

Kafka-docker-composer

Simple tool to create a docker-compose.yml file for failover testing

Introduction

Confluent has published official Docker containers for many years. They are the basis for deploying a cluster in Kubernetes using Confluent for Kubernetes (CFK), and one of the underpinning technologies behind Confluent Cloud.

For testing, the containers are convenient for quickly spinning up a local cluster with all the components required, such as the Confluent Schema Registry or the Confluent Control Center. You configure each container via a set of environment variables, using the internal docker network and hostname resolution to reference the various components such as the Kafka broker. This is easiest done using docker-compose, which will create a network for your cluster and allow you to keep all configurations in a single document.

This is easily done if you have a single ZooKeeper or KRaft controller and a single Kafka broker in your cluster.

The problem with a docker-compose file begins when you try to configure a larger cluster, for example, for failover testing. Suppose you want to convince yourself (or demonstrate) how the Kafka cluster gracefully handles the loss of a broker. In that case, you need to configure multiple Kafka brokers in your docker-compose.yml file.

You need to be careful not to reuse the same hostname or external port number, and when referencing, for instance, your bootstrap servers, you need to ensure to use the correct hostnames and ports. If you copy/paste and edit your entries, you may mistype or forget an entry. The cluster might still come up but will be degraded.

Of course, you can create a template for a 3 ZooKeeper / 3 Broker cluster and reuse that, but what if you want to test 4 or 6 brokers instead?

This is the situation I found myself in a few years ago with a customer who wanted to test out various configurations. It became very tedious to reconfigure the various docker-compose files for the different scenarios. So, I wrote a Python script that generates the docker-compose file for me: [kafka-docker-composer](#).

Installation

The tool requires a working Python 3 environment (I have tested with Python 3.8 - 3.11) with jinja2 installed, which comes automatically in many distributions. You can verify this with the following command

```
> python3 -c "import jinja2"
```

If this comes up with a `ModuleNotFoundError`, install jinja2 with pip3, for example

```
> pip3 install jinja2
```

You also will need to have Docker and docker-compose installed, of course.

Then clone the repository and you are ready to start creating docker-compose files.

```
> git clone https://github.com/sknop/kafka-docker-composer
> cd kafka-docker-composer
```

How?

The templating engine [Jinja2](#) allows me to create a template document with variables in it I can populate through my application.

This [template](#) simply loops through all configured services and populates the required field in the final docker-compose.yml file for each container I want to run. Each container encapsulates a service such as a Confluent Server or a ZooKeeper instance, a Schema Registry, and so on.

For each entry, I define the container image to be run as well as the host and container name. Each instance also has a whole set of optional parameters that are defined in the application:

- Health checks so I can define dependencies that will wait for the prerequisites to be healthy.
- Dependencies, which define the order in which containers are created. If health checks are defined as a dependency, this will also be listed here.
- The environment variables that are used for configuring the service within the container.
- Ports that are mapped to the docker host.
- Volumes that are mounted to inject additional files, such as Kafka Connect Connector plugins or metrics configurations for monitoring using Prometheus.

The application [kafka_docker_composer.py](#) takes a list of arguments and calculates how the template should be populated. It creates the dependencies between the different components,

ensures that names and ports are unique, the advertised listeners are correctly set up, and that dependent services like the Schema Registry of Kafka Connect point to the corresponding Confluent Server brokers.

