store connector & task status

As with Consumer Groups, the group.id will determine which Workers will cooperate as a team.

Connect Worker Groups will use unique administrative topics to hold connector configurations, offsets for both source and sinks, and status (if necessary) of the external data sources/sinks. Different worker groups can use different administrative topics if you need to keep their configurations separate.



Make sure that all Workers with the same group.id are using the same names for the administrative topics or you will get unpredictable results.

конфиг-топик создается автоматически но
я могу вручную создать конфиг-топик для того чтобы явно задать replication-factor, partitions-count

- используется компактед-топик для хранения настроек

Creating Kafka Connect's Topics Manually

- The Topics used in Kafka Connect distributed mode will be automatically created
- They are created with recommended values for
 - replication factor
 - partition counts
 - cleanup policy

The following number of Partitions:

Topic	Partitions
config.storage.topic	1
offset.storage.topic	25
status.storage.topic	5

Notice that the recommendations for offset.storage.topic are similar to the __consumer_offsets Topic since they do similar jobs.



In general, it's usually ok to let Connect create these topics for us. Connect will create these topics with the recommended replication factor, partition counts, and cleanup policy. The only time you might want to create these topics manually would be when you need custom configuration, or Connect doesn't have appropriate ACL privileges to create them. See <https://docs.confluent.io/current/connect/userguide.html#distributed-worker-configuration>

Configuring the Connector

Parameter	Description
name	Connector's unique name
connector.class	Connector's Java class
tasks.max	Maximum tasks to create. The Connector may create fewer if it cannot achieve this level of parallelism (Default: 1)
key.converter	(optional) Override the worker key converter
value.converter	(optional) Override the worker value converter
topics (Sink connectors only)	List of input Topics (to consume from)

The number of tasks that a Connector should start will be dependent on the number of cores

The number of tasks that a Connector should start will be dependent on the number of cores and workers available.

The reason that the `topics` parameter only exists for sink connectors is that source connectors typically have logic built in that will create new topic names based on a specific naming convention.

Key and value converters were added to the connector configurations in AK 0.10.1. Prior to that, there was a single converter setting that was specified in the worker configurations. Additional converters can be installed via Confluent Hub (when available).

Config Definitions



- Config Def are the way to communicate to the user how you want your Configuration to be
- You can have mandatory or optional parameters (with default), and validate the types and constraints of your parameters. You should also document the parameters

```
public static final String OWNER_CONFIG = "github.owner";
private static final String OWNER_DOC = "Owner of the repository you'd like to follow";

public static ConfigDef config() {
    return new ConfigDef()
        .define(OWNER_CONFIG, Type.STRING, Importance.HIGH, OWNER_DOC)
        .define(BATCH_SIZE_CONFIG, Type.INT, 100, new BatchSizeValidator(), Importance.LOW, BATCH_SIZE_DOC)
```

```
GitHubSourceConnectorConfig TOPIC_DOC
1 import ...
2
3 public class GitHubSourceConnectorConfig extends AbstractConfig {
4
5     public static final String TOPIC_CONFIG = "topic";
6     private static final String TOPIC_DOC = "Topic to write to";
7
8     public static final String OWNER_CONFIG = "github.owner";
9     private static final String OWNER_DOC = "Owner of the repository you'd like to follow";
10
11    public static final String REPO_CONFIG = "github.repo";
12    private static final String REPO_DOC = "Repository you'd like to follow";
13
14    public static final String SINCE_CONFIG = "since.timestamp";
15    private static final String SINCE_DOC =
16        "Only issues updated at or after this time are returned.\n"
17        + "This is a timestamp in ISO 8601 format: YYYY-MM-DDTHH:MM:SSZ.\n"
18        + "Defaults to a year from first launch.";
19
20    public static final String BATCH_SIZE_CONFIG = "batch.size";
21    private static final String BATCH_SIZE_DOC = "Number of data points to retrieve at a time";
22    Defaults to 100 (max value);
23
24    public static final String AUTH_USERNAME_CONFIG = "auth.username";
25    private static final String AUTH_USERNAME_DOC = "Optional Username to authenticate calls";
26
27    public static final String AUTH_PASSWORD_CONFIG = "auth.password";
28
29    public GitHubSourceConnectorConfig(Map<String, String> parsedConfig) { this(config(),
30        parsedConfig); }
31
32    public static ConfigDef config() {
33        return new ConfigDef()
34            .define(TOPIC_CONFIG, Type.STRING, Importance.HIGH, TOPIC_DOC)
35            .define(OWNER_CONFIG, Type.STRING, Importance.HIGH, OWNER_DOC)
36            .define(REPO_CONFIG, Type.STRING, Importance.HIGH, REPO_DOC)
```

```
46     public GitHubSourceConnectorConfig(Map<String, String> parsedConfig) { this(conf(),  
47         parsedConfig); }  
48  
49     public static ConfigDef conf() {  
50         return new ConfigDef()  
51             .define(TOPIC_CONFIG, Type.STRING, Importance.HIGH, TOPIC_DOC)  
52             .define(OWNER_CONFIG, Type.STRING, Importance.HIGH, OWNER_DOC)  
53             .define(REPO_CONFIG, Type.STRING, Importance.HIGH, REPO_DOC)  
54             .define(BATCH_SIZE_CONFIG, Type.INT, defaultValue: 100, new BatchSizeValidator(),  
55                 Importance.LOW, BATCH_SIZE_DOC)  
56             .define(SINCE_CONFIG, Type.STRING, ZonedDateTime.now().minusYears(1).toInstant()  
57                 .toString(),  
58                 new TimestampValidator(), Importance.HIGH, SINCE_DOC)  
59             .define(AUTH_USERNAME_CONFIG, Type.STRING, defaultValue: "", Importance.HIGH,  
60                 AUTH_USERNAME_DOC)  
61             .define(AUTH_PASSWORD_CONFIG, Type.PASSWORD, defaultValue: "", Importance.HIGH,
```

```
1 name=GitHubSourceConnectorDemo  
2 tasks.max=1  
3 connector.class=com.simplesteph.kafka.GitHubSourceConnector  
4 topic=github-issues  
5 github.owner=kubernetes  
6 github.repo=kubernetes  
7 since.timestamp=2017-01-01T00:00:00Z  
8 # I heavily recommend you set those two fields:  
9 # auth.username=your_username  
0 # auth.password=your_password
```



Schemas

- Schema is an abstraction that allow you to define how your data structure will look like. It is using primitive types, and then the Kafka Connect framework will convert it into JSON or Avro as needed.
- It is necessary to correctly design your schema before you program

```
public static Schema USER_SCHEMA =
SchemaBuilder.struct().name(SCHEMA_VALUE_USER)
.version(1)
.field(USER_URL_FIELD, Schema.STRING_SCHEMA)
.field(USER_ID_FIELD, Schema.INT32_SCHEMA)
.field(USER_LOGIN_FIELD, Schema.STRING_SCHEMA)
.build();
```

```
GitHubSchemas
1 package com.simplesteph.kafka;
2
3 import ...
4
5 public class GitHubSchemas {
6
7     public static String NEXT_PAGE_FIELD = "next_page";
8
9     // Issue fields
10    public static String FULL_REPO_FIELD = "owner/repository";
11    public static String OWNER_FIELD = "owner";
12    public static String REPOSITORY_FIELD = "repository";
13    public static String CREATED_AT_FIELD = "created_at";
14    public static String UPDATED_AT_FIELD = "updated_at";
15    public static String NUMBER_FIELD = "number";
16    public static String URL_FIELD = "url";
17    public static String HTML_URL_FIELD = "html_url";
18    public static String TITLE_FIELD = "title";
19    public static String STATE_FIELD = "state";
20
21    // User fields
22    public static String USER_FIELD = "user";
23    public static String USER_URL_FIELD = "url";
24    public static String USER_HTML_URL_FIELD = "html_url";
25    public static String USER_ID_FIELD = "id";
26    public static String USER_LOGIN_FIELD = "login";
27
28    // PR fields
29    public static String PR_FIELD = "pull_request";
30}
```

```

47
48     // Value Schema
49     public static Schema USER_SCHEMA = SchemaBuilder.struct().name(SCHEMA_VALUE_USER)
50         .version(1)
51         .field(USER_URL_FIELD, Schema.STRING_SCHEMA)
52         .field(USER_ID_FIELD, Schema.INT32_SCHEMA)
53         .field(USER_LOGIN_FIELD, Schema.STRING_SCHEMA)
54         .build();
55
56     // optional schema
57     public static Schema PR_SCHEMA = SchemaBuilder.struct().name(SCHEMA_VALUE_PR)
58         .version(1)
59         .field(PR_URL_FIELD, Schema.STRING_SCHEMA)
60         .field(PR_HTML_URL_FIELD, Schema.STRING_SCHEMA)
61         .optional()
62         .build();
63
64     public static Schema VALUE_SCHEMA = SchemaBuilder.struct().name(SCHEMA_VALUE_ISSUE)
65         .version(1)
66         .field(URL_FIELD, Schema.STRING_SCHEMA)
67         .field(TITLE_FIELD, Schema.STRING_SCHEMA)
68         .field(CREATED_AT_FIELD, Schema.INT64_SCHEMA)
69         .field(UPDATED_AT_FIELD, Schema.INT64_SCHEMA)
70         .field(NUMBER_FIELD, Schema.INT32_SCHEMA)
71         .field(STATE_FIELD, Schema.STRING_SCHEMA)
72         .field(USER_FIELD, USER_SCHEMA) // mandatory
73         .field(PR_FIELD, PR_SCHEMA) // optional
74         .build();
75     }

```

Issues | GitHub Developer Guide kubernetes/kubernetes: Production https://api.github.com/repos/... jcustenborder/kafka-connect: ...