There are many different configurations, such as using ZooKeeper or KRaft, all controlled by a set of arguments or a configuration file. Here is the help overview

```
> python3 kafka_docker_composer.py -h
usage: kafka_docker_composer.py [-h] [-r RELEASE] [--with-tc] [--shared-mode] [-b BROKERS]
[-z ZOOKEEPERS] [-c CONTROLLERS] [-s SCHEMA_REGISTRIES] [-C CONNECT_INSTANCES] [-k
KSQLDB_INSTANCES] [--control-center] [--uuid UUID] [-p] [--kafka-container KAFKA_CONTAINER]
[--racks RACKS]
                                [--zookeeper-groups ZOOKEEPER_GROUPS] [--docker-compose-file
DOCKER_COMPOSE_FILE] [--config CONFIG]

Kafka docker-compose Generator

options:
  -h, --help                show this help message and exit
  -r RELEASE, --release RELEASE
                                Docker images release [7.6.0]
  --with-tc                 Build and use a local image with tc enabled
  --shared-mode             Enable shared mode for controllers
  -b BROKERS, --brokers BROKERS
                                Number of Brokers [1]
  -z ZOOKEEPERS, --zookeepers ZOOKEEPERS
                                Number of ZooKeepers [0] - mutually exclusive with controllers
  -c CONTROLLERS, --controllers CONTROLLERS
                                Number of Kafka controller instances [0] - mutually exclusive with
ZooKeepers
  -s SCHEMA_REGISTRIES, --schema-registries SCHEMA_REGISTRIES
                                Number of Schema Registry instances [0]
  -C CONNECT_INSTANCES, --connect-instances CONNECT_INSTANCES
                                Number of Kafka Connect instances [0]
  -k KSQLDB_INSTANCES, --ksqldb-instances KSQLDB_INSTANCES
                                Number of ksqldb instances [0]
  --control-center          Include Confluent Control Center [False]
  --uuid UUID              UUID of the cluster [Nk018hRAQFytWskYqtQduw]
  -p, --prometheus          Include Prometheus [False]
  --kafka-container KAFKA_CONTAINER
                                Container used for Kafka, default [cp-server]
  --racks RACKS            Number of racks among which the brokers will be distributed evenly
[1]
  --zookeeper-groups ZOOKEEPER_GROUPS
                                Number of zookeeper groups in a hierarchy [1]
  --docker-compose-file DOCKER_COMPOSE_FILE
                                Output file for docker-compose, default [docker-compose.yml]
  --config CONFIG          Properties config file, values will be overridden by command line
```

arguments

Examples and Use Cases

This is all easier explained through some use cases.

Simplest cluster

```
> python3 kafka_docker_composer.py -b 1 -z 1
Generated docker-compose.yml
```

We just created a docker-compose file for the simplest case: 1 ZooKeeper and 1 Broker using the latest Confluent Docker images, as time of writing 7.6.0.

Not convinced? Try it out

```
> docker compose up -d
> docker compose ps --format "table {{.Name}}\t{{.Image}}\t{{.Status}}"
NAME                IMAGE                                STATUS
kafka-1             confluentinc/cp-server:7.6.0       Up 9 minutes (healthy)
zookeeper-1         confluentinc/cp-zookeeper:7.6.0    Up 9 minutes
```

You might have to use “docker-compose” instead of “docker compose” on your platform or upgrade your Docker environment to use [Compose V2](#).

The broker ports start with 9091 for the first broker, use kafka-topics to create and list topics:

```
> kafka-topics --bootstrap-server localhost:9091 --list
> kafka-topics --bootstrap-server localhost:9091 --create --topic test-topic \
  --replication-factor 1 --partitions 1
```

After you are done, with the cluster, shut it down again with

```
> docker compose down -v
[+] Running 3/2
✓ Container kafka-1 Removed
✓ Container zookeeper-1 Removed
✓ Network kafka-docker-composer_default Removed
```

The '-v' option removes the volumes as well, avoiding the potential problem of reusing stale data.

KRaft Controller

ZooKeeper is deprecated; modern versions of Kafka prefer KRaft. Just change the option from zookeeper to controller:

```
> python3 kafka_docker_composer.py --brokers 1 --controllers 1
Generated docker-compose.yml
```

Look inside the generated docker-compose.yml file to see which environment variables you must set to create a Controller-Broker pair successfully. Do you really want to set this by hand?

A more realistic cluster

Creating single broker setups is nice, but more is needed to show the power of this tool. A minimum standard cluster consists of three controllers and three brokers:

```
> python3 kafka_docker_composer.py -b 3 -c 3
Generated docker-compose.yml
> docker compose up -d
> docker compose ps --format "table {{.Name}}\t{{.Status}}\t{{.Ports}}"
```

NAME	IMAGE	STATUS
controller-1	Up About a minute	9092/tcp, 0.0.0.0:19091->19091/tcp
controller-2	Up About a minute	9092/tcp, 0.0.0.0:19092->19092/tcp
controller-3	Up About a minute	9092/tcp, 0.0.0.0:19093->19093/tcp
kafka-1	Up About a minute (healthy)	0.0.0.0:9091->9091/tcp, 0.0.0.0:10001->10001/tcp, 9092/tcp, 0.0.0.0:10101->8091/tcp
kafka-2	Up About a minute (healthy)	0.0.0.0:9092->9092/tcp, 0.0.0.0:10002->10002/tcp, 0.0.0.0:10102->8091/tcp
kafka-3	Up About a minute (healthy)	0.0.0.0:9093->9093/tcp, 9092/tcp, 0.0.0.0:10003->10003/tcp, 0.0.0.0:10103->8091/tcp

Note that each broker has its own externally visible port mapped to the host so that you can access each broker individually.

What about the other ports? They are for the JMX agents if you want to configure Prometheus and Grafana for your little cluster:

```
> python3 kafka_docker_composer.py -b 3 -c 3 -p
Generated docker-compose.yml
> docker compose up -d
> docker compose ps --format "table {{.Name}}\t{{.Status}}\t{{.Ports}}"
NAME                IMAGE                STATUS
controller-1        Up About a minute    9092/tcp, 0.0.0.0:19091->19091/tcp
controller-2        Up About a minute    9092/tcp, 0.0.0.0:19092->19092/tcp
controller-3        Up About a minute    9092/tcp, 0.0.0.0:19093->19093/tcp
grafana              Up 8 seconds         0.0.0.0:3000->3000/tcp
kafka-1              Up About a minute (healthy) 0.0.0.0:9091->9091/tcp,
0.0.0.0:10001->10001/tcp, 9092/tcp, 0.0.0.0:10101->8091/tcp
kafka-2              Up About a minute (healthy) 0.0.0.0:9092->9092/tcp,
0.0.0.0:10002->10002/tcp, 0.0.0.0:10102->8091/tcp
kafka-3              Up About a minute (healthy) 0.0.0.0:9093->9093/tcp, 9092/tcp,
0.0.0.0:10003->10003/tcp, 0.0.0.0:10103->8091/tcp
prometheus           Up 8 seconds         0.0.0.0:9090->9090/tcp
```

As you can see, Prometheus is exposed on port 9090 and Grafana on port 3000. Try out Grafana by pointing your browser to <http://localhost:3000>. The user and password are admin/admin. The dashboards need updating; you can always find the latest versions in [this repository](#).

The exporter configuration files and dashboards are in the **volumes** directory, as are the exporter jar, so you do not have to download anything separately.

Other components

In addition to the brokers and controllers, you can add Confluent Schema registries, Kafka Connect worker nodes, ksqldb nodes, and the Confluent Control Center to the mix. You will probably up the memory of your docker environment, specifically, if you run this on your notebook with Docker Desktop.

My Docker Desktop is configured with 8 cores and 16 GB of memory, which gives me ample room to run a large cluster in Docker compose.

The health checks are built for this purpose. If, for example, the schema registry starts up before the broker finishes booting, it will fail since it cannot create its topic, and it will not try again. Ensuring the brokers are up and running and ready to receive clients ensures that the dependent components do not fail on startup.

Connector Plugins

One specific component is the Connect cluster. This cluster comes with a bare-bone set of connector plugins installed. Still, it turned out that by mapping a volume containing unzipped connector plugin jar files before starting up the cluster, I could easily use the same setup to test various connectors cheaply. I have added a few connector plugins as examples, such as the Datagen Connector, which is useful to act as a producer for testing without external dependencies.