<https://developer.github.com/v3/issues/#list-issues-for-a-repository>

Secure https://developer.github.com/v3/issues/#list-issues-for-a-repository

desc	
since	string
Only issues updated at or after this time are returned. This is a timestamp in ISO 8601 format: YYYY-MM-DDTHH:MM:SSZ .	

Response

```

Status: 200 OK
Link: <https://api.github.com/resource?page=2>; rel="next",
      <https://api.github.com/resource?page=5>; rel="last"
[ {
  "id": 1,
  "url": "https://api.github.com/repos/octocat>Hello-World/issues/1347",
  "repository_url": "https://api.github.com/repos/octocat>Hello-World",
  "labels_url": "https://api.github.com/repos/octocat>Hello-World/issues/1347/labels",
  "comments_url": "https://api.github.com/repos/octocat>Hello-World/issues/1347/comments",
  "events_url": "https://api.github.com/repos/octocat>Hello-World/issues/1347/events",
  "html_url": "https://github.com/octocat>Hello-World/issues/1347",
  "number": 1347,
  "state": "open",
  "title": "Found a bug",
  "body": "I'm having a problem with this.",
  "user": {
    "login": "octocat",
    "id": 1,
    "avatar_url": "https://github.com/images/error/octocat_happy.gif",
    "gravatar_id": "",
    "url": "https://api.github.com/users/octocat",
    "html_url": "https://github.com/octocat",
    "followers_url": "https://api.github.com/users/octocat/followers",
    "following_url": "https://api.github.com/users/octocat/following{/other_user}",
    "gists_url": "https://api.github.com/users/octocat/gists{/gist_id}",
    "starred_url": "https://api.github.com/users/octocat/starred{/owner}{/repo}",
    "subscriptions_url": "https://api.github.com/users/octocat/subscriptions",
    "organizations_url": "https://api.github.com/users/octocat/orgs",
    "repos_url": "https://api.github.com/users/octocat/repos",
    "events_url": "https://api.github.com/users/octocat/events{/privacy}",
    "received_events_url": "https://api.github.com/users/octocat/received_events",
    "type": "User",
    "site_admin": false
  },
  "labels": [
  ]
}
]

```

Data Model



- It is always really good to have your data exist in classes in a “model” package.
- I modelled Issue, User and Pull Request as POJOs

```
public class PullRequest {  
  
    private String url;  
    private String htmlUrl;  
}
```

The screenshot shows an IDE interface with two main panes. On the left is the 'Project' view, which lists several Java files under the 'java/com.simplesteph.kafka/model' directory. The 'User.java' file is currently selected. On the right is the 'User' code editor pane, displaying the following Java code:

```
11  public class User {  
12      private String login;  
13      private Integer id;  
14      private String avatarUrl;  
15      private String gravatarId;  
16      private String url;  
17      private String htmlUrl;  
18      private String followersUrl;  
19      private String followingUrl;  
20      private String gistsUrl;  
21      private String starredUrl;  
22      private String subscriptionsUrl;  
23      private String organizationsUrl;  
24      private String reposUrl;  
25      private String eventsUrl;  
26      private String receivedEventsUrl;  
27      private String type;  
28      private Boolean siteAdmin;  
29      private Map<String, Object>  
30          additionalProperties = new HashMap<~>();  
31  
32      /**  
33      * No args constructor for use in  
34      * serialization  
35      */  
36      public User() {  
37  
38          super();  
39      }  
40  
41      public static User fromJson(JSONObject jsonObject) {  
42          User user = new User();  
43          user.setUrl(jsonObject.getString(USER_URL_FIELD));  
44          user.setHtmlUrl(jsonObject.getString(USER_HTML_URL_FIELD));  
45          user.setId(jsonObject.getInt(USER_ID_FIELD));  
46          user.setLogin(jsonObject.getString(USER_LOGIN_FIELD));  
47          return user;  
48      }  
49  }
```


это HE partition и offset в кафке

Source partition, Source Offsets



- Source Partition allows Kafka Connect to know which source you've been reading
- Source Offsets allow Kafka Connect to Track until when you've been reading for the Source Partition you chose
- **They are different than partition and offsets for Kafka.**
- Source Partition and Source Offsets are for Kafka Connect Source

```
private Map<String, String> sourcePartition() {  
    Map<String, String> map = new HashMap<>();  
    map.put(OWNER_FIELD, config.getOwnerConfig());  
    map.put(REPOSITORY_FIELD, config.getRepoConfig());  
    return map;  
}
```

```

111     @Override
112     public void stop() {
113         // Do whatever is required to stop your task.
114     }
115
116     private Map<String, String> sourcePartition() {
117         Map<String, String> map = new HashMap<>();
118         map.put(OWNER_FIELD, config.getOwnerConfig());
119         map.put(REPOSITORY_FIELD, config.getRepoConfig());
120         return map;
121     }
122
123     private Map<String, String> sourceOffset(
124         Instant updatedAt) {
125         Map<String, String> map = new HashMap<>();
126         map.put(UPDATED_AT_FIELD, DateUtils
127             .MaxInstant(updatedAt, nextQuerySince).toString());
128         map.put(NEXT_PAGE_FIELD, nextPageToVisit
129             .toString());
130         return map;
131     }
132
133     private Struct buildRecordKey(Issue issue){
134         // Key Schema

```

при каждом чтении из источника я буду сдвигать source offset

```

private Map<String, String> sourcePartition() {
    Map<String, String> map = new HashMap<>();
    map.put(OWNER_FIELD, config.getOwnerConfig());
    map.put(REPOSITORY_FIELD, config.getRepoConfig());
    return map;
}

private Map<String, String> sourceOffset(Instant updatedAt) {
    Map<String, String> map = new HashMap<>();
    map.put(UPDATED_AT_FIELD, DateUtils.MaxInstant(updatedAt, nextQuerySince).toString());
    map.put(NEXT_PAGE_FIELD, nextPageToVisit.toString());
    return map;
}

private Struct buildRecordKey(Issue issue){
    // Key Schema
    Struct key = new Struct(KEY_SCHEMA)
        .put(OWNER_FIELD, config.getOwnerConfig())
        .put(REPOSITORY_FIELD, config.getRepoConfig())
        .put(NUMBER_FIELD, issue.getNumber());
}

```

```

GitHubSourceTask generateSourceRecord()
94     }
95     return records;
96 }
97
98 private SourceRecord generateSourceRecord(Issue issue) {
99     return new SourceRecord(
100         sourcePartition(),
101         sourceOffset(issue.getUpdatedAt()),
102         config.getTopic(),
103         partition: null, // partition will be inferred by the framework
104         KEY_SCHEMA,
105         buildRecordKey(issue),
106         VALUE_SCHEMA,
107         buildRecordValue(issue),
108         issue.getUpdatedAt().toEpochMilli());
109 }
110
111 @Override
112 ...

```

Find Usages of sourcePartition() in Project and Libraries

- Unclassified usage 2 usages
 - kafkaconnectgithubsource 2 usages
 - com.simplesteph.kafka 2 usages
 - GitHubSourceTask 2 usages
 - generateSourceRecord(Issue) 1 usage
 `100 sourcePartition(),`
 - initializeLastVariables() 1 usage

при инициализации таски она будет начинать с последнего source offset

```

GitHubSourceTask initializeLastVariables()
    @Override
    public String version() { return VersionUtil.getVersion(); }

    @Override
    public void start(Map<String, String> map) {
        //Do things here that are required to start your task. This could be open a connection
        to a database, etc.
        config = new GitHubSourceConnectorConfig(map);
        initializeLastVariables();
        gitHubHttpAPIClient = new GitHubAPIHttpClient(config);
    }

    private void initializeLastVariables(){
        Map<String, Object> lastSourceOffset = null;
        lastSourceOffset = context.offsetStorageReader().offset(sourcePartition());
        if( lastSourceOffset == null){
            // we haven't fetched anything yet, so we initialize to 7 days ago
            nextQuerySince = config.getSince();
            lastIssueNumber = -1;
        } else {
            Object updatedAt = lastSourceOffset.get(UPDATED_AT_FIELD);
            Object issueNumber = lastSourceOffset.get(NUMBER_FIELD);
            Object nextPage = lastSourceOffset.get(NEXT_PAGE_FIELD);
            if(updatedAt != null && (updatedAt instanceof String)){
                nextQuerySince = Instant.parse((String) updatedAt);
            }
            if(issueNumber != null && (issueNumber instanceof String)){
                lastIssueNumber = Integer.valueOf((String) issueNumber);
            }
        }
    }
}

```

```
GitHubSourceTask poll()
69     @Override
70     public List<SourceRecord> poll() throws InterruptedException {
71         gitHubHttpClient.sleepIfNeed();
72
73         // fetch data
74         final ArrayList<SourceRecord> records = new ArrayList<>();
75         JSONArray issues = gitHubHttpClient.getNextIssues(nextPageToVisit, nextQuerySince);
76         // we'll count how many results we get with i
77         int i = 0;
78         for (Object obj : issues) {
79             Issue issue = Issue.fromJson((JSONObject) obj);
80             SourceRecord sourceRecord = generateSourceRecord(issue);
81             records.add(sourceRecord);
82             i += 1;
83             lastUpdatedAt = issue.getUpdatedAt();
84         }
85         if (i > 0) log.info(String.format("Fetched %s record(s)", i));
86         if (i == 100){
87             // we have reached a full batch, we need to get the next one
88             nextPageToVisit += 1;
89         }
90         else {
91             nextQuerySince = lastUpdatedAt.plusSeconds(1);
92             nextPageToVisit = 1;
93             gitHubHttpClient.sleep();
94         }
95     }
96 }
```

```
GitHubSourceTask generateSourceRecord()
88     if(nextPageToVisit != 1,
89     }
90     else {
91         nextQuerySince = lastUpdatedAt.plusSeconds(1);
92         nextPageToVisit = 1;
93         gitHubHttpClient.sleep();
94     }
95     return records;
96 }

98     private SourceRecord generateSourceRecord(Issue issue) {
99         return new SourceRecord(
100             sourcePartition(),
101             sourceOffset(issue.getUpdatedAt()),
102             config.getTopic(),
103             partition: null, // partition will be inferred by the framework
104             KEY_SCHEMA,
105             buildRecordKey(issue),
106             VALUE_SCHEMA,
107             buildRecordValue(issue),
108             issue.getUpdatedAt().toEpochMilli());
109     }

111     @Override
112     public void stop() {
113         // Do whatever is required to stop your task.
114     }

116     private Map<String, String> sourcePartition() {
117 }
```

Task

- Our Task is what does the actual job
- It's supposed to initialize, then find where to resume from, and finally poll the source for records

```
public class GitHubSourceTask extends SourceTask {  
    public String version()  
    public void start(Map<String, String> map)  
    public List<SourceRecord> poll() throws  
        InterruptedException  
    public void stop()  
}
```

Example: FileStreamSourceConnector STANDALONE MODE

- Goal:

- Read a file and load the content directly into Kafka
- Run in a connector in **standalone mode** (useful for development)



- Learning:

- Understand how to configure a connector in standalone mode
- Get a first feel for Kafka Connect Standalone

```
# we launch the kafka connector in standalone mode:
```

```
cd /tutorial/source/demo-1
```

```
root@fast-data-dev / $ cd tutorial/source/demo-1/
```

```
root@fast-data-dev demo-1 $ ls
```

```
demo-file.txt  file-stream-demo-standalone.properties  worker.properties
```

```
# create the topic we write to with 3 partitions
```

```
kafka-topics --create --topic demo-1-standalone --partitions 3 --replication-factor 1 --zookeeper 127.0.0.1:2181
```

```
# Usage is connect-standalone worker.properties connector1.properties [connector2.properties connector3.properties]  
connect-standalone worker.properties file-stream-demo-standalone.properties
```

```
# write some data to the demo-file.txt !
```

```
# shut down the terminal when you're done.
```

```
#####
```

worker.properties -worker infrastructure config for standalone mode

```
# from more information, visit: http://docs.confluent.io/3.2.0/connect/userguide.html#common-worker-configs
```

```
bootstrap.servers=127.0.0.1:9092
```

```
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
```

```
# we always leave the internal key to JsonConverter
```

```
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
```

file-stream-demo-standalone.properties file connector config

```
# These are standard kafka connect parameters, need for ALL connectors
name=file-stream-demo-standalone
connector.class=org.apache.kafka.connect.file.FileStreamSourceConnector
tasks.max=1 # число задач которые можно запускать паралельно

# Parameters can be found here:
https://github.com/apache/kafka/blob/trunk/connect/file/src/main/java/org/apache/kafka/connect/file/FileStreamSourceConnector.java
file=demo-file.txt # файл откуда будем читать
topic=demo-1-standalone # топик в который будем писать
```

при реалтайм добавлении строки в файл она будет автоматически подхватываться и передаваться в ТОПИК

Example: FileStreamSourceConnector DISTRIBUTED MODE

- Goal:

- Read a file and load the content directly into Kafka
- Run in **distributed mode** on our already set-up Kafka Connect Cluster



- Learning:

- Understand how to configure a connector in distributed mode
- Get a first feel for Kafka Connect Cluster
- Understand the schema configuration option

```
#####
# B) FileStreamSourceConnector in distributed mode:
# create the topic we're going to write to
docker run --rm -it --net=host landoop/fast-data-dev bash
kafka-topics --create --topic demo-2-distributed --partitions 3 --replication-factor 1 --zookeeper 127.0.0.1:2181
# you can now close the new shell

# head over to 127.0.0.1:3030 -> Connect UI
# Create a new connector -> File Source
# Paste the configuration at source/demo-2/file-stream-demo-distributed.properties

# Now that the configuration is launched, we need to create the file demo-file.txt
docker ps
docker exec -it <containerId> bash
touch demo-file.txt
echo "hi" >> demo-file.txt
echo "hello" >> demo-file.txt
echo "from the other side" >> demo-file.txt
```

законектился на докер и создадим файл

```
20:53 $ docker ps
CONTAINER ID        IMAGE               COMMAND             CREATED            STATUS              PORTS
NAMES
cclb88f71d31      landoop/fast-data-dev   "/usr/local/bin/du..."   11 minutes ago    Up 11 minute
$ 0.0.0.0:2181->2181/tcp, 0.0.0.0:3030->3030/tcp, 0.0.0.0:8081-8083->8081-8083/tcp, 0.0.0.0:9092->9092/tcp, 3031/tcp temp-2
✓ ~

20:53 $ docker exec -it cclb88f71d31 bash
root@fast-data-dev / $ touch demo-file.txt
root@fast-data-dev / $ ls
bin connectors dev extra-connect-jars lib mnt proc run srv tmp var
build.info demo-file.txt etc home media opt root sbin sys usr
```

```
# Read the topic data
docker run --rm -it --net=host landoop/fast-data-dev bash
kafka-console-consumer --topic demo-2-distributed --from-beginning --bootstrap-server 127.0.0.1:9092
# observe we now have json as an output, even though the input was text!
#####
```

как создать distributed connector через UI

Kafka Development Environment
docker container powered by Landoop

SCHEMAS	TOPICS	CONNECTORS	BROKERS
5 SCHEMA REGISTRY UI manage avro schemas ENTER	12 KAFKA TOPICS UI browse topics and data ENTER	5 KAFKA CONNECT UI setup, manage connectors ENTER	1 KAFKA BROKERS management and monitoring COMING SOON
5 KAFKA CONNECT UI setup, manage connectors ENTER	1 KAFKA BROKERS management and monitoring COMING SOON		

RUNNING SERVICES
fast-data-dev

KAFKA CONNECT

5 Connectors NEW

Search connectors

- logs-broker
- logs-connect-distributed
- logs-rest-proxy
- logs-schema-registry
- logs-zookeeper

Kafka Connect : /api/kafka-connect
Kafka Connect Version : 0.10.1.1
Kafka Connect UI Version : 0.9.0

New Connector

Search

Sources

- twitter**
Subscribe to feeds using the Twitter API and stream data into a kafka topic
- ftp**
Tail directories in remote FTP location and bring messages in Kafka
- file**
Read files and stream data into Kafka Topics
- bloomberg**
A connector to subscribe to Bloomberg feeds via the Bloomberg Ibis open API and write data to Kafka

connector properties SHOW REQUIRED FIELDS

Working task

```
# These are standard kafka connect parameters, need for ALL connectors
name=file-stream-demo-distributed
connector.class=org.apache.kafka.connect.file.FileStreamSourceConnector
tasks.max=1
```

Parameters can be found here:

```
https://github.com/apache/kafka/blob/trunk/connect/file/src/main/java/org/apache/kafka/connect/file/FileStreamSourceConnector.java
file=demo-file.txt # файл откуда будем читать
topic=demo-2-distributed # топик в который будем писать
```

```
# Added configuration for the distributed mode: так как здесь нету worker.properties файла то все настройки worker прям здесь
key.converter=org.apache.kafka.connect.json.JsonConverter
key.converter.schemas.enable=true
value.converter=org.apache.kafka.connect.json.JsonConverter
value.converter.schemas.enable=true
```