I have added an example plugin configuration [here](#). Bring up a cluster with

```
> python3 kafka_docker_composer.py -b 3 -c 3 -p -s 1 -C 2 --control-center
Generated docker-compose.yml
> docker compose up -d
> curl localhost:8083/connector-plugins | jq
```

This will list all installed connector plugins. You might have to install [jq](#), a handy JSON formatting and filtering tool.

Create a target topic, then install the Datagen connector:

```
> kafka-topics --bootstrap-server localhost:9092 --create --topic users
Created topic users.
> curl -X PUT -H "Content-Type: application/json" \
  --data @connect/datagen.json localhost:8083/connectors/datagen-users/config | jq
{
  "name": "datagen-users",
  "config": {
    "connector.class": "io.confluent.kafka.connect.datagen.DatagenConnector",
    "key.converter": "org.apache.kafka.connect.storage.StringConverter",
    "kafka.topic": "users",
    "confluent.topic.replication.factor": "3",
    "quickstart": "users",
    "max.interval": "1000",
    "iterations": "10000000",
    "tasks.max": "2",
    "name": "datagen-users"
  },
  "tasks": [],
  "type": "source"
}
> curl -s localhost:8083/connectors/datagen-users/status | jq
```

```
{
  "name": "datagen-users",
  "connector": {
    "state": "RUNNING",
    "worker_id": "kafka-connect-1:8083"
  },
  "tasks": [
    {
      "id": 0,
      "state": "RUNNING",
      "worker_id": "kafka-connect-1:8083"
    },
    {
      "id": 1,
      "state": "RUNNING",
      "worker_id": "kafka-connect-2:8084"
    }
  ],
  "type": "source"
}
```

Other connector plugins can be installed by downloading the zip file from the connector hub, for example, the ElasticSearch Sink connector from <https://www.confluent.io/hub/confluentinc/kafka-connect-elasticsearch>. Unzip the file into the volumes/connect-plugin-jars directory and restart your connect clusters:

```
> docker compose restart kafka-connect-1 kafka-connect-2
```

Multi-DC scenarios

The original purpose of the kafka-docker-composer was to test failover scenarios in multiple data centres. For this purpose, there are two additional options: racks and zookeeper groups.

Three data centres and KRaft

Specify the number of racks for the number of unique racks or data centres you want to test with. A typical example is to choose 3 racks, which is a common setup for production clusters. The configured brokers will then be assigned round-robin to the racks. This is particularly useful if you configure more brokers than the number of racks to test the distribution of partitions across the different racks.


```
> python3 kafka_docker_composer.py -b 6 -c 3 -racks 3
Generated docker-compose.yml
> docker compose up -d
```

Have a look at the generated docker-compose.yml file. You will notice that both the controllers and the brokers have a new environment variable, for example,

KAFKA_BROKER_RACK: rack-2

Create a new topic called products, but add the option --partitions 6, then run

```
> kafka-topics --bootstrap-server localhost:9092 --describe --topic products
Topic: products    PartitionCount: 6  ReplicationFactor: 3    Configs:
min.insync.replicas=2
    Topic: products  Partition: 0 Leader: 5    Replicas: 5,6,7
    Topic: products  Partition: 1 Leader: 9    Replicas: 9,4,8
    Topic: products  Partition: 2 Leader: 7    Replicas: 7,5,6
    Topic: products  Partition: 3 Leader: 8    Replicas: 8,9,4
    Topic: products  Partition: 4 Leader: 6    Replicas: 6,7,5
    Topic: products  Partition: 5 Leader: 4    Replicas: 4,8,9
```

If you go through every single partition and compare their placement on the individual replicas, you will notice that each replica is in a different rack. Even in the case of a loss of one rack each partition will still have 2 replicas online. Since min.insync.replicas is set to 2, producers and consumers will still be able to work.

You can then test what happens if a data centre goes down by using standard docker-compose methods to kill the containers. Attach a producer and consumer to the cluster to convince yourself that the cluster is still accessible and capable of processing requests. This is how I demonstrate to colleagues, partners, and customers the resilience features of Apache Kafka.