используем landoop UI

version: '2'

services:

this is our kafka cluster.

```

kafka-cluster:
image: landoop/fast-data-dev:cp3.3.0
environment:
  ADV_HOST: 127.0.0.1      # Change to 192.168.99.100 if using Docker Toolbox
  RUNTESTS: 0              # Disable Running tests so the cluster starts faster
ports:
  - 2181:2181              # Zookeeper
  - 3030:3030              # Landoop UI
  - 8081-8083:8081-8083   # REST Proxy, Schema Registry, Kafka Connect ports
  - 9581-9585:9581-9585   # JMX Ports
  - 9092:9092              # Kafka Broker

# we will use elasticsearch as one of our sinks.
# This configuration allows you to start elasticsearch
elasticsearch:
image: itzg/elasticsearch:2.4.3
environment:
  PLUGINS: appbaseio/dejavu
  OPTS: -Dindex.number_of_shards=1 -Dindex.number_of_replicas=0
ports:
  - "9200:9200"

# we will use postgres as one of our sinks.
# This configuration allows you to start postgres
postgres:
image: postgres:9.5-alpine
environment:
  POSTGRES_USER: postgres  # define credentials
  POSTGRES_PASSWORD: postgres # define credentials
  POSTGRES_DB: postgres    # define database
ports:
  - 5432:5432              # Postgres port

```

* local file connector (only for dev)

14 декабря 2020 г. 13:33

коннектор отслеживает только новопоявившиеся строки в конце файла и не может определить изменения уже закаченных строк

FileStream Connector

можно накрутить этот коннектор на линке, куда будут поступать внешние данные

- FileStream Connector acts on a local file
 - Local file Source Connector: tails local file and sends each line as a Kafka message
 - Local file Sink Connector: Appends Kafka messages to a local file

The FileSource Connector reads data from a file and sends it to Kafka. This connector will read only one file and send the data within that file to Kafka. It will then watch the file for appended updates only. Any modification of file lines already sent to Kafka will not be reprocessed.

The FileSink Connector reads data from Kafka and outputs it to a local file. Multiple topics may be specified as with any other sink connector. As messages are added to the topics specified in the configuration, they are produced to a local file as specified in the configuration.

The FileStream Connector examples on docs.confluent.io are intended to show how a simple connector runs for those first getting started with Kafka Connect as either a user or developer. It is not recommended for production use. Instead, we encourage users to use them to learn in a local environment. The examples include both a file source and a file sink to demonstrate an end-to-end data flow implemented through Kafka Connect. The FileStream Connector examples are also detailed in the developer guide as a demonstration of how a custom connector can be implemented.

* twitter source connector (distributed)

4 декабря 2020 г. 18:33

TwitterSourceConnector



• Goal:

- Gather data from Twitter in Kafka Connect Distributed mode



• Learning:

- Gather real data using <https://github.com/Eneco/kafka-connect-twitter>

Kafka Connect Sink/Source for Twitter

Branch: master	New pull request	Find file	Clone or download
82 commits	1 branch	2 releases	8 contributors
Latest commit af63e4c on 13 Jan			
src	Fix previous erroneous commit	6 months ago	
.gitignore	simplify pom.xml & make it working with eclipse	a year ago	
.travis.yml	Remove ulimit	3 months ago	
Dockerfile	Update base for jmx	3 months ago	
LICENSE	Initial commit	a year ago	
README.md	Update README.md	6 months ago	
connect-simple-source-standalone.properties	Provide example properties; fix README	a year ago	
connect-sink-standalone.properties	Sink (#20)	a year ago	
connect-source-standalone.properties	Sink (#20)	a year ago	

для того чтобы использовать API нужен twitter DEVELOPER account

Twitter Apps

Please [sign in](#) with your Twitter Account to create and maintain Twitter Apps.



Twitter Apps

Create New App



Udemy-Kafka-Connect-Demo

Feed tweets into Kafka in real time!



Your application has been created. Please take a moment to review and adjust your application's settings.

Kafka-Connect-Demo

Test OAuth



Source data from Twitter to put it directly into Kafka, using the Kafka Connect framework

http://udemy.com/



Organization

Information about the organization or company associated with your application. This information is optional.

Organization None

Organization website None

Application Settings

Your application's Consumer Key and Secret are used to **authenticate** requests to the Twitter Platform.

Access level Read and write (modify app permissions)

Consumer Key (API Key) 9jT7rYU09ELMSibxgVK7tRpRs (manage keys and access tokens)

Callback URL None

Callback URL Locked No

Keys and Access Tokens

Token Actions

Create my access token



Application Settings

Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.

Consumer Key (API Key) 9jT7rYU09ELMSibxgVK7tRpRs

Consumer Secret (API Secret) ggOg21pZLFcexZ0PjTX1zRtG0bHHuWFhMx4DxGfRAVD2xwrHR4

Access Level Read and write (modify app permissions)



Application Actions

Regenerate Consumer Key and Secret

Change App Permissions

Application Settings

Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.

Consumer Key (API Key) 9jT7rYU09ELMSibxgVK7tRpRs

Consumer Secret (API Secret) ggQg21pZLFcexZ0PjTX1zRtG0bHHuWFhMx4DxGfRAVD2xwrHR4

Access Level Read and write ([modify app permissions](#))



Application Actions

[Regenerate Consumer Key and Secret](#)

[Change App Permissions](#)

Your Access Token

This access token can be used to make API requests on your own account's behalf. Do not share your access token secret with anyone.

Access Token 2932095118-245ltf5j2H01MDrRq0MMPNtPHTlxER4e3Cj2hoD

Access Token Secret 7kAn2qXA7wJJ8k46ySj1Iri9Reey9h2NLAxLPtFUkjVgb

Access Level Read and write

как создать distributed connector через UI

The screenshot shows the Kafka Connect UI interface. On the left, there is a list of existing connectors: 'file-stream-demo-distributed', 'logs-broker', and 'logs-connect-distributed'. On the right, a 'New Connector' form is open, specifically for a 'twitter' source connector. The 'Sources' section shows a single entry for 'twitter', which is described as 'Subscribe to feeds using the Twitter API and stream data into a Kafka topic'. A large blue 'NEXT' button is at the bottom of the form.

#####

C) TwitterSourceConnector in distributed mode:

create the topic we're going to write to

docker run --rm -it --net=host landoop/fast-data-dev bash

kafka-topics --create --topic demo-3-twitter --partitions 3 --replication-factor 1 --zookeeper 127.0.0.1:2181

Start a console consumer on that topic

kafka-console-consumer --topic demo-3-twitter --bootstrap-server 127.0.0.1:9092

Follow the instructions at: <https://github.com/Eneco/kafka-connect-twitter#creating-a-twitter-application>

To obtain the required keys, visit <https://apps.twitter.com/> and Create a New App. Fill in an application name & description & web site and accept the developer agreement. Click on Create my access token and populate a file twitter-source.properties with consumer key & secret and the access token & token secret using the example file to begin with.

Setup instructions for the connector are at: <https://github.com/Eneco/kafka-connect-twitter#setup>

```
# fill in the required information at demo-3/source-twitter-distributed.properties  
# Launch the connector and start seeing the data flowing in!
```

The screenshot shows the Kafka Connect UI interface. On the left, there's a list of existing connectors: file-stream-demo-distributed, logs-broker, logs-connect-distributed, logs-rest-proxy, logs-schema-registry, and logs-zookeeper. Each item has a green circular icon, a checkmark, a '1 x' label, a 'file' button, and a '8s' button. On the right, a modal window titled 'New Connector (Source): TwitterSourceConnector' is open. It displays the configuration for the 'TwitterSourceConnector' class: com.eneco.trading.kafka.connect.twitter.TwitterSourceConnector. Below this, there's a checkbox labeled 'Show cURL command' which is unchecked. The configuration code is listed below:

```
1 name=source-twitter-distributed  
2 connector.class=com.eneco.trading.kafka.connect.twitter.TwitterSourceConnector  
3 tasks.max=1  
4 topic=demo-3-twitter  
5 key.converter=org.apache.kafka.connect.json.JsonConverter  
6 key.converter.schemas.enable=true  
7 value.converter=org.apache.kafka.connect.json.JsonConverter  
8 value.converter.schemas.enable=true  
9 twitter.consumerkey=9jT7rYU09ELMSibxgVK7tRpRs  
10 twitter.consumersecret=ggQg21pZLFcexZ0PjTX1zRtG0bHHuWFhMx4DxGfRAVD2xwrHR4  
11 twitter.token=2932095118-2451tf5j2H01MDrRq0MMPNtPHTIxER4e3Cj2hoD  
12 twitter.secret=7kAn2qXA7wJJ8k46ySj1lr9Reey9h2NLAxLptFUKjVgb  
13 track.terms=programming,java,kafka,scala  
14 language=en
```

```
# Basic configuration for our connector  
name=source-twitter-distributed  
connector.class=com.eneco.trading.kafka.connect.twitter.TwitterSourceConnector  
tasks.max=1  
topic=demo-3-twitter  
key.converter=org.apache.kafka.connect.json.JsonConverter  
key.converter.schemas.enable=true  
value.converter=org.apache.kafka.connect.json.JsonConverter  
value.converter.schemas.enable=true  
  
# Twitter connector specific configuration Настройки с twitter API  
twitter.consumerkey=9jT7rYU09ELMSibxgVK7tRpRs  
twitter.consumersecret=ggQg21pZLFcexZ0PjTX1zRtG0bHHuWFhMx4DxGfRAVD2xwrHR4  
twitter.token=2932095118-2451tf5j2H01MDrRq0MMPNtPHTIxER4e3Cj2hoD  
twitter.secret=7kAn2qXA7wJJ8k46ySj1lr9Reey9h2NLAxLptFUKjVgb  
  
track.terms=programming,java,kafka,scala  
language=en
```

* Elastic sink connector (distributed)

4 декабря 2020 г. 18:33

ElasticSearch



- We're going to start an ElasticSearch instance to Sink the data to
- We will use Docker to start our ElasticSearch instance
- This will serve as our Sink for our first Sink Connector

- ElasticSearch is an easy way to store json data and search across it
- Used by many companies around the world to power search engines and advanced analytics capabilities
- No knowledge of ElasticSearch is required for this course

- Goal:

- Start an ElasticSearch instance using Docker
- Sink a topic with multiple partitions to ElasticSearch
- Run in distributed mode with multiple tasks



- Learning:

- Learn about the `tasks.max` parameter
- Understand how Sink Connectors work

The screenshot shows a web browser displaying the Confluent Kafka Connect Configuration Options for the Confluent Elasticsearch Connector. The URL is docs.confluent.io/3.1.1/connect/connect-elasticsearch/docs/configuration_options.html. The page has a dark header with the Confluent logo and navigation links for Product, Services, Resources, About, Blog, and Download. The main content area is titled "Configuration Options" and includes sections for "connection.url", "type.name", and "key.ignore". Each section provides a description, type information, default values, and importance levels. On the left sidebar, there is a navigation menu with links to various Confluent documentation pages like "What is the Confluent Platform?", "Confluent Platform 3.1.1 Release Notes", etc.

используем landoop UI для управления кафкой

```
version: '2'

services:
  # this is our kafka cluster.
  kafka-cluster:
    image: landoop/fast-data-dev:cp3.3.0
    environment:
      ADV_HOST: 127.0.0.1      # Change to 192.168.99.100 if using Docker Toolbox
      RUNTESTS: 0              # Disable Running tests so the cluster starts faster
    ports:
      - 2181:2181            # Zookeeper
      - 3030:3030            # Landoop UI
      - 8081-8083:8081-8083  # REST Proxy, Schema Registry, Kafka Connect ports
      - 9581-9585:9581-9585  # JMX Ports
      - 9092:9092            # Kafka Broker

  # we will use elasticsearch as one of our sinks.
  # This configuration allows you to start elasticsearch
  elasticsearch:
    image: itzg/elasticsearch:2.4.3
    environment:
      PLUGINS: appbaseio/dejavu
      OPTS: -Dindex.number_of_shards=1 -Dindex.number_of_replicas=0
    ports:
      - "9200:9200"

  # we will use postgres as one of our sinks.
  # This configuration allows you to start postgres
  postgres:
    image: postgres:9.5-alpine
    environment:
      POSTGRES_USER: postgres  # define credentials
      POSTGRES_PASSWORD: postgres # define credentials
      POSTGRES_DB: postgres    # define database
    ports:
      - 5432:5432            # Postgres port
```

запустим докер

```
#!/bin/bash

# Make sure you change the ADV_HOST variable in docker-compose.yml
# if you are using docker Toolbox

# 1) Source connectors
# Start our kafka cluster
docker-compose up kafka-cluster elasticsearch postgres

# Wait 2 minutes for the kafka cluster to be started

#####
# A) ElasticSearch Sink
# Info here: http://docs.confluent.io/3.2.0/connect/connect-elasticsearch/docs/elasticsearch\_connector.html
# We make sure elasticsearch is working. Replace 127.0.0.1 by 192.168.99.100 if needed
http://127.0.0.1:9200/
```



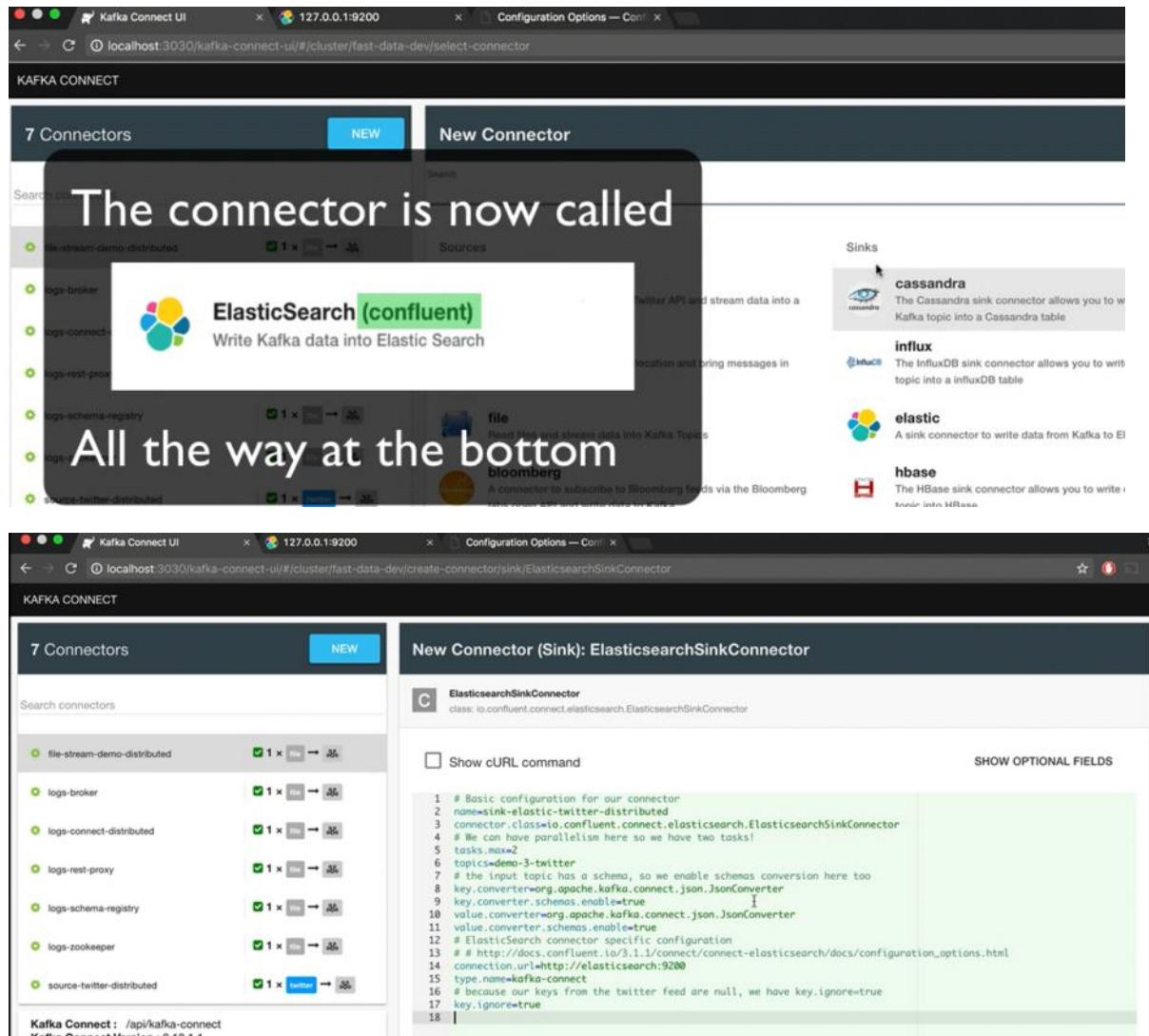
```
# Go to the connect UI and apply the configuration at :  
sink/demo-elasticsearch/twittersink-distributed.properties
```

```
# Visualize the data at:  
http://127.0.0.1:9200/\_plugin/dejavu
```

```
# http://docs.confluent.io/3.1.1/connect/connect-elasticsearch/docs/configuration\_options.html  
# Counting the number of tweets:  
http://127.0.0.1:9200/demo-3-twitter/\_count
```

```
# You can download the data from the UI to see what it looks like  
# We can query elasticsearch for users who have a lot of friends, see query-high-friends.json  
#####
```