In our example, we have the following distribution:

Container	Rack	Broker Id
controller-1	rack-0	1
controller-2	rack-1	2
controller-3	rack-2	3
kafka-1	rack-0	4

kafka-2	rack-1	5
kafka-3	rack-2	6
kafka-4	rack-0	7
kafka-5	rack-1	8
kafka-6	rack-2	9

This is because the tool assigned the rack id round-robin.

We need to produce and consume some data. Start with the consumer readily waiting for some messages:

```
> kafka-console-consumer --bootstrap-server \
localhost:9091,localhost:9092,localhost:9093 --topic products
```

Note that I am specifying three brokers across all three racks here. For this test, it is not strictly necessary because we will receive the full list of all brokers upon connection anyway, but it is good practice in case a broker or a whole rack (data centre) is down while we are starting our application.

Then we start the producer:

```
> seq -f "%0.0f" 1000000 10000000 | kafka-console-producer --bootstrap-server \
localhost:9091,localhost:9092,localhost:9093 --topic products
```

The tool **seq** generates a sequence of numbers separated by a newline. I added the formatting to avoid presenting the numbers in scientific notation. The result is a new message for each number, nicely sequenced for easy verification. Note that there is no guarantee that the consumer will return the numbers in the same order since we created the topic with 6 partitions.

On my Mac, the producer will work through this sequence quite quickly, so we need to hurry with the next step. To simulate an outage of one data centre, we need to kill all containers in one rack. Let's pick rack-2:

```
> docker compose kill controller-3 kafka-3 kafka-6
```

Note that the consumer will stumble for a moment while the active controller sorts out the leadership for each partition, but it will then pick up again. When the consumer does not find any new messages anymore it will wait and we can shut it down. It should print out the total number of messages: 9000001.

You can vary the sequence to test out the loss of a rack for producers as well. You should notice a lot of error messages when the producer complains about brokers not being reachable anymore, but it will sort itself out after a while. This is also a good example for the idempotent producer since no messages will be duplicated during this time, as you will be able to verify with your consumer.

Bring the cluster back up again with a simple

```
> docker compose up -d
```

Depending on how long the producer was working while the cluster was in a degraded state, this last command might take a while to finish, because the brokers have to catch up first before they respond to the health-check command.

Two data centres (DC) and ZooKeeper Groups

The other option is for configuring ZooKeeper groups, specifically for a 2-DC scenario with hierarchical groups, for which you will need 6 ZooKeeper instances minimum. The tool will calculate the distribution and set up the docker-compose file accordingly.

```
> python3 kafka_docker_composer.py -b 4 -z 6 --zookeeper-groups 2
```

If you look at the generated docker-compose.yml file, you will find these lines

```
ZOOKEEPER_GROUPS: 1:2:3;4:5:6
```

These are the generated groups, with ZooKeeper instances distributed between them. The cluster will come up just fine with these settings, but when you shut one (simulated) data centre down, you will see why a 2-DC solution is always inferior to the 3-DC setup.

The ZooKeeper instances do not automatically failover to the degraded state of 3 ZooKeepers. Instead, they will refuse to accept all connections until the situation is resolved. This also means brokers will refuse to acknowledge producers and even consumers will fail because the brokers need to update the ZooKeeper about their progress.

You can resolve the situation by starting the second DC up again, in which case the cluster will recover. If you want to continue in the degraded state, you need to manually remove the groups from the configuration files and restart the remaining ZooKeeper instances.

The required procedure goes beyond the scope of this post. In production environments, we advise administrators to keep a second configuration file handy that can be swapped in after the ZooKeeper instances have been shut down. This is not something you can and want to

automate, but you can script it and invoke the failover script manually should the need arise. The procedure should be well documented and tested in a non-production environment to ensure that administrators know during the panic of an outage what to do. Note that this is not possible without downtime (RTO > 0).

Further features and outlook

Looking through the help list, you might notice a few other arguments that have not been discussed.