настроим коннектор через UI



The screenshot shows the Kafka Connect UI interface. At the top, there's a navigation bar with tabs for 'Configuration Options' and 'Connectors'. Below the navigation, a banner says 'The connector is now called' followed by the connector name 'ElasticSearch (confluent)'. The main area displays the 'Sources' and 'Sinks' sections. Under 'Sources', there are several options like 'file', 'logstash', 'kafka', etc. Under 'Sinks', there are options for 'cassandra', 'influx', 'elastic', and 'hbase'. In the center, there's a large box for creating a new connector, titled 'New Connector (Sink): ElasticsearchSinkConnector'. It contains fields for 'Name' (set to 'ElasticsearchSinkConnector'), 'Class' (set to 'io.confluent.connect.elasticsearch.ElasticsearchSinkConnector'), and a 'Show cURL command' checkbox. Below these, there's a code editor with the basic configuration for the connector.

```
# Basic configuration for our connector  
name=sink-elasticsearch-tweet-distributed
```

```

connector.class=io.confluent.connect.elasticsearch.ElasticsearchSinkConnector
# We can have parallelism here so we have two tasks!
tasks.max=2
topics=demo-3-twitter

# the input topic has a schema, so we enable schemas conversion here too
key.converter=org.apache.kafka.connect.json.JsonConverter
key.converter.schemas.enable=true
value.converter=org.apache.kafka.connect.json.JsonConverter
value.converter.schemas.enable=true

# ElasticSearch connector specific configuration
# # http://docs.confluent.io/3.3.0/connect/connect-elasticsearch/docs/configuration\_options.html
connection.url=http://elasticsearch:9200
type.name=kafka-connect
# because our keys from the twitter feed are null, we have key.ignore=true
key.ignore=true
# some (dummy) settings we need to add to ensure the configuration is accepted (see https://www.udemy.com/apache-kafka-series-kafka-connect-hands-on-learning/learn/v4/questions/2250394). The bug is tracked at https://github.com/Landoop/fast-data-dev/issues/42
topic.index.map="demo-3-twitter:index1"
topic.key.ignore=true
topic.schema.ignore=true

```

посмотрим как работает elasticsearch

The screenshot shows a browser window titled 'Kafka Connect UI' with the URL '127.0.0.1:9200/_plugin/dejavu/'. The main content area is titled 'Déjà vu' and 'The Missing Web UI for Elasticsearch'. Below this is a search bar with the placeholder 'Popname (aka index) goes here' and the value 'demo-3-twitter' selected. A 'Start Browsing' button is visible below the search bar.

The screenshot shows a browser window titled 'Kafka Connect UI' with the URL '127.0.0.1:9200/_plugin/dejavu/#?input_state=XQAAAALHAAAAAAA9ilqny-B2BnTZGEQzY8iwklFBT5R2DHWCav-XotPkazMCCc3CAO70LDBpnI6t8so05hY-POMA3hTgVK2DW04...'. The main content area is titled 'Déjà vu' and 'The Missing Web UI for Elasticsearch'. At the top, there's a navigation bar with 'Query View' (selected), 'demo-3-twitter', 'http://127.0.0.1:9200', and 'Disconnect'. Below this is a table with the following data:

Type	Queries	Reload	Showing 20 of total 262	Types: 1	Disconnect
<input checked="" type="checkbox"/> kafka-connect	type / id	id	created_at	user	text
	kafka-connect / demo-3...	851404744677576700	2017-04-10T12:01:24.000+0000	...	RT @KOYCHEV: Reactive Angular - An Introduction https://t.co/eAX0uPzNaZ #angularjs #javascript #p https://t...
	kafka-connect / demo-3...	851404771646984200	2017-04-10T12:01:30.000+0000	...	Backend Software Engin

The screenshot shows the Kafka Connect UI interface. A modal window titled "Add Query" is open in the center. The "Name" field contains "Query name". The "Type" dropdown is set to "kafka table". The "Query body" field contains the following JSON code:

```
1 ↵ {  
2 ↵   "query":{  
3 ↵     "term" : {  
4 ↵       "is_retweet" : true  
5 ↵     }  
6 ↵   }  
7 ↵ }
```

At the bottom right of the modal is a blue "Add" button. In the background, there's a "Query View" tab with "demo-3-twitter" selected. On the left, there are tabs for "Types", "Queries", and "Shortcuts". Below the tabs, there are buttons for "+ Add Query", "type / id", "kafka-connect", and "kafka-connect". To the right, there's a table with columns "user" and "text", showing some sample data. At the bottom of the screen, there's a footer with the text "kafka-connect / demo-3... 851404801996972000 2017-04-10T12:01:37.000+0000".

* jdbc source connector

4 декабря 2020 г. 22:42

на каждую таблицу БД используется свой топик

- ?? новые таблицы отслеживаются автоматически -> и автоматом создается топик??

JDBC Source Connector: Overview

- JDBC Source periodically polls a relational database for new or recently modified rows
 - Creates a record for each row, and Produces that record as a Kafka message
- Records from each table are Produced to their own Kafka Topic
- New and deleted tables are handled automatically

The JDBC source connector allows you to import data from any relational database with a JDBC driver into Kafka topics. By using JDBC, this connector can support a wide variety of databases without requiring custom code for each one.

Data is loaded by periodically executing a SQL query and creating an output record for each row in the result set. By default, all tables in a database are copied, each to its own output topic. The database is monitored for new or deleted tables and adapts automatically. When copying data from a table, the connector can load only new or modified rows by specifying which columns should be used to detect new or modified data.



Here is a good place to discuss CDC versus query based source connectors (Pros and Cons of each)!

query based CDC connector - в отличие от log-based CDC коннекторов, этот коннектор не может определять deleted rows (а также проблемы с инсертами и апдейтами)

- log-based означает что коннектору надо влезть в кишку/log самой БД
- для определения апдейта только таймстампа недостаточно, тк могут быть записи с одинаковым таймстампом (тогда должна быть колонка с автоинкрементом и если такой нет то проблема, или должна быть возможность изменить схему таблицы источника)
-

Detecting New and Updated Rows

Incremental query mode	Description
Incrementing column определяет только insert не может определить update	Check a single column where newer rows have a larger, auto-incremented ID. Does not support updated rows
Timestamp column	Checks a single 'last modified' column. Can't guarantee reading all updates
Timestamp and incrementing column определяет апдейты	Combination of the two methods above. Guarantees that all updates are read
Custom query	Used in conjunction with the options above for custom filtering



bulk mode for one-time load, not incremental, unfiltered

The Connector can detect new and updated rows in several ways as described on the slide.

For the reasons stated on the slides, many environments will use both the timestamp and the incrementing column to capture all updates.

Note: Because timestamps are not necessarily unique, the timestamp column mode cannot guarantee all updated data will be delivered. If two rows share the same timestamp and are returned by an incremental query, but only one has been processed before a crash, the second update will be missed when the system recovers.

The custom query option can only be used in conjunction with one of the other incremental modes as long as the necessary WHERE clause can be appended to the query. In some cases, the custom query may handle all filtering itself.

В **bulk-mode** мы каждый раз будем запрашивать содержимое таблицы целиком

- апдейты 100% определяются
- годится только для маленьких таблиц
- для миграции: годится для первоначальной-одноразовой загрузки содержимого БД

как настроить

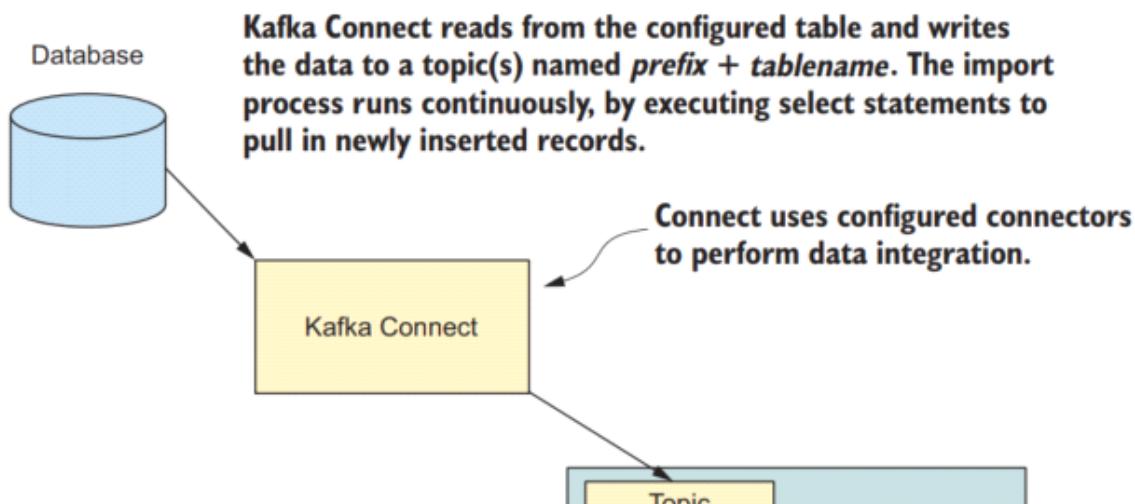
JDBC Source Connector: Configuration

Parameter	Description
<code>connection.url</code>	The JDBC connection URL for the database
<code>topic.prefix</code>	The prefix to prepend to table names to generate the Kafka Topic name
<code>mode</code>	The mode for detecting table changes. Options are <code>bulk</code> , <code>incrementing</code> , <code>timestamp</code> , <code>timestamp+incrementing</code>
<code>query</code>	The custom query to run, if specified
<code>poll.interval.ms</code>	The frequency in milliseconds to poll for new data in each table
<code>table.blacklist</code>	A list of tables to ignore and not import. If specified, <code>tables.whitelist</code> cannot be specified
<code>table.whitelist</code>	A list of tables to import. If specified, <code>tables.blacklist</code> cannot be specified



See <https://docs.confluent.io/current/connect/connect-jdbc/docs/index.html> for a complete list

Setting both `table.whitelist` and `table.blacklist` does not fail any upfront configuration validation checks but will fail when starting the connector at runtime.



The Kafka Streams application now processes the incoming data.

Because the import process is continuous, you are essentially stream processing a database table.

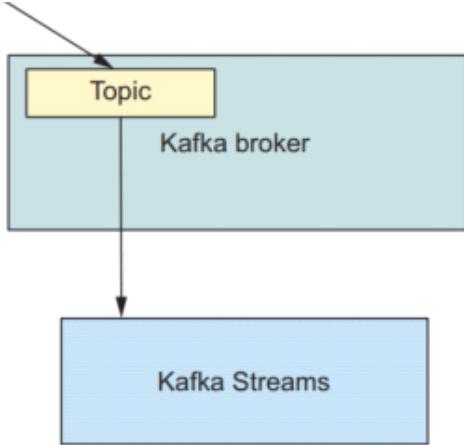


Figure 9.1 Kafka Connect integrating a database table and Kafka Streams

- ❑ `bootstrap.servers` — разделенный запятыми список используемых Connect брокеров Kafka;
- ❑ `key.converter` — класс преобразователя, отвечающего за сериализацию ключей из формата Connect в формат, в котором данные записываются в Kafka;
- ❑ `value.converter` — класс преобразователя, отвечающего за сериализацию значений из формата Connect в формат, в котором данные записываются в Kafka. В нашем примере будет применяться встроенный класс `org.apache.kafka.connect.json.JsonConverter`;
- ❑ `value.converter.schemas.enable` — может быть равен `true` или `false`. Указывает, должен ли Connect включать схему для значения. В этом примере мы установим его равным `false`. В следующем разделе я объясню почему;
- ❑ `plugin.path` — содержит местоположение коннектора для Connect и его зависимостей. Можно указать один каталог верхнего уровня, содержащий один или несколько файлов JAR. Можно также указать несколько путей в виде разделенного запятыми списка;
- ❑ `offset.storage.file.filename` — файл, содержащий сохраненные потребителем Connect смещения.

Вам придется также указать несколько параметров конфигурации для JDBC-коннектора. Вот их список:

- ❑ `name` — название коннектора;
- ❑ `connector.class` — класс коннектора;
- ❑ `tasks.max` — максимальное число задач этого коннектора;
- ❑ `connection.url` — URL подключения к базе данных;
- ❑ `mode` — метод отслеживания изменений JDBC-коннектором для источника;
- ❑ `incrementing.column.name` — название столбца, по которому отслеживаются изменения;
- ❑ `topic.prefix` — Kafka Connect записывает все таблицы в топики с названиями топик.префикс + название_таблицы.

* log-based CDC source connector (debezium)

14 декабря 2020 г. 13:12

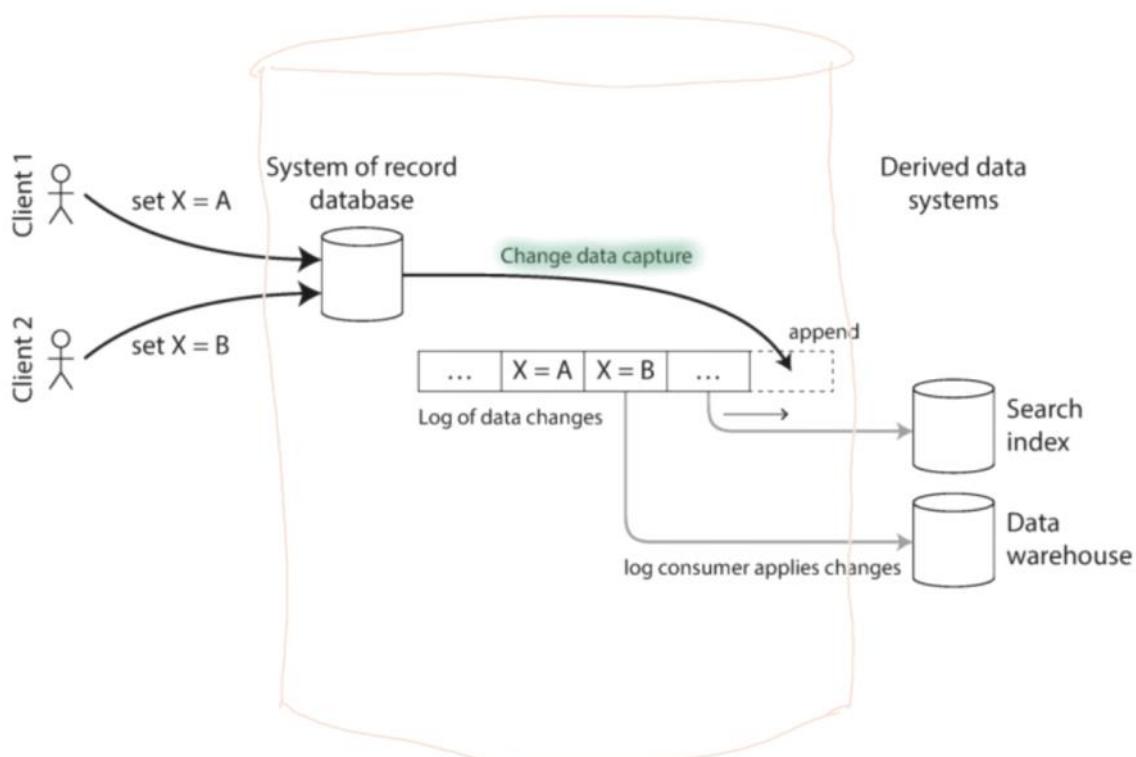
log based CDC connector может даже лезть в журналы логов БД (тесно интегрирован с БД)

- log-based означает что коннектору надо влезть в кишки/log самой БД

- так как у каждой БД свой подход к transaction logs and change logs
- Debezium подходит для нескольких механизмов БД

например, query CDC connector - в отличие от log-based CDC коннекторов, этот коннектор не может определять deleted rows (а также проблемы с инсертами и апдейтами)

- CDC использует собственный «событийный» интерфейс,
- (а) очень эффективен, поскольку коннектор отслеживает файл или запускается непосредственно при возникновении изменений, а не отправляет запросы через основной API базы данных, и
- (б) очень точен, поскольку выдача запросов через основной API базы данных часто создает возможность пропуска операций, если несколько операций поступают для одной и той же строки в течение периода опроса.



Как и брокеры сообщений, перехват данных обычно является асинхронным:

- система записи данных в базу не ждет, пока изменения будут применены у потребителей, прежде чем завершить операцию. Эксплуатационные преимущества такой архитектуры заключаются в том, что появление медленного потребителя не очень сильно влияет на систему записи, но имеет такой недостаток: ей свойственны все проблемы, связанные с задержкой репликации

- Kafka Connect [41] — это попытка объединить Kafka с инструментами сбора информации об изменениях данных для широкого спектра баз
 - те позволяет создать стрим изменений из БД

Представление баз данных в виде потоков открывает широкие возможности для интеграции систем

- Мы также увидели, что полезно осуществлять запись в базу данных в виде потока: можно фиксировать историю всех изменений, сделанных в базе, — либо неявно, путем захвата данных об изменениях, либо явно, через источники событий. Уплотнение журнала позволяет потоку сохранять полную копию содержимого БД.
- Можно постоянно поддерживать актуальность производных информационных систем, таких как индексы поиска, кэши и аналитические системы, потребляя журнал изменений и применяя их к производной системе. Можно даже создавать новые представления для уже существующих данных, потребляя журнал изменений с самого начала вплоть до настоящего времени

КТО ИСПОЛЬЗУЕТ

- Эта идея, будучи масштабированной, используется в
- LinkedIn Databus [25],
- Facebook Wormhole [26] и
- Yahoo! Sherpa [27].
- B Bottled Water реализован CDC для PostgreSQL, где специальный API декодирует журнал с упреждающей записью [28].
- B Maxwell и Debezium реализовано нечто похожее для MySQL, с синтаксическим анализом двоичного журнала [29–31].
- Mongoriver читает операционный журнал MongoDB [32, 33],
- a GoldenGate предоставляет аналогичные возможности для Oracle [34, 35]