Shared mode

This is an experimental feature for upgrading KRaft Controllers to full brokers. This is not officially supported for production environments, but if you want to play with it, you can use the “--shared-mode” argument. If nothing else, it will show you which additional environment variables you need for this setup.

```
> python3 kafka_docker_composer.py -c 3 -b 3 --shared-mode
> docker compose up -d
```

When you start up this cluster, you will see that you have 6 brokers active rather than the original 3 since the controllers act as brokers as well. If you have [kcat](#) (kafkacat) installed, you can use

```
kcat -L -b localhost:9094
```

to verify this.

UUID

A cluster running with KRaft needs a UUID to identify membership for controllers and brokers. I have pre-created and hardcoded such a UUID, but if you want to change this value, use the --uuid option.

```
> kafka-storage random-uuid
CzJxwb_zS0CSAEyEZn8G_A
> python3 kafka_docker_composer.py -c 3 -b 3 --uuid CzJxwb_zS0CSAEyEZn8G_A
```

TC

I have built the infrastructure for it but have yet to do any testing with it. This option enables a build for all components to create a new image that contains the **tc** tool that can be used to inject latency. There is a [Medium article](#) that explains a little more about this feature if you are

interested, but most of the ideas for this tool come from the [Confluent tutorial](#) on multi-region clusters.

I have verified that tc is running successfully in each image, but I have not done any latency testing yet. If you want to try this out and share your setup and experience, please drop me a line.

Build steps

Building the docker images with tc enabled involves multiple steps:

- Adjust the .env file in the root directory to the release you want to base your test on. Apologies for the repetition: these are straight from the underlying source and have a lot of redundancies. I have just updated this to 7.6.0.
- Run the build script `script/build_docker_images.sh`. This will download the configured base images and build new images stored in your local image cache.
- Run `kafka-docker-compose` with the `--with-tc` option enabled to create your docker-compose file. Ensure the version matches the version you have just created. Use the `--release` option if necessary.

You now have a configuration with Docker images that have been enhanced with the tc utility.

To start this up and test it, do the following:

- Bring the cluster up as before
- Log into one instance and run a ping against another instance to verify the normal latency
- Enable latency injection in each node
- Run ping again to verify the effect
- Perform your tests with a multi-region simulator

Here is the simplest example:

```
> scripts/build_docker_images.sh
> python3 kafka_docker_composer.py -c1 -b1 --with-tc
> docker compose up -d
> docker compose exec controller-1 ping kafka-1
> docker compose exec -u0 controller-1 \
    tc qdisc add dev eth0 root netem delay 100ms
> docker compose exec -u0 kafka-1 \
    tc qdisc add dev eth0 root netem delay 100ms
> docker compose exec controller-1 ping kafka-1
```

You should now observe a difference in the ping round time of around 200ms, with 100ms latency injected from each side.

Future work

Security

I have yet to add any authentication or authorization features to this tool, mostly because I have other tools to test and demonstrate security. Still, it might be a worthwhile project, certainly for SASL/PLAIN or even SASL/SCRAM. TLS certificates are a bit trickier because I'd need a new image to generate these. I want the whole project to start with a single **docker compose up**, not a script.

The same is true for Kerberos and LDAP for RBAC, which would probably require a Samba service in a separate container and some configuration. Let me know in the comments or file an issue on GitHub.

Enable Open Source Kafka

Apache Kafka 3.7 is out and comes with an official Docker image as of [KIP-975](#).

I have successfully tested the image separately from kafka-docker-compose and I think it might be useful to enable the swapping out of the Confluent Server (cp-server) image to test out the new features in the next Apache Kafka release — since these releases are published typically 3-6 months before the corresponding Confluent release. The main difficulty will be that open-source Apache Kafka has no built-in REST interface, making health checks more challenging.

Conclusion

I have been using and extending kafka-docker-composer for the last 5 years whenever the need arose to demonstrate how to set up a cluster in Docker. The main purpose of this tool was to show how resilient a Confluent Cluster is even during a large outage.

Lately, I have used the same tool to teach myself KRaft, experiment with it, and use the setup for connector testing and development.

I'd like to know what you will use this tool for. Let me know in the comments or by contacting me directly.

Happy hacking!