PostgreSQL



- We're going to start an PostgreSQL instance to Sink the data to
- We will use Docker to start our PostgreSQL instance
- This will serve as our Sink for our second Sink Connector

- PostgreSQL is a very popular relational database
- Used by many companies around the world to back their website or run BI reports onto
- No knowledge of PostgreSQL is required for this course

Example: JDBCConnector DISTRIBUTED MODE



- Goal:
 - Start an PostgreSQL instance using Docker
 - Run in distributed mode with multiple tasks



- Learning:
 - Learn about the JDBC Sink Connector

The screenshot shows a browser window with the following tabs:

- Kafka Connect UI
- 127.0.0.1:9200
- Managing Connectors — Confluent Platform 3.2.0
- Releases · sqlectron/sqlectron
- JDBC Sink Configuration Options

The main content area has a header "JDBC Sink Configuration Options". On the left, there's a sidebar with navigation links:

- What is the Confluent Platform?
- Confluent Platform 3.2.0 Release Notes
- Installation
- Upgrade
- Confluent Platform Quickstart
- Application Development
- Operations
- Docker
- Confluent Control Center
- Kafka Security
- Kafka Streams
- Kafka Connect**
- Introduction

The right side contains detailed configuration options:

- connection.url**: JDBC connection URL.
 - Type: string
 - Importance: high
- connection.user**: JDBC connection user.
 - Type: string
 - Default: null
 - Importance: high
- connection.password**: JDBC connection password.
 - Type: password
 - Default: null
 - Importance: high

At the bottom of the page, there are several hashtags and a note:

```
#####
# C) PostgreSQL demo
# Examples are taken from here: http://docs.confluent.io/3.2.0/connect/connect-jdbc/docs/sink\_connector.html#quickstart
```

ИСПОЛЬЗУЕМ landoop UI

```
version: '2'

services:
  # this is our kafka cluster.
  kafka-cluster:
    image: landoop/fast-data-dev:cp3.3.0
    environment:
      ADV_HOST: 127.0.0.1      # Change to 192.168.99.100 if using Docker Toolbox
      RUNTESTS: 0              # Disable Running tests so the cluster starts faster
    ports:
      - 2181:2181      # Zookeeper
      - 3030:3030      # Landoop UI
      - 8081-8083:8081-8083  # REST Proxy, Schema Registry, Kafka Connect ports
      - 9581-9585:9581-9585  # JMX Ports
      - 9092:9092      # Kafka Broker

  # we will use elasticsearch as one of our sinks.
  # This configuration allows you to start elasticsearch
  elasticsearch:
    image: itzg/elasticsearch:2.4.3
    environment:
      PLUGINS: appbaseio/dejavu
      OPTS: -Dindex.number_of_shards=1 -Dindex.number_of_replicas=0
    ports:
```

```
- "9200:9200"
```

```
# we will use postgres as one of our sinks.  
# This configuration allows you to start postgres  
postgres:  
image: postgres:9.5-alpine  
environment:  
  POSTGRES_USER: postgres # define credentials  
  POSTGRES_PASSWORD: postgres # define credentials  
  POSTGRES_DB: postgres # define database  
ports:  
- 5432:5432 # Postgres port
```

настроим коннектор через UI

The screenshot shows the Kafka Connect UI interface. On the left, there is a list of existing connectors: file-stream-demo-distributed, logs-broker, logs-connect-distributed, logs-rest-proxy, logs-schema-registry, logs-zookeeper, sink-elastic-twitter-distributed, and source-twitter-distributed. A 'NEW' button is visible at the top right of this list. On the right, a detailed configuration page for a 'JdbcSinkConnector' is displayed. The connector class is set to 'io.confluent.connect.jdbc.JdbcSinkConnector'. Below this, a code editor contains the configuration JSON. At the bottom right of the configuration page, there are 'VALIDATE' and 'VALIDATE AND CREATE' buttons.

```
# Basic configuration for our connector  
name=sink-postgres-twitter-distributed  
connector.class=io.confluent.connect.jdbc.JdbcSinkConnector  
# We can have parallelism here so we have two tasks!  
tasks.max=1  
topics=demo-3-twitter  
  
# the input topic has a schema, so we enable schemas conversion here too  
key.converter=org.apache.kafka.connect.json.JsonConverter  
key.converter.schemas.enable=true  
value.converter=org.apache.kafka.connect.json.JsonConverter  
value.converter.schemas.enable=true  
  
# JDBCConnector specific configuration  
# http://docs.confluent.io/3.2.0/connect/connect-jdbc/docs/sink\_config\_options.html  
connection.url=jdbc:postgresql://postgres:5432/postgres  
connection.user=postgres  
connection.password=postgres  
insert.mode=upsert  
# we want the primary key to be offset + partition  
pk.mode=kafka  
# default value but I want to highlight it:  
pk.fields=__connect_topic,__connect_partition,__connect_offset  
fields.whitelist=id,created_at,text,lang,is_retweet  
auto.create=true  
auto.evolve=true
```

? нужно сделать первоначальную настройку

```
# Load the recent and long term stats into Postgres using Kafka Connect Sink!
confluent load SinkTopics -d kafka-connectors/SinkTopicsInPostgres.properties
```

? auto table creation

* HDFS sink connector

14 декабря 2020 г. 13:26

HDFS Sink Connector: Overview

- Continuously polls from Kafka and writes to HDFS (Hadoop Distributed File System)
- Integrates with Hive
 - Auto table creation
 - Schema evolution with Avro
- Works with secure HDFS and the Hive Metastore, using Kerberos
- Provides exactly once delivery
- Data format is extensible
 - Avro, Parquet, custom formats
- Pluggable Partitioner, supporting:
 - Kafka Partitioner (default)
 - Field Partitioner
 - Time Partitioner
 - Custom Partitioners

The HDFS connector allows you to export data from Kafka topics to HDFS files in a variety of formats and integrates with Hive to make data immediately available for querying with HiveQL.

The connector periodically polls data from Kafka and writes them to HDFS. The data from each Kafka topic is partitioned by the provided partitioner and divided into chunks. Each chunk of data is represented as an HDFS file with topic, kafka partition, start and end offsets of this data chunk in the filename. If no partitioner is specified in the configuration, the default partitioner which preserves the Kafka partitioning is used. The size of each data chunk is determined by the number of records written to HDFS, the time written to HDFS and schema compatibility.

The HDFS connector integrates with Hive and when it is enabled, the connector automatically creates an external Hive partitioned table for each Kafka topic and updates the table according to the available data in HDFS.

kafka connect security

18 декабря 2020 г. 11:39

Enabling Security with Connect

- If you have security enabled in the Kafka cluster (e.g. SSL or SASL), you need to enable it in Connect as well
 - Configure security on a per-worker basis
 - Override the default Producer or Consumer configurations that the worker uses
 - Configurations apply to *all* connectors running on the worker

You need the configuration at the worker top level because Connect leverages Kafka's group membership protocol to coordinate the workers and distribute work between them (as it does for Consumer Groups). The configuration at the producer/consumer level is required to interact with the Kafka cluster as a client.

If using SASL for authentication, you must use the same principal for workers and connectors as only a single JAAS is currently supported on the client side.

Securing passwords used for individual connectors is handled differently and will be discussed later in the module.

* syslog connector

18 декабря 2020 г. 14:40

- Once the connect service is ready, run the following command to configure the Syslog source connector in Kafka Connect:

```
$ curl -s -X POST \
-H "Content-Type: application/json" \
--data '{
  "name": "Syslog-Connector",
  "config": {
    "connector.class": "io.confluent.connect.syslog.SyslogSourceConnector",
    "tasks.max": "1",
    "syslog.port": "5454",
    "syslog.listener": "TCP",
    "confluent.topic.bootstrap.servers": "kafka:9092",
    "confluent.topic.replication.factor": "1",
    "confluent.license": ""
  }
}' http://connect:8083/connectors | jq
```

- In a new browser tab navigate to http://localhost:9021 to access Confluent Control Center:

- Select cluster **Cluster 1**
- then select the **Connect** tab. You will be shown the list of Connect clusters associated with your environment. There should only be a single one such cluster called `connect-default`
- select the `connect-default` cluster and verify that there is indeed a source connector `Syslog-Connector` defined and that it is in status running:

Status	Name	Category	Type	Topics	Number of tasks	
Running	Syslog-Connector	...	Source	SyslogSourceConnector	-	1

- Now let's test the connector by simulating a (syslog) message on port 5454 using the `netcat` (or `nc`) tool:

- Output a test log entry to port 5454:

```
$ echo "<34>1 2003-10-11T22:14:15.003Z host1.acme.com su - ID47 - Engine started" \  
| nc -w 1 connect 5454
```

8. Let's verify that the entry has been written to the `syslog` topic:

```
$ kafka-console-consumer \  
--bootstrap-server kafka:9092 \  
--property schema.registry.url=http://schema-registry:8081 \  
--topic syslog \  
--from-beginning | jq .'
```

You should see this: