

---

# Queuing and scaling in Kubernetes infrastructure

---

PS6 - Computer Science

Report

Fribourg, January 2021

Romain Agostinelli

[romain.agostinelli@edu.hefr.ch](mailto:romain.agostinelli@edu.hefr.ch)

Supervisors:

Mäder Michael

[michael.maeder@hefr.ch](mailto:michael.maeder@hefr.ch)

Supcik Jacques

[jacques.supcik@hefr.ch](mailto:jacques.supcik@hefr.ch)

Version: 4.0



Haute école d'ingénierie et d'architecture Fribourg  
Hochschule für Technik und Architektur Freiburg

## Contents

<b>1</b>	<b>Document's history</b>	<b>1</b>
<b>2</b>	<b>Introduction</b>	<b>2</b>
<b>3</b>	<b>Project's environment</b>	<b>2</b>
3.1	Infrastructure . . . . .	2
3.1.1	Kubernetes . . . . .	2
3.1.2	Minikube . . . . .	4
3.2	Technologies . . . . .	4
3.2.1	Programming Language . . . . .	4
3.2.2	Frameworks . . . . .	4
3.2.3	Queuing . . . . .	4
<b>4</b>	<b>Technical research</b>	<b>5</b>
4.1	Scaling . . . . .	5
4.2	K8s HPA . . . . .	5
4.2.1	How it works . . . . .	6
4.2.2	Installation and tests . . . . .	7
4.3	Custom Metrics Server . . . . .	8
4.3.1	Installation and tests . . . . .	9
4.4	Kubernetes Jobs . . . . .	11
4.5	Sidecar pattern . . . . .	11
4.6	Implementation and tests . . . . .	12
4.7	Forwarding / Redirecting requests . . . . .	12
4.8	Passing K8s environment variable to Spring Boot container . . . . .	14
4.8.1	K8s vs Docker vs Spring environment variables . . . . .	14
4.8.2	Installation and tests . . . . .	15
4.9	Queuing with KubeMQ . . . . .	17
4.9.1	Queues . . . . .	17
4.9.2	Publisher / Subscriber . . . . .	18
4.10	Monitoring KubeMQ . . . . .	20
4.11	Reverse Proxies . . . . .	21
<b>5</b>	<b>Finding solution iterations</b>	<b>21</b>
5.1	Jobs . . . . .	21
5.2	Custom metrics server . . . . .	22
5.3	Proxy and Adapter . . . . .	23
<b>6</b>	<b>Kueuer</b>	<b>24</b>
6.1	Conception - Design . . . . .	24
6.1.1	Architecture . . . . .	24
6.1.2	Communication . . . . .	25
6.2	Proxy . . . . .	28
6.2.1	Design . . . . .	28
6.2.2	Implementation . . . . .	29
6.2.3	Deployment . . . . .	29
6.3	QAdapter . . . . .	30
6.3.1	Design . . . . .	30
6.3.2	Implementation . . . . .	32
6.3.3	Deployment . . . . .	34
6.4	Responder . . . . .	34
6.4.1	Design . . . . .	34
6.4.2	Implementation . . . . .	36

6.4.3	Deployment . . . . .	37
6.5	Queues . . . . .	38
6.6	Monitoring and HPA . . . . .	38
6.6.1	Prometheus Adapter . . . . .	39
6.6.2	HPA . . . . .	39
6.6.3	Implementation . . . . .	39
6.6.4	Deployment . . . . .	40
6.7	Evolution . . . . .	42
6.7.1	Automation / Injection . . . . .	42
6.7.2	Language Refactoring . . . . .	42
<b>7</b>	<b>Tests</b>	<b>42</b>
7.1	Tools . . . . .	43
7.1.1	Siege . . . . .	43
7.1.2	K6 . . . . .	43
7.2	Tests and results . . . . .	43
7.2.1	Environment simulation . . . . .	43
7.2.2	Threshold . . . . .	44
7.2.3	Functional testing . . . . .	45
7.3	Stress test . . . . .	47
7.4	Bugs reported / fixed . . . . .	48
7.4.1	Responder closes pending requests . . . . .	48
7.4.2	Responder stores too much events . . . . .	49
7.4.3	QAdapter auto-scale loss of messages . . . . .	50
7.5	Tests synthesis . . . . .	51
<b>8</b>	<b>Conclusion</b>	<b>52</b>
8.1	Objectives and achievement . . . . .	52
8.2	State of the project . . . . .	52
8.3	Evolution of the project . . . . .	52
8.4	Personal conclusion . . . . .	53
<b>9</b>	<b>Resources</b>	<b>53</b>
<b>10</b>	<b>Déclaration d'honneur - Declaration of honour</b>	<b>53</b>
	<b>References</b>	<b>55</b>
<b>11</b>	<b>Appendixes</b>	<b>56</b>
11.1	Simple App with HPA based on CPU consumption . . . . .	56
11.2	Environment deployment configuration files . . . . .	57
11.2.1	App with QAdapter Deployment file . . . . .	57
11.2.2	App - HPA Based on Queue metrics . . . . .	58
11.2.3	Proxy Deployment file . . . . .	59
11.2.4	Reponder deployment configuration . . . . .	60
11.3	Metrics configurations files . . . . .	60
11.3.1	Prometheus Configuration . . . . .	60
11.3.2	Prometheus adapter configuration . . . . .	65
11.4	K6 Tests configurations . . . . .	66
11.4.1	Functional testing K6 configuration . . . . .	66
11.4.2	K6 Stress test configuration . . . . .	67
<b>12</b>	<b>Appendix - Cahier des charges - Specification</b>	<b>67</b>
12.1	Introduction . . . . .	67
12.2	Context . . . . .	68

12.3	Objective . . . . .	68
12.4	Infrastructure topologies . . . . .	68
12.4.1	First part . . . . .	68
12.4.2	Second part . . . . .	69
12.5	Activities . . . . .	69
12.5.1	Kickoff and specs . . . . .	69
12.5.2	Setting up a development environment . . . . .	70
12.5.3	First micro-service . . . . .	70
12.5.4	Solution . . . . .	71
12.5.5	Complex infrastructure . . . . .	71
12.6	Planning . . . . .	71

## 1 Document's history

Sous cette section se trouve l'historique du document avec ce qui a été ajouté, supprimé ou modifier en fonction des versions.

Version	Modifications
v1.0	Document creation
	Introduction
	Project's environment
v1.1	Technical research
	Scaling
	K8s HPA
v1.2	Custom Metrics Server
	K8s Jobs
	Sidecar pattern
	Forwarding - Receiving
v1.3	Passing K8s env. var.
	Queuein with KubeMQ
v.2	Monitor KubeMQ
v2.1	Finding solution iterations
	Kueuer - Conception design
	Kueuer - Proxy
	Kueuer - QAdapter
v.2.2	Kueuer - Responder
v.3	Kueuer - Queues/Monitoring HPA
	Kueuer Evolutions
v.3.1	Tests
v.3.2	Refactoring
	Research - Reverse Proxy
v.4	BibTex
	Glossary
	Final check

## 2 Introduction

This project consists in developing a system that provide mechanisms to handle a high rate of request/second to a service considering that this service can take up to 30sec to answer. The goal is to limit the "waiting" time of each requests and do not create a bottleneck as the requests arrives.

For example, let's consider a movie theater with theses constraints:

- A new film has just been released and everyone wants to see it.
- Before entering the movie theater, people must give their ticket to a ticket booth.
- The duration of the film is 30sec.
- It can only have one person in the movie theater room to see the film.
- People don't want to wait more than 30sec before being redirected to a room.
- A cinema room costs money and we try to minimize the total cost.

0798912175 In the example, the service is the movie theater that offers different path (films) and the people coming to the cinema are the requests. The duration of the film refers to the computation time of the service for this path.

The enumerated constraints above describe the problem we must resolve. How to make all the people happy? How to satisfy all these constraints?

In this project, we'll investigate how to solve this problem by using queues and (auto-) scaling.

Queuing in IT is the same idea of queuing in real life. If we take back the described situation with the cinema, queuing is setting a queue before the ticket booth so people can wait here and be taken in the order they arrived.

Scaling is the concept of increasing the number of cinema room regarding the number of people waiting and decreasing this number if nobody is here (at midnight for example) because an opened cinema room costs money.

## 3 Project's environment

This section describes in which environment the project is done, for which use case and so on. It contains two subsections referring to the main infrastructure aspects: first the infrastructure in which the project takes place and then the technologies used for external components.

### 3.1 Infrastructure

The project takes place in a Kubernetes (K8s) infrastructure. The services are in a K8s environment and the requests are directed to the K8s cluster.

For containerization technology, Docker is used. Be careful because K8s just announce (at the time this project has begun) that Docker will be deprecated in the future<sup>1</sup>.

#### 3.1.1 Kubernetes

On [K8s' website](#)<sup>2</sup> we found a sentence that describe really well what Kubernetes is: "Kubernetes, also known as K8s, is an open-source system for automating deployment, scaling,

<sup>1</sup><https://kubernetes.io/blog/2020/12/02/dont-panic-kubernetes-and-docker/>

<sup>2</sup><http://kubernetes.io/>

and management of containerized applications".

K8s was initially developed by Google to manage a lot of services deployed in container. K8s is not a container platform as Docker, but it uses it.

Kubernetes allows you to manage your container by dispatching them on different machines, by creating namespaces, clusters and so on. Bellow, we discuss some features or terms related to K8s and necessary to understand for the project. Please refer to the graphic on figure 1 :

- **Cluster:** A cluster is a logic component representing a whole infrastructure. The cluster contains all the nodes of the infrastructure. It is the biggest component.
- **Node:** A node represents a machine. This machine can be a physical machine or a virtual machine that runs different pods and/or service, etc. It must have a container run-time.
- **Pod:** A pod is the smallest object of K8s. It's a logical object that contains containers. Generally, there is one container per pod but for some use cases, a pod can contains multiple containers. When a pod contains multiple container, we can call them "sidecars" container. There is generally one primary container and multiple sidecars that can run in the same pod to help/interact with the primary container.
- **Service:** A service is component running on the node (same logical level than a pod) and creates a gateway between the Web and the cluster for one pod (generally). A service is most of the time linked to only one pod and provides access to the pod via web. If a pod is replicated (two or more same pods running in the same time), the service is capable of load balancing the traffic between all the replicas of the pod.

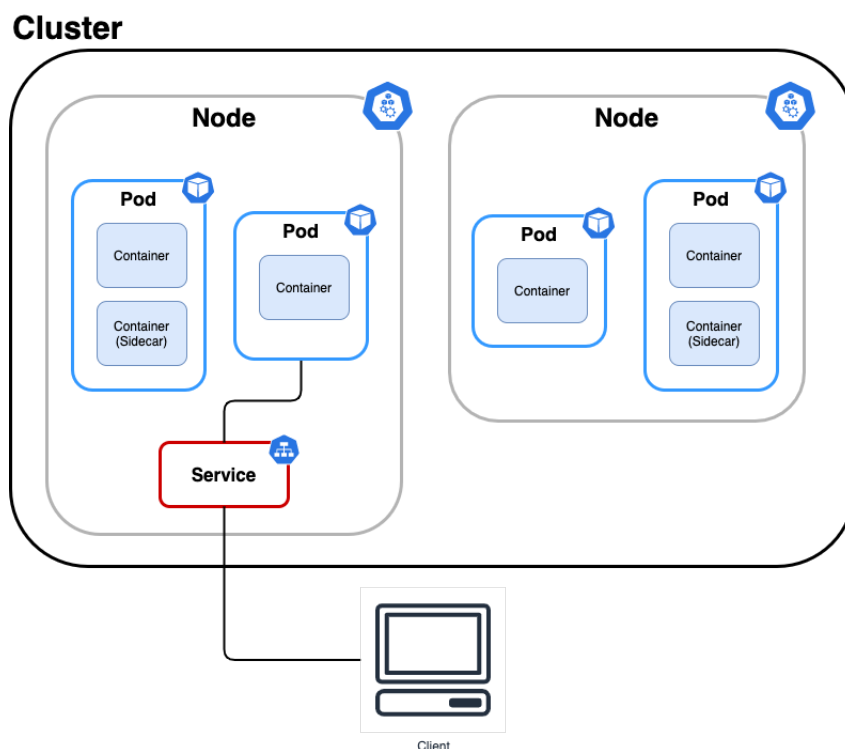


Figure 1: Basic K8s objects

### 3.1.2 Minikube

Minikube is a small local K8s that can be easily installed on different OS like Windows, macOS and Linux. It is used to deploy a Kubernetes cluster with one node, which is the machine on which Minikube is launched. Minikube provides a fast installed Kubernetes development environment and allows you to choose between different container providers. It is adapted for development and not for production.

More can be seen on [Minikube's website](https://minikube.sigs.k8s.io/docs/)<sup>3</sup>.

As an alternative for Minikube, there is [Kind](#)[1] that provides a quick way to setup a running cluster on our machine. This solution is also proposed by Kubernetes.

For this project, Minikube is used because at the beginning of the project I wasn't familiar with Kubernetes and all the demonstration on Kubernetes' website are with Minikube. So I decided to use it for simplicity.

## 3.2 Technologies

This chapter discusses different technologies that is used in the project, like programming language, framework and queuing solution.

### 3.2.1 Programming Language

In this project, Java is used to develop different solutions or tests. The Java language has the advantages to be highly known and fast for API.

Kubernetes is developed in Go and all that is developed for is in Go. Unfortunately, I (Romain Agostinelli) don't know Go and even if it would have been the right choice, I decided to not use Go as it can slow me down during the project.

Java seems a good alternative, not as fast as Go but more known and easy to setup.

### 3.2.2 Frameworks

As Java is used as programming language, the fastest way to setup Rest application is to use Spring Boot. Spring Boot provides a fast way to setup Java application and it has a big community with lot of tutorials and solutions.

### 3.2.3 Queuing

For the queuing part, there are multiple technologies that can be used. Just to cite the most well known and used:

- [Apache Kafka](#)
- [RabbitMQ](#)
- [Redis](#)

Instead of the three cited queuing solutions, there is one that is really recent (2018-2019) and seems amazing on the paper. It is a Kubernetes Native Queuing System named [KubeMQ](#). Compared to the others queuing solutions, it has a feature that is important in this project: Exactly one Delivery. This feature that the others don't have means that each message will be delivered exactly once so it guarantees us that each message putted in a queue will be given to a service only once. So, if multiple services/threads are connected at the same time to a message queue, we are sure that only one will received the message and then it will be next processed only once too.

---

<sup>3</sup><https://minikube.sigs.k8s.io/docs/>



For more comparisons between KubeMQ and other solutions as Kafka, Redis and RabbitMQ, there is a page dedicated to it on [KubeMQ's website](https://kubemq.io/compare-kubemq/)<sup>4</sup>.

For this project, KubeMQ will be used as queuing solution.

## 4 Technical research

This section discusses all the researches and tests done with different technology before finding a concrete solution.

### 4.1 Scaling

There are two types of scaling: Vertical scaling and Horizontal scaling.

Vertical scaling consists of adding resources (such as RAM size or CPU core). Horizontal scaling consists in increasing the number of total machines by duplicating the existing ones and distributing the load between the replicas. We can see a graphical representation on figure 2

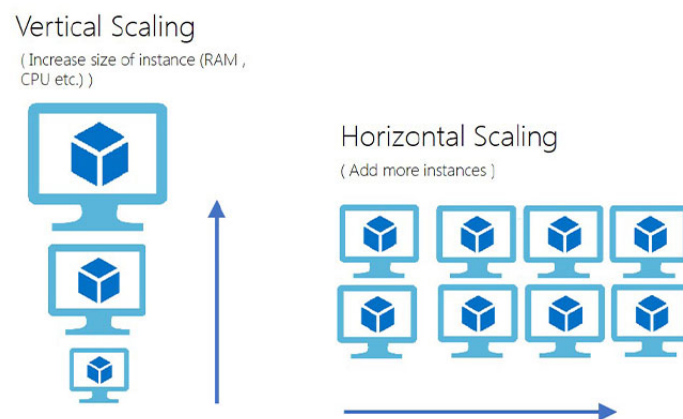


Figure 2: Horizontal vs Vertical scaling

Source: [webairy.com](https://webairy.com)

The vertical scaling is not appropriate for our problem because the issue is not that there aren't enough resources for computing but the computation just take times and there is too much requests.

### 4.2 K8s HPA

In Kubernetes environment, HPA stands for "Horizontal Pod Autoscaler". It is a system that can scale up (or down) the number of a pod regarding different metrics. A HPA is configured for only one pod, so, if there are several pods (not replicas of pod but clearly different pods who run various containers) that must be auto-scaled, it will have several HPA.

<sup>4</sup><https://kubemq.io/compare-kubemq/>

#### 4.2.1 How it works

There is multiple Kubernetes components involved into the HPA working flow (figure 3):

- CAdvisor: this is a lightweight component that run on each Node by default. It is capable of retrieving the CPU and memory usage of each pod inside its node.
- Kubelet: running on each Node, it offers the possibility to communicate with its node via API.
- Metric-Server: cluster-wide component ([aggregator](#))[2] that is not installed by default. Its job is to collect and aggregate all metrics (memory and CPU) from all nodes (via Kubelet) and provide these information through the Kubernetes API : `"/apis/metrics.k8s.io/"`.
- HPA: component working as controller (more info on controller pattern [here](#)[3]) that use the Kubernetes API to know the metrics for the pod it is linked to.
- Control Plane: logical view of various components which control the Kubernetes architecture. More info [here](#)[4].

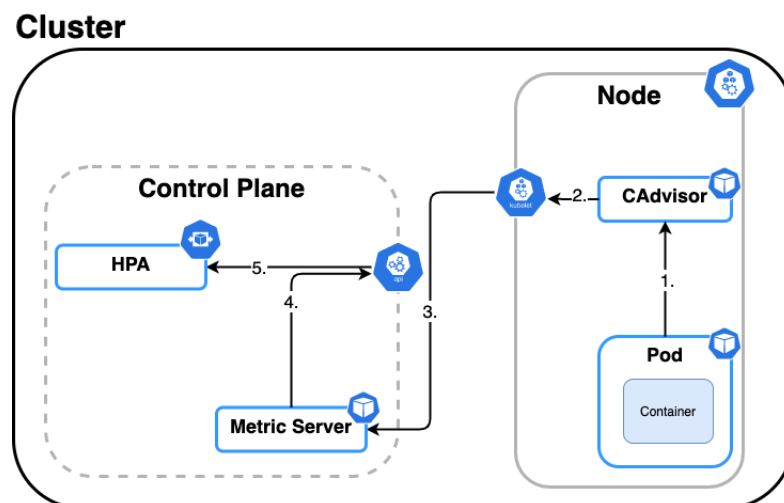


Figure 3: HPA working flow

On the figure 3 above, there is multiple steps for a HPA to works properly. Here their description:

1. CAdvisor collect information for each pod in its node.
2. CAdvisor provide collected information for its node via the kubelet component.
3. Metric-Server aggregate the metrics for each node in the cluster using the kubelet of each node.
4. Metric-Server provide metric for the whole cluster (each node/container of the cluster) via Kubernetes API.
5. HPA retrieves information about the pod that it must scale via the Kubernetes API.

By default, the CAdvisor and Metric-Server can retrieve only two types of metrics:

- Memory
- CPU

So the HPA, in its base configuration, can only be configured on these two metrics.

When setting up a HPA for a pod, the K8s infrastructure must know how many resources a container is allowed to use (limit) and how many it requests to run. If a node has enough resources, a container can use more than it requests, but not more than its limit. Resources can be CPU consumption and/or memory consumption. It is mandatory to setup the limit of each container in the pod deployment to use HPA on this pod. This way, the HPA can know when to scale up/down a pod. For example, if a container inside the pod begins to use more CPU than its limit, the HPA will scale up the pod to distribute the load between multiple replicas of the same pod to maintain a CPU consumption between request and limit for each pod. The full documentation is available [here](#)[5].

The next chapter (ch. 4.3) will discussed about the possibility to setup a custom metric server that can allow to provide more than these two type of metrics and maybe offers the HPA to be based on different custom metrics.

#### 4.2.2 Installation and tests

The documentation about the installation is available [here](#)[6].

As seen above, the K8s infrastructure needs a metric server to aggregate all the metrics from all the nodes. The metric server deployment file can be found on the [official repository](#). Warning: If you are running it on Minikube or Kind, please add two lines as follow in the metrics-server Deployment under the container spec (inside the template):

```
1 | command:
2 |   - /metrics-server
3 |   - --kubelet-insecure-tls
```

There is an example on the [project repository](#).

When the metric server is running, it's time to configure the application we want to automatically scale. For this, the application needs to tells its limits. To do so, a simple array can be added to the Deployment file of the pod under each container presents inside the template:

```
1 | resources:
2 |   limits:
3 |     cpu: 100m
4 |   requests:
5 |     cpu: 100m
```

A full example of a Deployment file is available [here](#)<sup>5</sup> and can be used as it is.

Finally, when everything is up deployed, it is time to setup a HPA. For this, a simple command can be used which contains the metric to look at for autoscaling, the minimum replicas of the pod and the maximum:

```
1 | #Generic command:
2 | kubectl autoscale deployment <name_of_deployment> <metric> --min=<min>
   | ↪ --max=<max>
3 | #Example for the deployment cited above
4 | kubectl autoscale deployment nodejs --cpu-percent=20 --min=1 --max=10
```

<sup>5</sup><https://gitlab.forge.hefr.ch/ps6-2021-kubemq/deployment/-/blob/cbae9c6728c859e68891d41214121b4468749291/simple-rest-deployment.yaml>

A HPA can also be configured to use multiple metrics and it can be done using a deployment file like this:

```

1  apiVersion: autoscaling/v2beta2
2  kind: HorizontalPodAutoscaler
3  metadata:
4    name: php-apache
5  spec:
6    scaleTargetRef:
7      apiVersion: apps/v1
8      kind: Deployment
9      name: php-apache
10   minReplicas: 1
11   maxReplicas: 10
12   metrics:
13   - type: Resource
14     resource:
15       name: cpu
16       target:
17         type: Utilization
18         averageUtilization: 50

```

More informations available on the HPA setup walkthrough available [here](#)<sup>[6]</sup>.

### 4.3 Custom Metrics Server

As discussed in the previous chapter (ch. 4.2), a HPA is based on metrics that the metric-server provide. By default, the metric server offers 2 types of metrics:

- CPU consumption
- Memory consumption

Theoretically, there is the possibility to create a custom metric server that provide information for a pod based on a custom worker. This way, it will become possible to deploy a custom metric server and base the HPA on this server and do not change the way all works.

The current version (v1beta2) of what a metric server must serve through its API for custom metrics is available [here](#)<sup>6</sup> and this uses the global "k8s.io/apimachinery/pkg/apis/meta/v1" for common (to all APIs in K8s) response information. This second file is available [here](#)<sup>7</sup>. These two files are written in Go.

After compiling and aggregating information, it is possible to create a Json structure of what a custom metric server must serve:

```

1  {
2    "kind": "MetricValueList",
3    "apiVersion": "custom.metrics.k8s.io/v1beta2",
4    "items": [
5      {
6        "metric": {
7          "name": "<name_of_custom_metric>"

```

<sup>6</sup>[https://github.com/kubernetes/metrics/blob/master/pkg/apis/custom\\_metrics/v1beta2/types.go](https://github.com/kubernetes/metrics/blob/master/pkg/apis/custom_metrics/v1beta2/types.go)

<sup>7</sup><https://github.com/kubernetes/apimachinery/blob/master/pkg/apis/meta/v1/types.go>

```

8         },
9         "describedObject": {
10             "kind": "<type_of_object>",
11             "apiVersion": "/v1beta2",
12             "name": "<name_of_object>",
13             "namespace": "<namespace_of_object>"
14         },
15         "timestamp": <timestamp_in_ISO_form>,
16         "value": "<value_of_metric>"
17     }
18 ]
19 }

```

The server must also serve this information under a defined path. There is no concrete list of path that a custom metric server must serve. There is [old proposal](#)<sup>8</sup> that is not up to date but on which it is possible to extract some information.

After investigation and tries, the custom metric server must serve this kind of path:

```

1 # The name in <> are correlated to the previous json description.
2 "/apis/custom.metrics.k8s.io/v1beta2/namespaces/<namespace_of_custom_metric_service>"
3 "/services/<name_of_custom_metric_service>/<name_of_custom_metric>"

```

With those information it is possible to create a custom metric server and compute our own metrics.

#### 4.3.1 Installation and tests

As discussed before, it is possible to setup a custom metric server than can be helpful if there is need to scale a kubernetes component regarding a queue length or load computed from something else than memory or cpu.

There is multiple steps to setup a custom metrics server:

First, create the server. It is possible to create a custom metrics server easily regarding what the research give us about what it must serve and on which path it must serve it. But, there is another thing: a service linked to the Kubernetes APIService (see ch. 4.2 to know how metric server works) must offer a certificate and TLS connection. It is possible to create self-signed certificate to give to our custom-metrics server.

Then, a service must be deployed to provide access to the custom-metrics server pods. The service must offer a https path.

After having a custom server running and a service that provide access to it, we need to register this service as provider for the K8s API: v1beta2.custom.metrics.k8s.io. This is done via deployment of an object of kind: APIService.

```

1 apiVersion: apiregistration.k8s.io/v1
2 kind: APIService
3 metadata:
4     name: v1beta2.custom.metrics.k8s.io
5 spec:

```

<sup>8</sup><https://github.com/kubernetes/community/blob/master/contributors/design-proposals/instrumentation/resource-metrics-api.md>

```

6   insecureSkipTLSVerify: true # Needed if using Minikube
7   group: custom.metrics.k8s.io
8   groupPriorityMinimum: 1000
9   versionPriority: 5
10  service:
11    name: custom-metrics-service
12    namespace: default
13  version: v1beta2

```

Finally, it is the time to setup a HPA based on our custom metric. To do this, instead of using the metric type "Resource" like in the chapter 4.2.2, we must use the type "Object" to configure exactly what we want. The "Object" kind contains multiple entries but the most interesting in a custom metrics server is the field named "describedObject". This field is the description of the custom metrics service (warning: service not server).

```

1  apiVersion: autoscaling/v2beta2
2  kind: HorizontalPodAutoscaler
3  metadata:
4    name: nodejs-hpa
5  spec:
6    scaleTargetRef: # Object on which the HPA will have effect on
7      apiVersion: apps/v1
8      kind: Deployment
9      name: nodejs
10   minReplicas: 1
11   maxReplicas: 30
12   metrics:
13     - type: Object
14       object:
15         target:
16           type: Value # use metric value, not average value
17           value: 12 # here the target value of the metric
18         describedObject: # description of the custom metrics service
19           apiVersion: v1
20           kind: Service
21           name: custom-metrics-service
22         metric: # metric to base the HPA on
23           name: kubemq-length

```

A custom metric server is available [here](https://gitlab.forge.hefr.ch/ps6-2021-kubemq/services/custom-metrics-server/-/tree/master)<sup>9</sup> and a complete deployment file for this custom-metrics server and HPA is available [here](https://gitlab.forge.hefr.ch/ps6-2021-kubemq/deployment/-/blob/e7a6c78602de8edec0b87f688cf28a5f8f24119/researches/custom-metrics-server.yaml)<sup>10</sup>.

If you want to setup the complete environment, you can execute these both commands:

```

1  # Pull the deployment repository of this project:
2  git clone git@gitlab.forge.hefr.ch:ps6-2021-kubemq/deployment.git
3
4  # Go to cloned project
5  cd deployment

```

<sup>9</sup><https://gitlab.forge.hefr.ch/ps6-2021-kubemq/services/custom-metrics-server/-/tree/master>

<sup>10</sup><https://gitlab.forge.hefr.ch/ps6-2021-kubemq/deployment/-/blob/e7a6c78602de8edec0b87f688cf28a5f8f24119/researches/custom-metrics-server.yaml>

```
6 |  
7 | # Deploy simple nodejs rest server+service (server available here:  
8 |   ↪ https://gitlab.forge.hefr.ch/ps6-2021-kubemq/services/nodejs-simple):  
9 | kubectl apply -f simple-rest-deployment.yaml  
10 |  
11 | # Deploy custom metrics server+service, APIService registration and HPA  
    ↪ for nodejs rest server:  
    kubectl apply -f custom-metrics-server.yaml
```

## 4.4 Kubernetes Jobs

[Kubernetes' Jobs](#)<sup>[7]</sup> is a task that creates various number of pods. The job tracks the number of pods successful completions and ends when a certain number of completion is reached.

So, the idea is that a job is created only when there is job to do or something to compute/process. It can be possible to link a job to a message queue and when a message appears on the queue, a job is triggered.

To do so, it is possible to setup a [CronJob](#)<sup>[8]</sup> that check for queue content and triggers Jobs when there is data inside it. This solution is proposed [here](#)<sup>11</sup>.

There is an example of implementation on the [K8s' website](#)<sup>[9]</sup> using a Redis queue.

K8s' Jobs might be a solution but required to transform the way an existing service works. This architecture is more suitable for parallels processing of a same data or processing chain of the data. More over, it takes each time a job is launched a bit of time to have it running. This can cause performance issue.

## 4.5 Sidecar pattern

After some research on [Istio](#)<sup>12</sup>, a service mesh for K8s, it is quasi impossible to not be confronted to Sidecar pattern as Istio used it to works.

Sidecar pattern consists of running multiple (generally lightweight and fast) containers along a main container inside a Pod. These sidecars are here to help the main container.

In the case of Istio, sidecars are injected automatically inside each pods along main container (remember: a Pod contains generally only one container). The injected sidecars work as proxy and intercept requests and redirect them to other proxies inside other pods to create a service mesh.

The sidecar pattern can be really useful in multitude of situation:

- Setting proxy in front of container.
  - For communication (as Istio do)
  - For logging (logging incoming request or outgoing response from a container)
  - For authentication and verification
  - For incoming data verification/cleanup/escaping
  - For caching

---

<sup>11</sup><https://stackoverflow.com/a/54045849>

<sup>12</sup><https://istio.io/>

- etc.
- Making an adapter for specific communication protocol
- etc.

## 4.6 Implementation and tests

A story on the website "Medium" has been written by "Bhargav Bachina" and is very useful as tutorial. This story is available [here](#)[10].

To create sidecars inside a deployment is really easy, it only needs to add more containers under the spec field of the template inside a deployment.

Below there is a deployment example of a main container named "nodejs" and its sidecar named "sidecar-container".

```
1 kind: Deployment
2 metadata:
3   name: nodejs
4   labels:
5     app: nodejs
6 spec:
7   selector:
8     matchLabels:
9       app: nodejs
10  replicas: 1
11  template:
12    metadata:
13      labels:
14        app: nodejs
15    spec:
16      containers:
17        - name: nodejs
18          image: agost/nodejs-starter:1.2
19          imagePullPolicy: Always
20        - name: sidecar-container
21          image: busybox
22          command: ["/bin/sh"]
23          args: ["-c", "while true; do echo echo \$(date -u) 'Hi I am from
    ↪ Sidecar container' >> /var/log/index.html; sleep 5;done"]
```

## 4.7 Forwarding / Redirecting requests

During the solution finding iteration process of the project, it comes the need of redirecting a request to another service for response. A service A (Forwarder) receive a request and answer directly by telling the client to go on another service (other URL pointing on service B: Receiver). See figure 4.



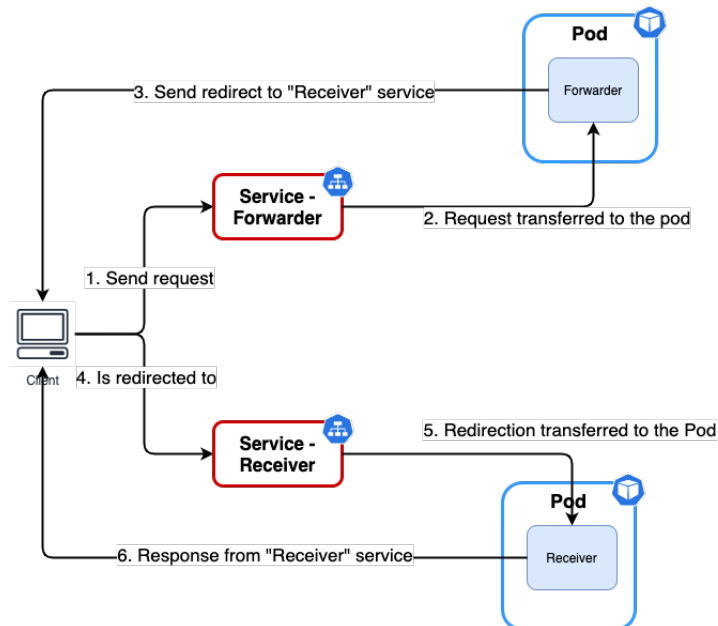


Figure 4: Redirecting workflow in K8s infrastructure

Additionally, the service "Receiver" must be capable of processing a lot of requests simultaneously (in parallel), so it requests to not have a "blocking API" that takes requests sequentially.

#### Implementation of a solution and tests:

For reminder: the project uses Java and Spring Boot.

For forwarding a request from a RestController in Spring Boot, it is possible to use the object "RedirectView" with the option "contextRelative" set to "false" because like this, the url put through the RedirectView is used as complete url and not additional path to the current relative address.

```

1 | @GetMapping("/")
2 | public RedirectView forward(HttpServletRequest req, HttpServletResponse
   |     resp) {
3 |     UUID requestUUID = UUID.randomUUID();
4 |     log.info("Forward request [" + requestUUID + "]");
5 |     return new RedirectView("http://localhost:8085/" + requestUUID,
   |         false);
6 | }
  
```

Now, as described above, the receiver must put each request in a separate thread. By putting the treatment each request into separated threads, we lose the symmetrical property of the communication. So it is mandatory to use a "DeferredResult" as result of the request. A "DeferredResult" object will trigger the ResponseServlet only when its method "setResult" is called. When this is done, the ResponseServlet will be notified and will respond to the request.

For putting the computation into a thread, instead of creating a Java Thread by using `new Thread()`, it is possible to use the ForkJoinPool used by Spring Boot. A ForkJoinPool acts like a big pool of jobs to run (normally used to fork and join threads). As Spring Boot use

the ForkJoinPool internally, when we put something in it, it will be executed when it is the best moment for Spring Boot. More information on ForkJoinPool [here](#)[11].

On the Java code bellow, we can see an implementation of the use of the ForkJoinPool along side the use of DeferredResult:

```

1  @GetMapping("/{id}")
2  public DeferredResult<String> getMethod(@PathVariable String id) {
3      DeferredResult<String> output = new DeferredResult<>();
4      log.info("Received request [" + id + "]");
5      ForkJoinPool.commonPool().submit(() -> {
6          log.info("Execution");
7          try {
8              Thread.sleep(5000);
9          } catch (InterruptedException e) {
10             e.printStackTrace();
11          }
12          output.setResult("voilà");
13      });
14      output.onCompletion(() -> log.info("Computation done"));
15      return output;
16  }

```

The both project (Receiver and Forwarder) used as test are available [here](#)<sup>13</sup>.

To test it with "curl", do not forget to use the "-L" parameter for following the redirect as written bellow:

```

1  | curl --request GET -L http://localhost:8086/

```

## 4.8 Passing K8s environment variable to Spring Boot container

This chapter discusses how the environment variables are shared and how they coexist. Finally, a test project with a simple implementation is presented.

### 4.8.1 K8s vs Docker vs Spring environment variables

In Kubernetes, there is the possibility to setup environment variables to a container and with Docker, it is possible to use environment variable when launching a container as the same as Spring Boot also support environment variables. Now the question is: How to pass environment variables from Kubernetes to Spring Boot via Docker?

The answer is: there is no such big deal.

As Spring Boot will use environment variables present on the machine it runs on, it will use the environment variables of the container.

When creating a Dockerfile, we have the possibility to setup environment variables that will be present on the container and set default values for them.

For each container inside a Pod, K8s allows to set environment variables that will be passed to the container. It works as the same way as if we launch the container on Docker using the "-env" parameter like this:

<sup>13</sup><https://gitlab.forge.hefr.ch/ps6-2021-kubemq/services/forwarder-receiver>

```

1 | # see:
2 | #https://docs.docker.com/engine/reference/commandline
3 |     # /run/#set-environment-variables--e---env---env-file
4 | docker run -e MYVAR1 --env MYVAR2=foo --env-file ./env.list ubuntu bash

```

The [official documentation](https://docs.docker.com/engine/reference/commandline)[12] of Docker says this: "Use the -e, -env, and -env-file flags to set simple (non-array) environment variables in the container you're running, or overwrite variables that are defined in the Dockerfile of the image you're running."

So it is confirmed that the environment variables created by Dockerfile will be overwritten by the environment variable declared in the K8s configuration.

On the figure 5 we have a simple representation of the environment variables and application.

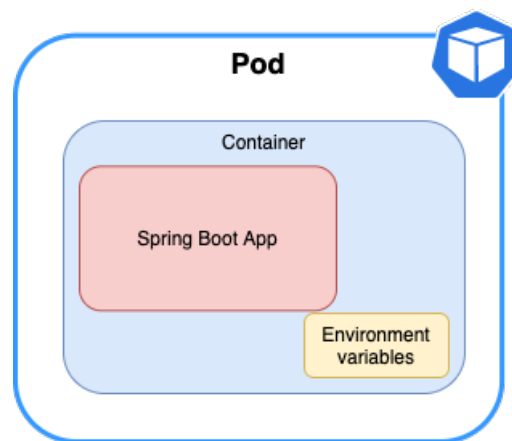


Figure 5: Spring Boot App and environment variables of the container

#### 4.8.2 Installation and tests

To transfer data from Kubernetes yaml file to a Spring Boot application, the app must succeed to retrieve the environment variables presents on the machine and use it in the code. To do so with Spring Boot, we can use Application Properties. Application Properties are the properties that we give to the application like the server port, the database url and so on. It is used to configure the application. There is a file named "application.properties" in each Spring Boot application under "/resources". Custom properties can be created here.

First, create custom properties inside the "application.properties" file and set these custom properties to environment variable name using the "\$<env\_var>:<default\_value>" form.

```

1 | kubemq.queue=${queue:null}
2 | kubemq.address=${queue_addr:null}

```

Now a configuration Spring component must be created to configure and interact with these properties inside our code. It is a class that is used only to store the value of these properties. It must have the same field names as the properties. Here a valid class for the both props declared in the "application.properties" file:

```

1  import lombok.Value;
2  import
   ↪  org.springframework.boot.context.properties.ConfigurationProperties;
3  import org.springframework.boot.context.properties.ConstructorBinding;
4
5  @ConfigurationProperties("kubemq")
6  @ConstructorBinding
7  @Value
8  public class KubeMQProps {
9      String address;
10     String queue;
11 }

```

Then, the custom properties configuration must be enable on the application to allow the configuration class to configure and access the properties. This is done by adding the annotation "@EnableConfigurationProperties" on the "Application class" and putting the configuration files for properties as parameter for the annotation. Here an coherent example for the both properties declared above:

```

1  @SpringBootApplication
2  @EnableConfigurationProperties({KubeMQProps.class})
3  public class PropertiesApplication {
4      public static void main(String[] args) {
5          SpringApplication.run(PropertiesApplication.class, args);
6      }
7  }

```

Now the Spring Boot application can interact with environment variables but there is not a single one in the container. This is time for adding environment variable via the Dockerfile. This can be done like this:

```

1  FROM openjdk:16-jdk-alpine
2  VOLUME /tmp
3  ARG JAR_FILE=target/*.jar
4  COPY ${JAR_FILE} app.jar
5
6  ENV queue "default_queue_name"
7  ENV queue_addr "default_queue_addr"
8  ENV port 8080
9
10 ENTRYPOINT ["java", "-jar", "/app.jar"]

```

Finally, when everything is set up, it only remains to set the environment variable in the K8s configuration file of the Pod/Deployment. The document on how to add environment variables to a container is provided at [this address](#)[13]. The simple way is to add the field "env" to container configuration and list under this field all the environment variables like this:

```

1  - name: env-var-spring-boot
2    image: agost/io.kueuer.proxy:latest
3    imagePullPolicy: Always

```

```
4 ports:
5   - containerPort: 3001
6 env:
7   - name: port
8     value: "3001"
9   - name: queue_addr
10     value: "kubemq-cluster-grpc.kubemq.svc.cluster.local:50000"
11   - name: queue
12     value: "s1-request"
```

## 4.9 Queuing with KubeMQ

KubeMQ documentation is available [here](#)<sup>14</sup>.

KubeMQ is a native Message Broker for Kubernetes. It is quite new with a lot of features. The big advantage of KubeMQ over other message broker is that it is very simple to use, it is fast and lightweight and native for Kubernetes. It also provides a complete CLI to help testing and developing.

Various SDK are available like .Net, Go, Node, Java, Python and Rest. More than queuing features, it allow to as connectors with KubeMQ Targets, Bridges and Sources. These functionalities are not useful for this project but are very interesting. More on that [here](#)<sup>15</sup>.

KubeMQ provide three different queuing patterns: CH 72 0026 0260 1180 7440. L

- Queues : Simple message queues (FIFO) with extended features.
- Pub/Sub : Publisher/Subscriber with event through channels.
- RPC : Remote Procedure Call pattern with commands and queries.

### 4.9.1 Queues

KubeMQ queues offers various useful features for a publisher or a consumer. For a publisher, he has the possibility to publish a single or batch messages to a queue. He can also setup Message expiration and Message delay (scheduled delay).

A consumer has the possibilities to:

- Make long polling (wait a certain amount of time on a message)
- Peek messages in the queue (without consuming them)
- Change the visibility of a message (consuming it, ack it or failed it)
- Ack all the messages in the queue
- Resend them

An interesting consumer feature is the Message Visibility. A consumer can set a visibility time to message and during this time the message will not be available for other consumers connected to the queue. The consumer can ACK the message during the visibility time and then the message will disappear definitely from the queue. If the consumer does not send an ACK during this period, the message is considered as "Failed" and it appears again for other consumers. The consumer can extend the visibility period at any time. This feature

<sup>14</sup><https://docs.kubemq.io/>

<sup>15</sup><https://docs.kubemq.io/learn/kubemq-connectors>

works like re-injection of the message on failure to process it. The figures 6 and 7 shows a graphical representation of this feature.

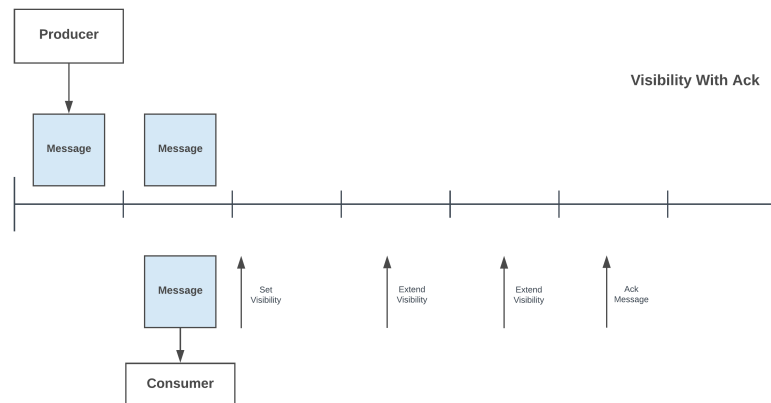


Figure 6: KubeMQ Visibility with Ack

Source: [KubeMQ.io](https://kubemq.io)

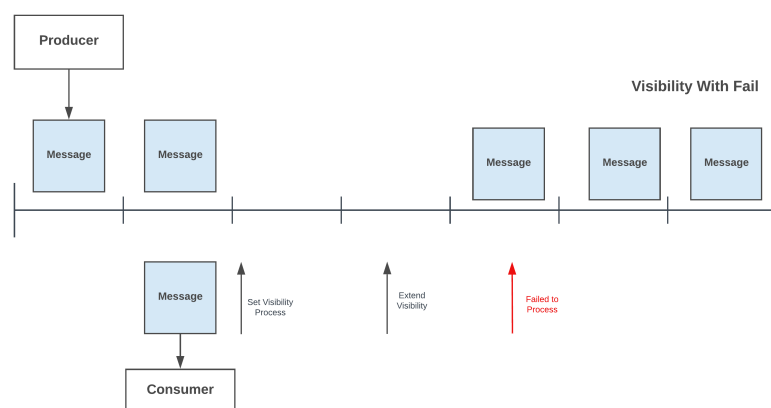


Figure 7: KubeMQ Visibility with Failure

Source: [KubeMQ.io](https://kubemq.io)

#### 4.9.2 Publisher / Subscriber

KubeMQ can be used to make a pub/sub pattern (figure 8). In KubeMQ, the messages are called "Events" and the queues are called "Channel" when it is related to pub/sub pattern.

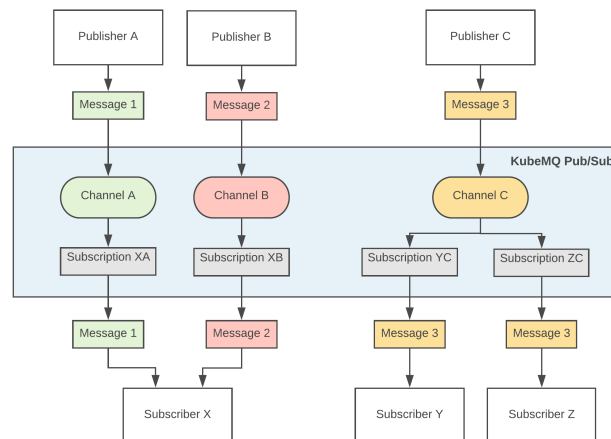


Figure 8: KubeMQ Basic Pub/Sub pattern

Source: [KubeMQ.io](https://kubemq.io)

Additionally, KubeMQ has another interesting feature that is the Event stores with Pub/Sub Pattern. This feature provides the possibility to store all the events into a channel dedicated database (figure 9). By storing all the events that are published into an Event Store, it allows a subscriber to read events that were published before the subscriber subscribe to the corresponding channel. With event stores, a subscriber has the possibility to read the events when it subscribes to the channel:

- From new events only.
- From the first event that were published into the channel.
- From the last event that were published into the channel.
- From sequence (from specific sequence)
- From Time (time in the past)
- From Time Delta (e.g. 30 sec back)

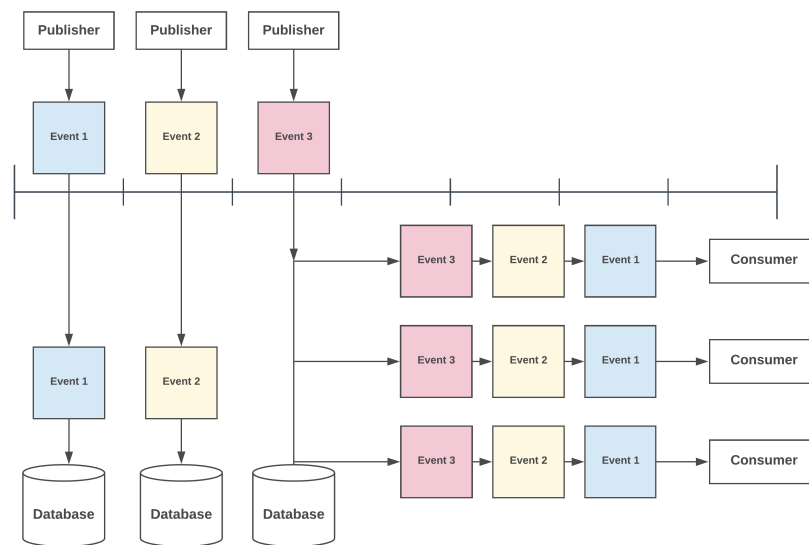


Figure 9: KubeMQ Pub/Sub with Event Stores

Source: [KubeMQ.io](https://kubemq.io)

## 4.10 Monitoring KubeMQ

KubeMQ provide "KubeMQ Dashboard" which is a pre-configured Grafana/Prometheus with all the metrics KubeMQ exports.

This deployment is available on [GitHub](https://github.com/kubemq-io/kubemq-dashboard)<sup>16</sup>.

There is the possibility of deploying Grafana/Prometheus simultaneously or either one or the other.

To deploy the whole package, kubemqctl can be used like this:

```
1 | kubemqctl get dashboard
```

To deploy only Prometheus:

```

1 | # Clone the project and move inside the directory
2 | git clone https://github.com/kubemq-io/kubemq-dashboard.git
3 | cd kubemq-dashboard
4 |
5 | # Apply the deployment file to the kubernetes environment
6 | # Warning, there is a typo in the repository
7 | kubectl apply -f promethues.yaml

```

The Prometheus configuration file with all the exported metrics is available in this [same repository](https://github.com/kubemq-io/kubemq-dashboard/blob/master/prometheus/kubemq.rule.yaml)<sup>17</sup>.

<sup>16</sup><https://github.com/kubemq-io/kubemq-dashboard>

<sup>17</sup><https://github.com/kubemq-io/kubemq-dashboard/blob/master/prometheus/kubemq.rule.yaml>



## 4.11 Reverse Proxies

This section discusses existing solution for queuing incoming requests.

It is important during the project to take a step back, see what we are doing and does it makes sense. The project is about queuing incoming requests to an application running in a Kubernetes infrastructure and scaling regarding metrics based on the queue.

A reverse proxy is a server placed before the application and used to route the requests from the client to the right application. There is several existing reverse proxies and some offers queuing features.

In the list of reverse proxies below, every proxy is capable of queuing incoming requests. They may be a potential solution for the project instead of creating a whole new system.

- NGINX is a well known reverse proxy, now part of the F5 company. It is capable of queuing incoming request to prevent requests from being dropped (see : [Advanced configuration for deploying in Kubernetes\[14\]](#)).
- F5 BIG-IP is used worldwide and is very powerful. It is also capable of queuing incoming requests as the manual documentation says[15].
- HAProxy is a load balancer/reverse proxy capable of high performance. HAProxy procedure for queuing incoming requests is presented in this publication [16]. It offers the possibility to set maximum requests accepted before queuing others.

In conclusion, solution already exists for queuing incoming requests which can be suitable for the project. Now, it remains to test in details and see if it offers the possibility to collect metrics based on their internal queues.

This check of existing solutions was made to late during the project and these solution have not been tested. It can be the next step if continuing this project.

## 5 Finding solution iterations

This section presents all the different iterations and researches/proposals of solutions made during the project. All of these proposals are linked to several tests described in chapter 4 to determine if it can be a suitable solution.

None of the solutions in this chapter is a suitable solution. To see why, check out the disadvantage listing of each solution.

### 5.1 Jobs

This is the first conception that was proposed. This solution (figure 10) works with Jobs and CronJobs (see ch. 4.4).

To queue incoming requests, it uses the idea of a "Dispatcher" that will works as proxy for the whole cluster. It intercepts incoming requests and regarding the service it was destined to, it will put the request into the queue corresponding to this service. There is one queue per service. The Dispatcher has a HPA that is based on the number of incoming requests per second to help handle a lot of incoming requests to the cluster.

For each queue, it will have a CronJob that will periodically check the length of the queue and if there is more than one element inside the queue, it will launch Jobs that will pull what's inside the queue, process it and put back the result. This solution does not use auto-scaler because it is the job of the CronJob to determine how many Jobs it must run.

With this solution, we say goodbye to the old implementation of existing service, and we convert them into executable jobs.

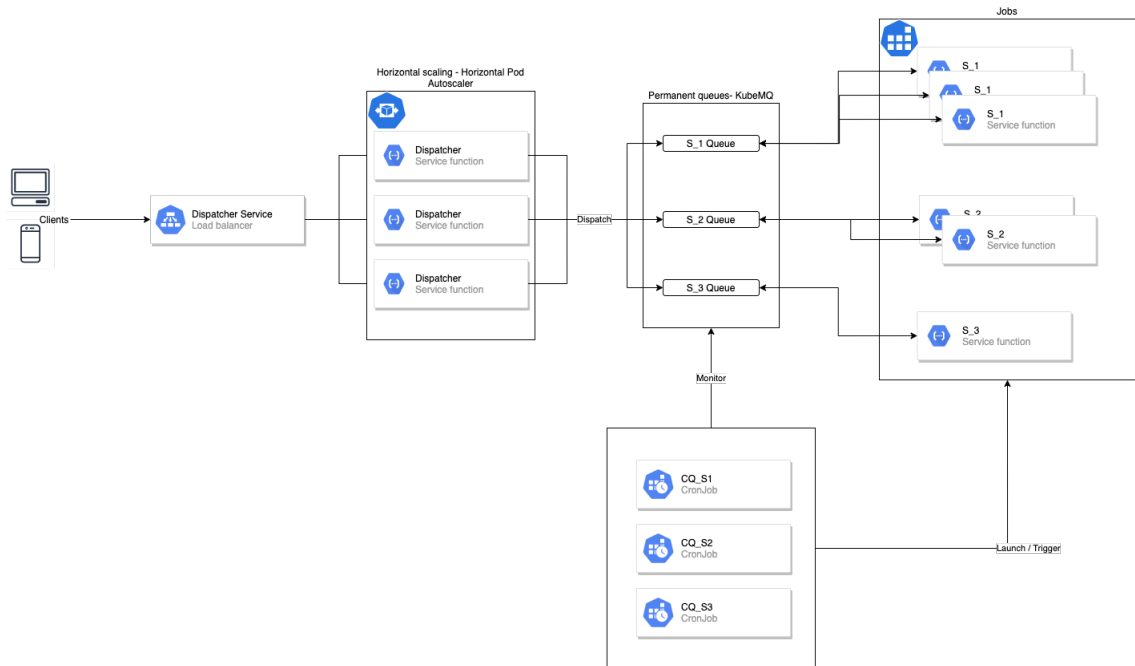


Figure 10: Jobs concept

There are disadvantages of using a "Dispatcher" pattern like this:

- How to intercept requests?
- There is no way of setting a proxy like this without changing the URL of each service. It requires not to access the service via its address but by the address of the dispatcher and so, modify each client to not call the address of the services anymore but the dispatcher instead.

Using Jobs have a lot of disadvantages:

- KubeMQ doesn't offer a way to get the queue length directly.
- All the existing services must be refactored to work as Job and not as basic containers anymore.
- CronJob can involve latency (not often enough checks) or load (too often checks) if they are not well tuned.

## 5.2 Custom metrics server

This other conception forgets the idea of using jobs because it is not suitable for the use case of the project.

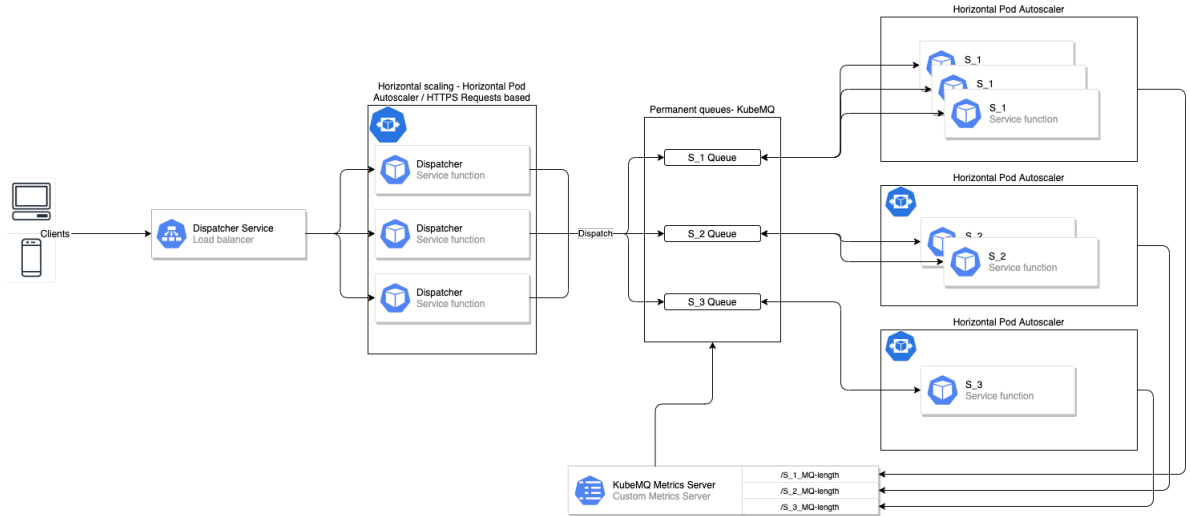


Figure 11: Custom metrics server concept

This concept (figure 11) still has the idea of a general "Dispatcher" which intercepts the request then dispatch them inside the corresponding queues of the service for which they were targeted.

The concept of this solution is to use a custom metrics server (see chap. 4.3) which is then used by all the HPA to scale up or down a service. For each queue, the custom metrics server serves a path that is used by the corresponding HPA. The custom metric served is named "kubemq-length" and it corresponds to the length of the queue.

In addition of the disadvantages and interrogations of the "Dispatcher" concept, this solutions have a lot of other disadvantages:

- KubeMQ doesn't offer a way to get the queue length directly.
- When a new service is added, the custom metrics server must be modified to serve a new path.
- Having a custom metrics server is not maintainable comparing of a Prometheus server.
- How existing services interact with the queues? They must be modified to work with queues.

### 5.3 Proxy and Adapter

This penultimate iteration is completely different form the previous ones.

This solution (figure 12) addresses the response confusion by creating a second queue dedicated to the responses.

It gets rid of the "Dispatcher" to replace it by a Proxy for each service. This way, no more needs to change the service path or clients to redirect to a dispatcher. The proxy must be invisible for the client. This proxy will intercept requests, put them in request queue and wait on the corresponding response appears in the response queue.

The most important thing is that this solution implements a sidecar (see ch. 4.5) named "QAdapter" that will communicate with the service and make the bridge between the service and the queues to avoid to change an existing service.

Finally, there is no more custom metrics server that are not maintainable and will use Prometheus (see ch. 4.10) for monitoring instead.

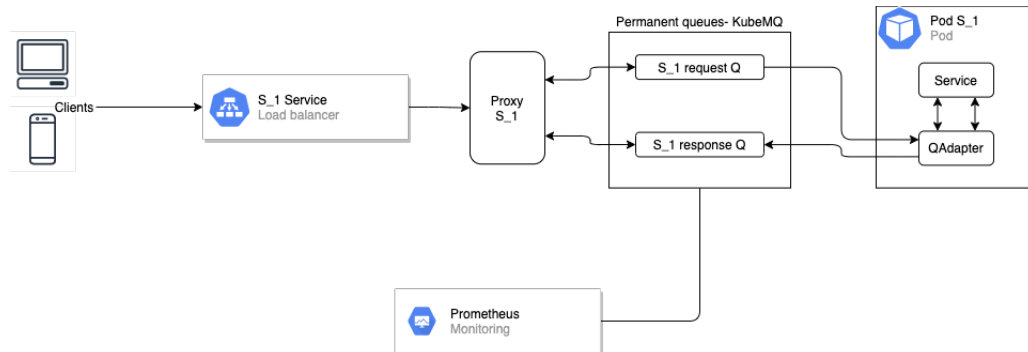


Figure 12: Proxy and Adapter concept

There are still some disadvantages to this solution:

- The Proxy must remember a list of requests in a table to know what to respond to whom.
- The Proxy works as two applications: one putting the data in the request queue and another one fetching continuously the data inside the response queue. This is not a good separation of work.

## 6 Kueuer

Kueuer is the name of the final solution proposal for this project. This final proposal is very similar to the one proposed in chapter 5.3(Proxy - Adapter).

This last proposed solution removes the disadvantages of the previous proposal (Proxy - Adapter) by dividing the Proxy into two different components but keeps the same idea.

### 6.1 Conception - Design

This chapter explains the general idea and the general concept behind the Kueuer system. It will not dive into the details for each parts. They all have their own chapter.

#### 6.1.1 Architecture

This solution is constituted by five parts:

- Proxy : intercept the request and put them into a request queue.
- QAdapter: adapter between the both queues and an existing application.
- Responder: Wait on responses published into the response queue to return them to the client.
- Queues: one queue for request (incoming data) and one channel for response (outgoing response).
- Monitoring and Scaling: Prometheus server for monitoring queues and HPA for scaling application (Pods).

All these parts will be discussed in detail in the following chapters.

The idea is that each existing application inside the Kubernetes environment will have its own Proxy, its own request and response queues. Additionally, each container will have its own QAdapter. The Responder can be used in two different configurations: 1. as application dependents (means one Responder per application like the Proxy) or 2. as cluster wide (it can be connected to several response Q of different applications)

The QAdapter is linked to a container (Sidecar pattern see ch. 4.5) unlike the Proxy which is linked to the K8s' service of the application.

On the figure 13 we can see the architecture in place for a single application K8s environment. This single application is named "S\_1".

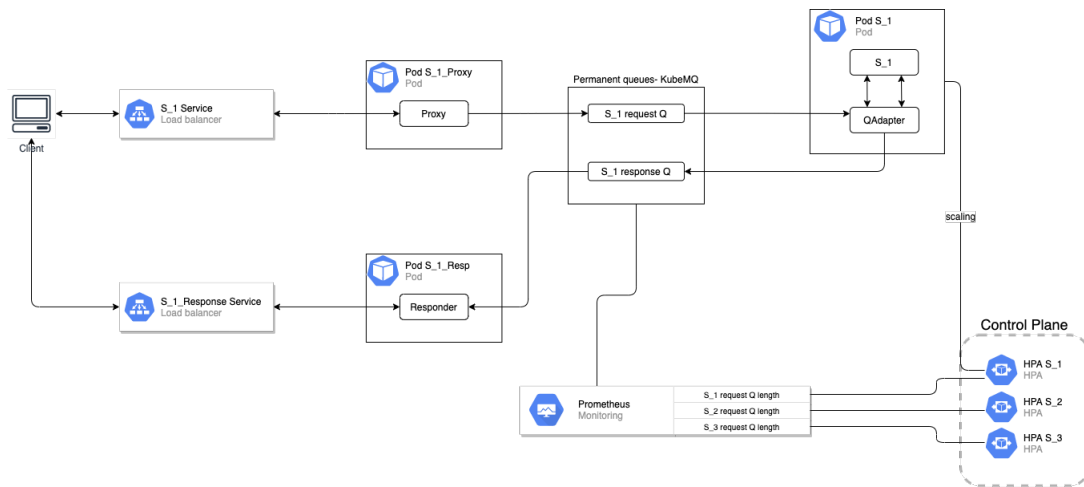


Figure 13: Kueuer architecture for one service S\_1

Each application has its own HPA (see chap. 4.2) that will automatically, horizontally (vertically vs horizontally: see ch. 4.1) scale up/down the application based on a shared Prometheus metrics server that collect metrics for each queues.

### 6.1.2 Communication

The easiest way to understand the communication between the components is to refers to the architecture schema above (figure 13) and imagining that the client wants to make a request to the "S\_1" application:

1. The client make a request as normal to the K8s' service of the application S\_1. No changes for him.
2. The K8s' service "S\_1 Service" redirect the request to the Proxy.
3. The Proxy "intercept" the request and put it in the queue "S\_1 request Q".
4. The Proxy send back a redirection link to the client. This redirection link points to the "S\_1 Responder Service".
5. The client follow the redirection to the Responder.
6. The Responder, which is subscribed to the "S\_1 response Q" wait until the corresponding response for the client's request appears, and so, the client wait with him.
7. Meanwhile, the QAdapter consume the request puttend into the "S\_1 request Q" by the Proxy (3)

8. The QAdapter make the call to the application "S\_1" and wait for the response.
9. When the QAdapter receives the response from the application, it put it in the "S\_1 response Q" and go back waiting on the "S\_1 request Q" to be fulfilled.
10. The Responder receives the response putted by the QAdapter and send it back to client.

On the more detailed sequence diagram below (figure 14), we can see that to make a relation between a request that were putted into the request queue and a response published into the response queue, the system use a unique identifier: one UUID. The Proxy will create a UUID for each request it intercepts and redirects the client on a path containing this generated UUID. This generated UUID acts like a token for the client, this token give him the right to access to its response via the Responder.

The UUID follow the request through all its processing, when the QAdapter put the response in the response queue, it uses the same UUID as the request.

Another important point, is that the Responder must be capable of handling numerous requests in parallel and must not have a "blocking API" (explained and tested in chapter 4.7) so it creates a thread each time it receives a request and this thread will just subscribe himself the internal observer and wait to be notified to send back the response to the client.

The Responder uses an internal machinery with internal observer to avoid subscribing multiple times to the same response queue. This observer is used to store all the response received (pub/sub pattern) and notify response threads waiting on a response (identification with UUID).

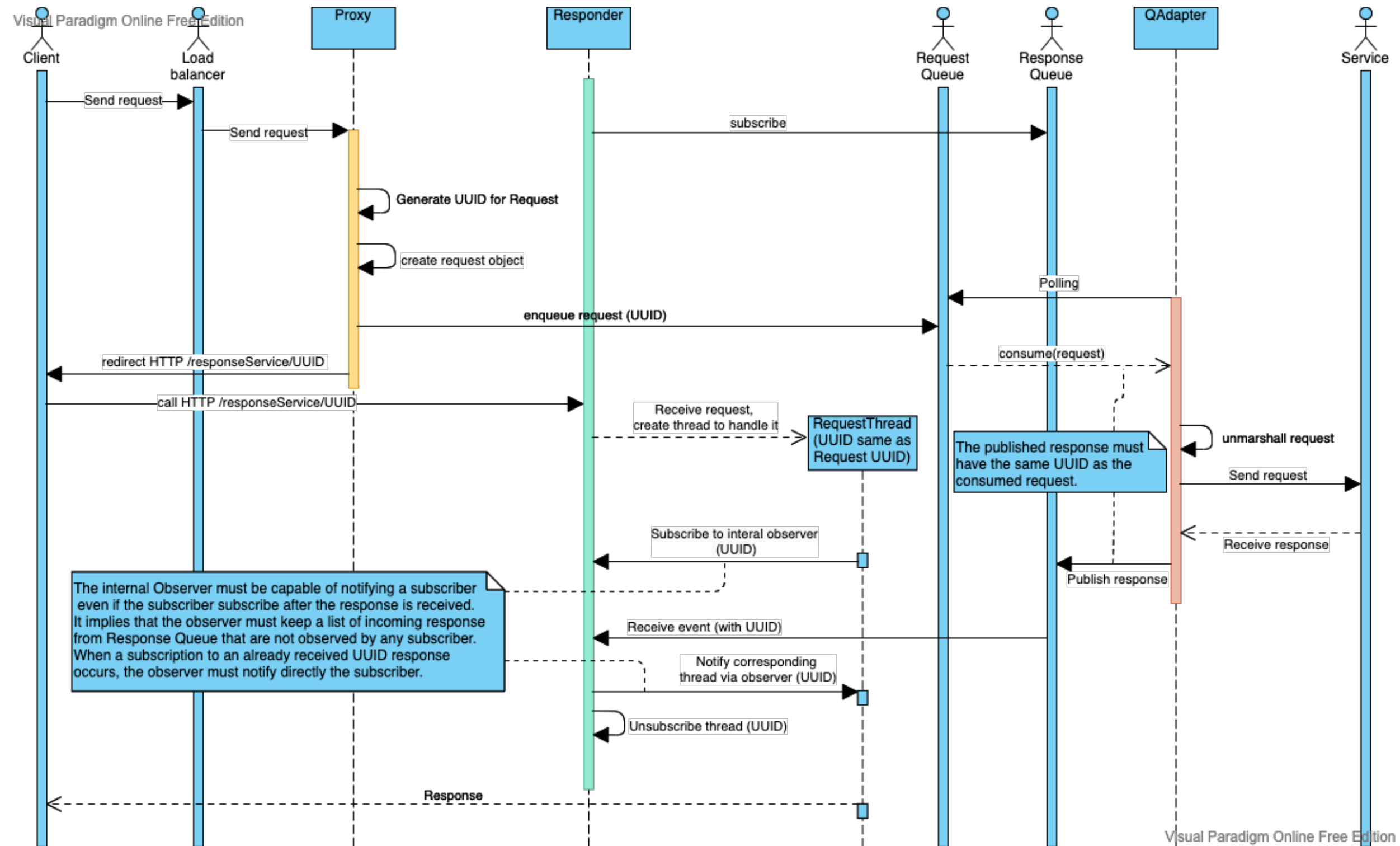


Figure 14: Sequence diagram of a client call on application named Service (actor on the right)

## 6.2 Proxy

The proxy component is used only for intercepting requests, converting them, put them into the request queue and redirecting the client to the Responder. It is a simple component having only one RestController.

### 6.2.1 Design

As said before, the Proxy component is a little component. It must:

1. Accept (Intercept) all kind of request.
2. Create a message corresponding to the request to put into the request queue.
3. Redirect to the Responder with the generated UUID.

As we can see on the class diagram of the application (figure 15), there is few component that are explained below:

- AppConfig + ResponderProps + KubeMQProps are used to configure environment variable to configure the proxy. The system used to retrieve environment variables from K8s config file to Spring Boot application is explained at chapter 4.8.
- ProxyInterceptor : RestController of the application that have only one method: "intercept". This method can handle all kind of request and return a RedirectView to redirect the user to the Responder. Same principle as explained in the chapter 4.7.
- Request : This pojo represent a request. It is the payload used to put into the request queue. It contains the uuid of the request (generated by the intercept method), the path, payload, incoming ip address and the HTTP Method of the request.

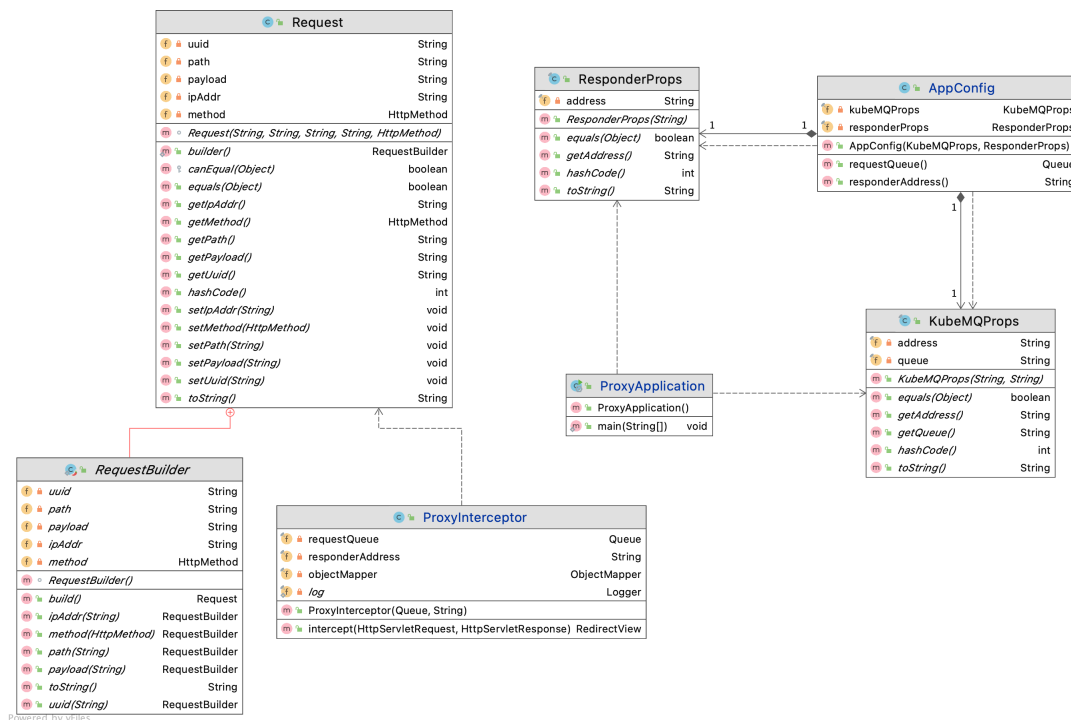


Figure 15: Proxy class diagram



### 6.2.2 Implementation

The implementation is not difficult, it is a mix of several technical tests described in these chapters:

- Forwarding-Redirect : ch. 4.7
- Environment Variables: ch. 4.8

The only interesting part of the implementation is how to put a message into a KubeMQ queue. All the documentation is available on the KubeMQ Java SDK on [GitHub\[17\]](#).

First a "Queue" object is needed. At the implementation is with Spring Boot, it is possible to create a single instance of a Queue and share it across the whole project using a Bean. As we can see on the class diagram (figure 15), the Bean is created inside the AppConfig class. When creating a Queue, it requires some parameters as the queue name, the clientId and the address of the queue:

```

1 | @Bean
2 | public Queue requestQueue() throws ServerAddressNotSuppliedException,
   |     ↳ SSLEException {
3 |     return new Queue(kubeMQProps.getQueue(), "proxy",
   |         ↳ kubeMQProps.getAddress());
4 | }

```

Then it is possible to create a message payload and push it into the queue using the "Message" object:

```

1 | Request request = Request.builder()
2 |     .path(req.getRequestURI())
3 |     .ipAddr(req.getRemoteAddr())
4 |     .method(HttpMethod.valueOf(req.getMethod()))
5 |     .payload("test")
6 |     .uuid(UUID.randomUUID().toString())
7 |     .build();
8 | Message msg = new Message()
9 |
   |     ↳ .setBody(Converter.ToByteArray(objectMapper.writeValueAsString(request)))
10 |     ↳ .setMetadata("empty");
11 | SendMessageResult res = requestQueue.SendQueueMessage(msg);

```

A Jackson Object mapper is used to convert the Request object into String. This provides a way to not have compatibility issues when un-marshalling the object on the QAdapter that uses the same object mapper.

### 6.2.3 Deployment

The proxy need to be configured when deployed. The service of the application it is related to must be changed to redirect requests to the proxy instead of the application.

Configuration files are available in appendix ch. 11.2.3

The Proxy component has three environment variables. Each environment variable annotated with "-" in the default value column are mandatory.

Name	Description	Example	default value
port	The port on which the Spring Boot application runs	3001	8080
queue_addr	The address KubeMQ grpc server. Use Kubernetes DNS.	kubemq-cluster-grpc.kubemq.svc.cluster.local:50000	-
responder_addr	The outside address of the responder	192.168.99.104:30001	-

### 6.3 QAdapter

The QAdapter is the only component connected to both queues simultaneously. It is a consumer and a publisher. The idea of the QAdapter is to avoid modifying existing application to run with queues. The adapter will pull the queued requests present into the request queue and make the corresponding HTTP request to the application running in the same Pod. When it receives the response for the HTTP request, it publishes the response into the response queue of the application. The QAdapter use the Sidecar pattern (ch. 4.5) so it runs along side an application.

#### 6.3.1 Design

The QAdapter works like this:

1. It pulls a request from the request queue.
2. It "process" the pulled request by creating a HTTP request to the application on the same Pod. This step is blocking, it means that the QAdapter cannot pull another request before the current one have not finished to process.
3. When it receives the response to its HTTP request, it publishes the response into the response queue.
4. After publishing the response into the response queue, it is again free to take another request from the request queue.

Below the class diagram of the QAdapter (figure 16):

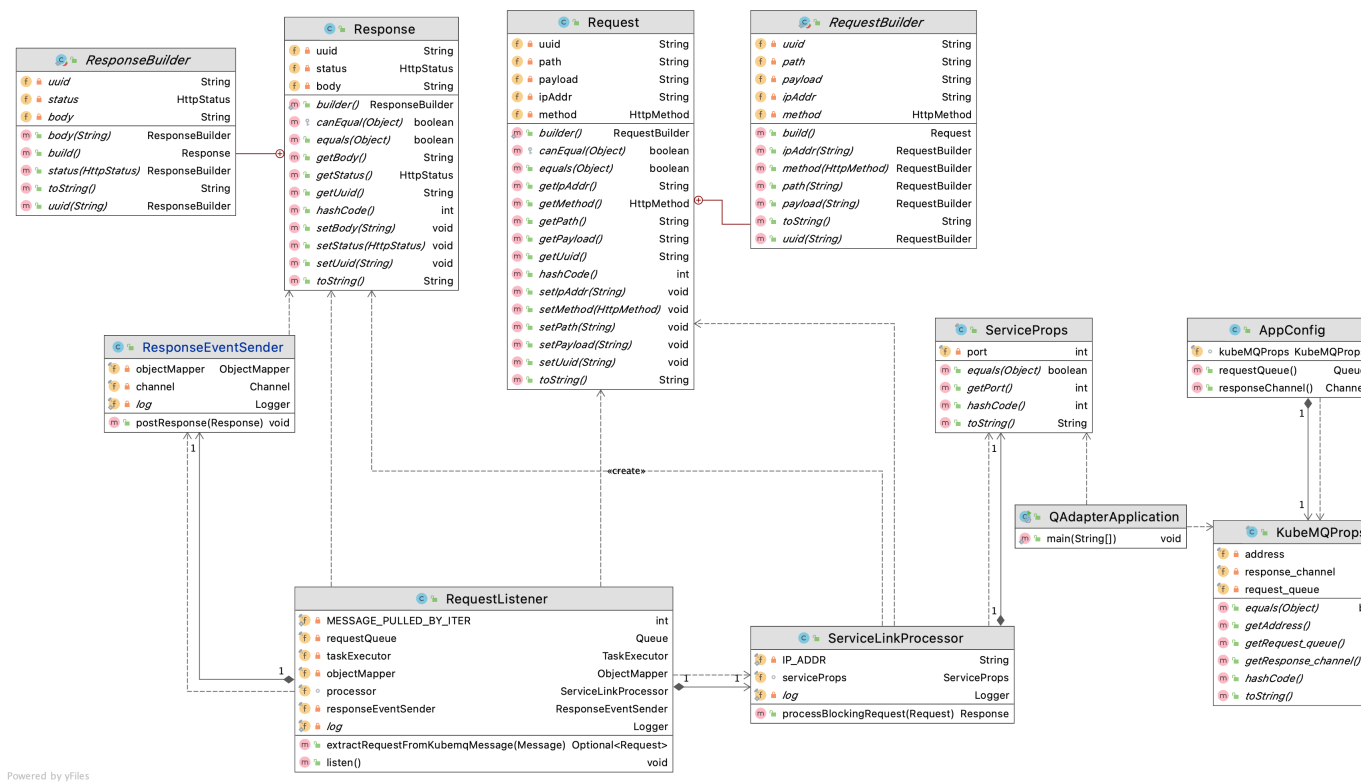


Figure 16: Proxy class diagram

On the class diagram we can see these classes:

- **ServiceProps + AppConfig + KubeMQProps** are used to configure and use environment variables given through K8s configuration of the container. (see ch. 4.8).
- **RequestListener**: this Spring Component is the core of the QAdapter. It is the listener connected to the request queue. It contains the method "listen" which is called automatically after initialization and contains the "while(true)" loop that check if something appears into the request queue and processes it.
- **ServiceLinkProcessor**: this Spring Component is used to process a request pulled from the request queue. It has only one method that process a Request object and return a Response. This method makes the call to the application the QAdapter is linked to. Its only one method is blocking and called by the method "listen" of the RequestListener when it finds something in the request queue. The fact this method is blocking forces the "while(true)" loop to stop during the processing and it avoids pulling other requests from the queue while already processing one.
- **ResponseEventSender**: this Spring Component is used to publish a Response object into the response queue. It is called by the method "listen()" of the RequestListener component after finishing processing a request.
- **Request + Response** are two POJO to store the data in objects.

The sequence diagram below (figure 17) shows graphically what is explained above:

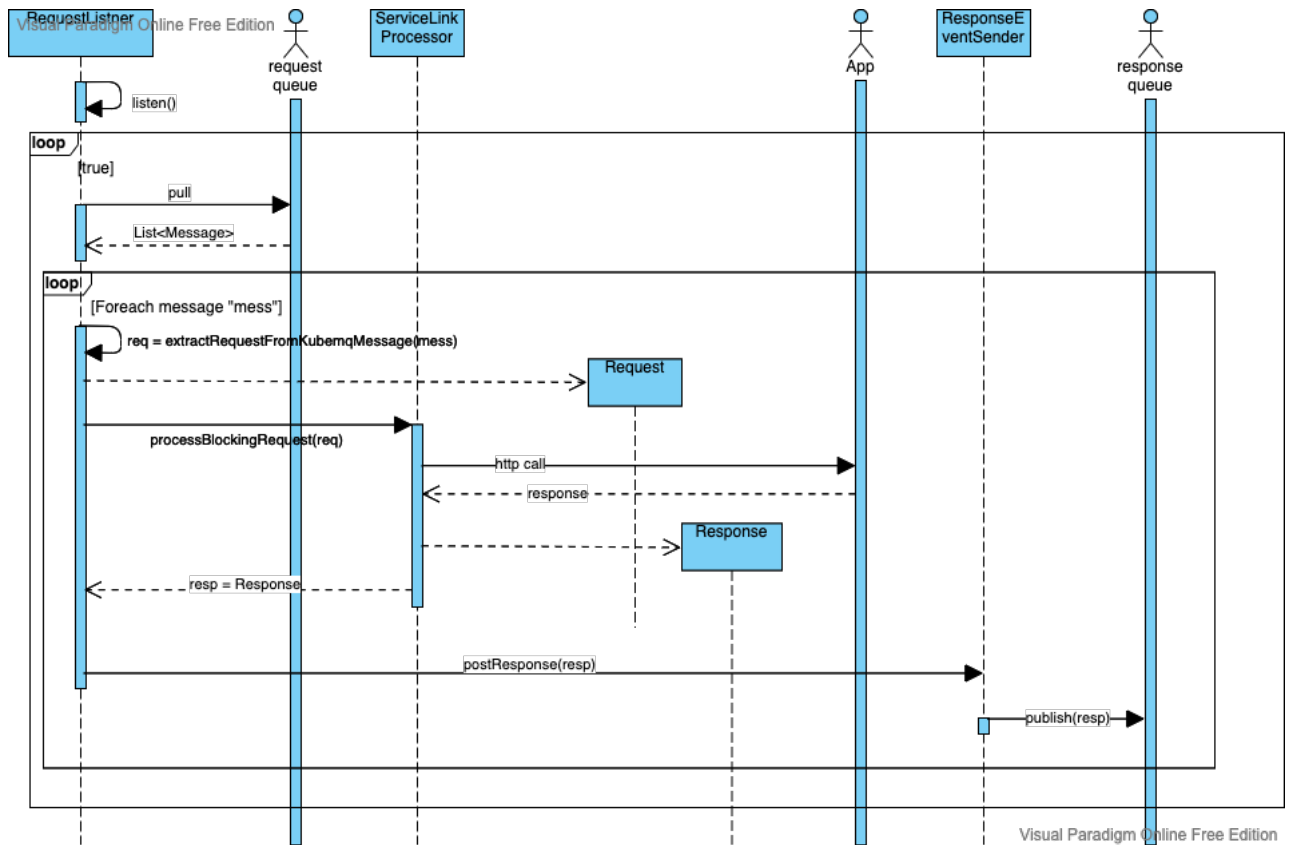


Figure 17: QAdapter sequence diagram

### 6.3.2 Implementation

There is no very complex implementation. There is just some point that might be important to understand:

First, how the method "ServiceLinkProcessor:processBlockingRequest" is effectively blocking: this method uses the Spring WebClient Object to send HTTP requests. There is multiple part of configuring the WebClient (see [documentation](https://docs.spring.io/spring-framework/docs/current/javadoc-api/org/springframework/web/reactive/function/client/WebClient.html)<sup>18</sup> or [tutorials](#)[18]) and after configuring a request and the client, it's time to ask for the response. Spring WebClient is by default an async client, but it is the possibility to ask to wait on the response using the method "block()":

```

1 // Declare the response here to be able to put the status code in it
2 // in the Mono parameterization
3 Response responseEvent = new Response();
4
5 Mono<String> response = req.exchangeToMono(clientResponse -> {
6     // put the status code in the response
7     responseEvent.setStatus(clientResponse.statusCode());
8     if (clientResponse.statusCode().is2xxSuccessful()) {
9         return clientResponse.bodyToMono(String.class);
10    } else if (clientResponse.statusCode()
  
```

<sup>18</sup><https://docs.spring.io/spring-framework/docs/current/javadoc-api/org/springframework/web/reactive/function/client/WebClient.html>

```

11         .is4xxClientError()) {
12             return Mono.just("Error response");
13         } else {
14             return clientResponse.createException()
15                 .flatMap(Mono::error);
16         }
17     });
18     String payload = response.block(); // <----- HERE

```

Also, as the QAdapter works as a sidecar, there is no need to know the address of the application is linked to because all application running in the same pods runs on the same machine. So the application linked to the QAdapter is reachable via the loopback address (127.0.0.1).

Then, it's also interesting to see how to consume message from a queue and how to publish an event to a channel. To consume a message from a queue using the KubeMQ Java SDK, the program needs to know the address of KubeMQ and the name of the queue. These information are given through the environment variables. After creating the Queue object as Spring Bean (see ch. 6.2.2), it is possible to pull message using the "ReceiveQueueMessages(numberOfPulled, wait\_time)" like this (the code above is inside the while(true) loop of the "RequestListner:listen" method):

```

1  ReceiveMessagesResponse res =
    ↳ requestQueue.ReceiveQueueMessages(MESSAGE_PULLED_BY_ITER, null); //null =
    ↳ default time
2  if (res.getIsError()) {
3      log.error("Cannot get a message in queue " + requestQueue.getQueueName() +
    ↳ " Error: " + res.getError());
4      continue;
5  }

```

For publishing event into a Channel, it just has to create an Event object, put body payload in it and send it with the "Channel:SendEvent" method. Example:

```

1  Event event = new Event();
2  event.setEventId(response.getUuid());
3  try {
4      ↳ event.setBody(Converter.ToByteArray(objectMapper.writeValueAsString(response)));
5  } catch (IOException ex) {
6      ↳ log.error("Cannot convert response object to byte[] for publishing event in channel {}");
    ↳ channel.getChannelName();
7      return; // TODO: Handle better than this
8  }
9  try {
10     channel.SendEvent(event);
11 } catch (ServerAddressNotSuppliedException | SSLException e) {
12     log.error("Error when trying to send event in channel {}");
    ↳ channel.getChannelName();
13 }

```

### 6.3.3 Deployment

As described before, the QAdapter works as sidecar. It must be integrated to the deployment of the application (explained in chapter 4.5).

Configuration file for deployment of an App with QAdapter attached to it is available in appendix ch. 11.2.1.

The QAdapter component has 5 environment variables. Each environment variable annotated with "-" in the default value column are mandatory.

Name	Description	Example	default value
port	The port on which the Spring Boot application runs	3001	8080
queue_addr	The address KubeMQ grpc server. Use Kubernetes DNS.	kubemq-cluster-grpc.kubemq.svc.cluster.local:50000	-
request_queue	The name of the request queue	s1-request	-
response_channel	The name of the response queue (=channel)	s1-response	-
service_port	The port of the main container of the pod.	8080	3000

## 6.4 Responder

The Responder component is used for two things. First it is used to handle a lot of pending requests from clients who wait on their response. Then, it is use to retrieve all the response from the response queues. We can compare it to a waiting room at the Post. You took a ticket and you received a number (here a UUID) and now you must go to the waiting room and wait until we call your number. The terminal on which you take the ticket is the Proxy component and the waiting room is the Responder component.

The Responder Component is the only component that can be cluster wide. You can imagine having only one Responder for the whole cluster and this Responder is connected to every response queues. If we take back the comparison with the Post waiting room, imagine having multiple Post offices with each their own terminal where you can take a ticket. All these Post offices share the same waiting room. This is possible because of the definition of the UUID: "A UUID is 128 bits long, and can guarantee uniqueness across space and time."<sup>19</sup>.

### 6.4.1 Design

The design of the Responder is done a way to limit the number of subscriptions to KubeMQ channels and to be capable of handling a lot of requests without blocking.

For the first constraints (limit the number of subscriptions), it has an internally observer pattern. This internal observer pattern allows the Responder to subscribe only once to each channel it must be connected to. The subject of the internal observer pattern is based on events coming from all the channels. This subject is a datastore of all events retrieved across the channels. Observers subscribe to this subject by giving the UUID they are waiting for. When a new event arrives into the subject, it checks if one observer is waiting on it, if yes, it gives the event to

<sup>19</sup><https://www.ietf.org/rfc/rfc4122.txt>

the corresponding observer. Regarding the uniqueness of a UUID, when an observer receives an event, it is automatically unsubscribed from the subject because it will never have another event with the same UUID. If there is no observer waiting on a event when the event arrives, the subject must be capable of storing the event until the corresponding observer comes. It can happen that an event arrives to the subject before the corresponding observer tries to subscribe because of communication delay or client error. For each event, it must have a corresponding observer.

For the second constraint (handling a lot of request without blocking), it uses the Java ForkJoinPool object to create a job each time a request arrives to the rest controller. These job work with DeferredResult to send the result back to client on a certain event. This implementation and concept are described in the chapter 4.7. To make the link with the internal observer pattern, the jobs created at each incoming request are the observers.

Finally, the Responder offers only one route and it is `"/<UUID>"`. The client must have a correct UUID to access to the server and a used UUID to access to data.

Warning, be aware that UUID do not protect against attacks or brute forcing (see [this article](#)[19]).

The Responder contains multiple classes (figure 18):

- **ResponseEventReceiver**: This class is used to automatically subscribe to the channels. When subscribing to a channel, it is mandatory to give a `grpc StreamObserver` which will consume the incoming events.
- **ResponseStreamObserver**: extends `grpc StreamObserver` and is used to consume events coming from the queues. When an event is received and has been correctly unmarshalled, it sends it to the `EventSubject`.
- **EventSubject**: The subject of the internal observer pattern.
- **ResponderController**: Rest controller serving the `"/<UUID>"` path. Add a job to the `ForkJoinPool.commonPool` when it receives a request. Its `"getResponse(String):DeferredResult<String>"` is the method which serves the `"/<UUID>"` path.
- **Response**: POJO for convertir an Event's body to a "Response".
- **KubeMQProps + AppConfig** are configuration for using environment variables.
- **CannotSubscribeException**: Java exception thrown when the `ResponseEventReceiver` cannot connect to a channel.

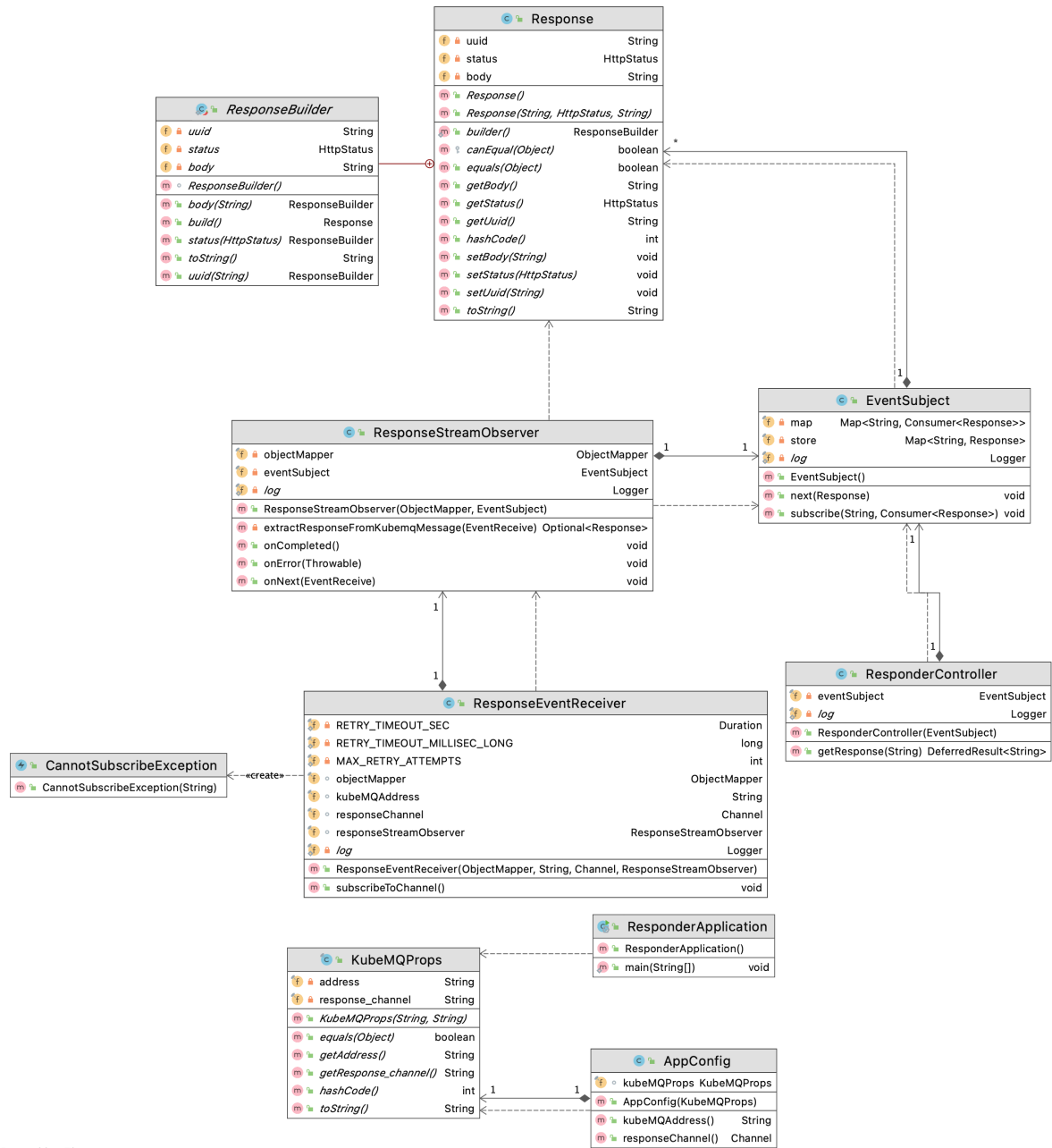


Figure 18: Responder class diagram

## 6.4.2 Implementation

For the controller, this is very similar to what is explained in the chapter about forwarding/redirecting requests (ch. 4.7). The creation of job is the same. The only difference is that a job consists of subscribing to the EventSubject and waiting for the response:

```

1 ForkJoinPool.commonPool().submit(() -> {
2     eventSubject.subscribe(uuid, response -> {

```



```

3         log.info("Putting response into deferred result");
4         deferredResult.setResult(response.getBody());
5     });
6 }

```

To implement the EventSubject, HashMap are used to store the observer and the event without observers. It uses HashMap because of the complexity to retrieve an information into Map that is generally  $O(1)$  and in the worst cases  $O(n)$ . This is important if we want to store and retrieve a lot of incoming requests.

The subject contains two HashMap: one for storing observers waiting on event and another one for storing events that have not yet an observer. You cannot unsubscribe from the subject because it is done automatically.

```

1 // Two maps:
2 private final Map<String, Consumer<Response>> map;
3 private final Map<String, Response> store;
4
5 public void subscribe(String uuid, Consumer<Response> responseConsumer) {
6     if (store.containsKey(uuid)) { // event already received and present in
7         ↪ store
8         log.info("Event already in store, executing..");
9         responseConsumer.accept(store.get(uuid));
10        log.info("Finish executing");
11        store.remove(uuid);
12    } else {
13        map.put(uuid, responseConsumer);
14    }
15
16 public void next(Response response) {
17     if (!map.containsKey(response.getUuid())) {
18         store.put(response.getUuid(), response);
19     } else {
20         map.get(response.getUuid()).accept(response);
21         map.remove(response.getUuid());
22     }
23 }

```

### 6.4.3 Deployment

The Responder can be deployed cluster wide because it can subscribe to multiple response queues simultaneously.

A configuration file for deploying the Responder component is available in the appendix chapter 11.2.4.

The Responder has 3 environment variables. Each environment variable annotated with "-" in the default value column are mandatory.

Name	Description	Example	default value
port	The port on which the Spring Boot application runs	3001	8080
queue_addr	The address KubeMQ grpc server. Use Kubernetes DNS.	kubemq-cluster-grpc.kubemq.svc.cluster.local:50000	-
response_channels	List of channel names separated by a comma ","	Only one: "s1-response" / two: "s1-response,s2-response"	-

The responder also need a service of type NodePort to be accessible from the outside. For example:

```

1  apiVersion: v1
2  kind: Service
3  metadata:
4    name: responder-service
5  spec:
6    selector:
7      app: responder
8    type: NodePort
9    ports:
10     - port: 8080
11       protocol: TCP
12       targetPort: 8080
13       nodePort: 30001

```

## 6.5 Queues

In the Kueuer system, queues are the central point. For queueing, KubeMQ is used.

The idea is simple: each application receives one queue and one channel:

- Request queue
- Response channel

For the request queue, it contains message transformed by the proxy. The producer is the proxy, who transform incoming request into message to put in queue and the consumer is a QAdapter connected to this queue (it may have more than one QAdapter if the POD is scaled up).

The response channel works with the publisher/subscriber pattern. The publisher is a QAdapter that finished to process a request and put the response in the response channel, and the subscribers are the multiple instances of the Responder component.

## 6.6 Monitoring and HPA

KubeMQ exposes metrics that can be used by a metric server as Prometheus.

The Kueuer system uses a custom Prometheus configuration to setup a Prometheus server. Below, a sample of the routes defined for the Prometheus server:

```

1 | - name: queues
2 |   rules:
3 |     - record: kubemq:messages:count:queues:send:sum
4 |       expr: sum(kubemq_messages_count{type="queues",side="send"}) OR on()
5 |         ↳ vector(0)
6 |     - record: kubemq:messages:count:queues:receive:sum
7 |       expr: sum(kubemq_messages_count{type="queues",side="receive"}) OR on()
8 |         ↳ vector(0)
9 |     - record: kubemq:messages:count:queues:sum

```

These routes allow to use pre-configured question to ask the Prometheus server.

### 6.6.1 Prometheus Adapter

Prometheus is a server capable of retrieving metrics from other services but it is not able to serve them the way a HPA configuration can use.

A HPA (Horizontal Pod Autoscaler) need to interact with the basic metrics server, the custom metrics server or the external metrics server defined in the K8s' API.

Prometheus is not suitable to be configured directly as custom or external metrics server.

This is why we use a component named "Prometheus-Adapter". This Prometheus adapter acts as external or custom metrics server (it implements the specification of the K8s custom/external metrics server) and is also capable of asking a Prometheus server to retrieve metrics.

In the case of the Kueuer system, the Prometheus Adapter will act as **external metrics** server because we want to base a HPA on metrics coming from another pod (metrics don't come from the pod the HPA is the controller but from another pod, here KubeMQ).

### 6.6.2 HPA

With the Prometheus Adapter, it is possible now to set up a HPA based on external metrics. A HPA can be based on multiple metrics at the same time and so, it is possible to declare several metrics in the Prometheus Adapter and base the HPA on all these metrics.

### 6.6.3 Implementation

The most difficult part of the implementation of working HPA based on KubeMQ metrics is the configuration. The most time consuming to understand is the Prometheus adapter configuration. Later, to install Prometheus Adapter, the package manager "HELM" is used.

The configuration discussed here is the helm configuration of the community version of the Prometheus Adapter available [here](#)<sup>20</sup>.

A whole explanation (see [Configuration Walkthrough](#)[20]) about the configuration's options of the Prometheus Adapter is available on the [official repository](#)<sup>21</sup>.

```

1 | prometheus:
2 |   url: http://prometheus-service.kubemq.svc
3 |   port: 8080
4 |   rules:

```

<sup>20</sup><https://artifacthub.io/packages/helm/prometheus-community/prometheus-adapter>

<sup>21</sup><https://github.com/kubernetes-sigs/prometheus-adapter>

```

5 external:
6   - seriesQuery: '{__name__=~"kubemq_messages_count.*"}'
7     resources:
8       overrides:
9         namespace:
10           resource: namespace
11         pod:
12           resource: pod
13       name:
14         matches: ^(.*)
15         as: "queue-length-increase-rate"
16     metricsQuery: rate(sum(kubemq_messages_count{<< range $key, $value :=
    → .LabelValuesByName >><< if ne $key "namespace" >><< $key >>="<< $value
    → >>",<< end >><< end >>})[10m:])

```

First thing to see in this configuration is the Prometheus server address.

Then we have a "rules" section containing several rules. Each rule is divided in 4 parts:

1. The Metric discovery (discover the metric available on the Prometheus server. Here it will check for all metrics available on the Prometheus server beginning with "kubemq\_message\_count." prefix). SeriesQuery
2. Association. Used to decide to which element in the Kubernetes environment this metric is available. Here, the metric will be available for POD elements and Namespaces. Resources section
3. Naming (define with which pattern the metrics must be served). Here the metric is served with the name "queue-length-increase-rate".
4. Querying. The query to execute on the Prometheus server to retrieve the metric.

On the first implementation of the configuration, I created a list of exact same queries with each time another property set in the kubemq\_messages\_count(<here>). This property changed regarding the queue we are looking for. This method was not good and another solution has been implemented: let the HPA select with a field on which queue it wants to retrieve metrics.

#### 6.6.4 Deployment

First, the Prometheus server must be deployed. To do so, a already existing configuration to retrieve the metrics from KubeMQ has been created and is available [here](#)<sup>22</sup>.

In the current state of the project, it is needed to pull the "Deployment" repository and apply the file "Promethues.yaml" chap. 11.3.1

```

1 | git clone git@gitlab.forge.hefr.ch:ps6-2021-kubemq/deployment.git
2 | cd deployment
3 | kubectl apply -f kubemq-dashboard/promethues.yaml

```

Now that Prometheus is running, it is time to setup the Prometheus adapter. To do so, it is more convenient to use "helm" as the [official repository](#)<sup>23</sup> describes. A configuration file has already

<sup>22</sup><https://gitlab.forge.hefr.ch/ps6-2021-kubemq/deployment/-/blob/29be92af3ed8f789406e5b3b37ca2b1ca5cbdcee/kubemq-dashboard/promethues.yaml>

<sup>23</sup><https://github.com/kubernetes-sigs/prometheus-adapter>

been implemented (chap. 11.3.2) in the deployment/kubemq-dashboard folder can be used like this:

```

1 | cd deployment
2 | # Ajout du repo helm
3 | helm repo add prometheus-community
   | → https://prometheus-community.github.io/helm-charts
4 | # Création
5 | helm install your_release --namespace=<namespace> -f
   | → kubemq-dashboard/prometheus-adapter-config.yaml
   | → prometheus-community/prometheus-adapter

```

A HPA is based on External metrics and each external metrics must contain in their "selector" properties these information: "channel=s1-request,type=queues,side=send"

Name	Description	Example	default value
channel	the KubeMQ queue on which the HPA must be based	s1-request	-

in addition, the HPA can set every KubeMQ properties for selection of metrics it wants. All the properties are listed on [KubeMQ documentation](#)[21].

Here an example of a configured HPA to use two different metrics exported by the Prometheus Adapter:

```

1 | apiVersion: autoscaling/v2beta2
2 | kind: HorizontalPodAutoscaler
3 | metadata:
4 |   name: s1-hpa
5 | spec:
6 |   scaleTargetRef:
7 |     apiVersion: apps/v1
8 |     kind: Deployment
9 |     name: s1
10 |   minReplicas: 1
11 |   maxReplicas: 10
12 |   metrics:
13 |   - type: External
14 |     external:
15 |       metric:
16 |         name: queue-length-increase-rate
17 |         selector: "channel=s1-request,type=queues,side=send"
18 |       target:
19 |         type: Value
20 |         value: 0.2

```

The final HPA used is available in the appendix chapter 11.2.2.

Warning: When setting HPA, do not forget to set limits in the concerned pod.

## 6.7 Evolution

There is a lot of possible evolution for this project, but two are the most important. The following sections describe the both main evolution/improvements to make to this system.

### 6.7.1 Automation / Injection

One of the main goals to a project like this is to remove all the configuration and make it a lot more simpler to deploy.

Istio (Service Mesh) can be taken as example. There is quasi no-need of heavy configuration to launch Istio. It automatically deploys proxy, sidecars and so on.

The first main evolution is to develop a system capable of injecting proxy component like Istio make and like this, the user will not have any configuration to do. Everything can be automated.

To follow this first possible evolution, here a list of maybe useful links:

- Explanation of the creation of a sidecar injector: <https://engineering.salesforce.com/a-generic-sidecar-injector-for-kubernetes-c05eede1f6bb>
- Open-source sidecar injector: <https://github.com/tumblr/k8s-sidecar-injector>
- Google explaining how to use Istio as sidecar injector: <https://cloud.google.com/service-mesh/docs/proxy-injection>
- Kubernetes Admission Controllers: <https://kubernetes.io/docs/reference/access-authn-authz/admission-controllers/>
- Kubernetes Dynamic admission control: <https://kubernetes.io/docs/reference/access-authn-authz/extensible-admission-controllers/>

### 6.7.2 Language Refactoring

When this project began, I (Romain Agostinelli), didn't know Kubernetes and the language chosen for developing solution was Java as it is a well know language and typed language.

It is during the development/project that I figured out that the language used everywhere for Kubernetes is Go (it makes sense as Go was developed inside Google).

A big improvement would be to re-code the entire project with Go to make it compatible with other existing components for Kubernetes.

The Go language has been defined as the default language to works on Kubernetes by Google and the community.

The more early this refactor will be made, the more easily it will be made. For now, the applications are not complexes.

## 7 Tests

The final stage of the project consists of testing the solution developed.

This chapter discusses the tools used for testing, the two kinds of test done with results, what issues/bugs the the tests shows up and finally a synthesis/conclusion.

As a reminder, the goal of the project is queueing incoming requests to a service inside a Kubernetes infrastructure and being capable of scaling the service regarding queue metrics.

## 7.1 Tools

For testing the solution developed, a lot of existing tools exist and are very useful.

### 7.1.1 Siege

A first one, developed especially for testing REST services, is siege. Siege is an open-source utility that can perform load testing. It is used via a terminal, and configuration is done via command options. Here an example of the use of siege that generate 10 parallel threads during 10 min with other output configurations:

```
1 | siege -q -v -c 10 -b --log=./siege.log -j -t 10m http://127.0.0.1:62569
```

The official repository for siege is accessible [here](#)<sup>24</sup>.

### 7.1.2 K6

K6 is well known tool for load testing. It is written in Go and allows to configure load tests with virtual users and complex user simulation.

It uses a NodeJS interface (facade) to setup tests but then use a core written in Go for performance.

It is really easy to use and offers a lot of functionalities, like virtual users, staging, checks and so on.

The official documentation is available [here](#)<sup>25</sup>.

The two kinds of tests made at the end of the project are using K6 as testing tool, utility.

## 7.2 Tests and results

This chapter 756.1579.7189.30 shows the two kinds of test that were made, with their configuration and procedure.

### 7.2.1 Environment simulation

For testing, a simulation environment has been setup and the tests are again one service randomly waiting between 0 and 30 seconds. There is no need to simulate CPU consumption because the goal is to scale based on queues metrics and not on resources used.

The queue metric used for auto-scaling is the queue incoming messages instant rate. This metric show a good status of number of incoming requests treated by the proxy.

The environment infrastructure is the same as presented on the figure 13 (chapter 6.1.1) and is based on one service.

The K8s cluster has these resources:

Resources	Amount
CPU	8 cores
Memory	10240 MB

<sup>24</sup><https://github.com/JoeDog/siege>

<sup>25</sup><https://k6.io/docs/>

To setup the environment, every "yaml" files are available on GitLab in the deployment project under the "kueuer" directory. The deployment project is available [here](#)<sup>26</sup>.

Below, the directory tree with what to deploy first:

```

1 |kueuer
2 | |-- merged_configs
3 | | |-- app-kueuer-proxy-outside.yaml
4 | | |-- app-kueuer-proxy-sidecar.yaml
5 | |-- metrics
6 | | |-- bck
7 | | | |-- prometheus-adapter-config-backup.yaml
8 | | |-- prometheus-adapter-config.yaml
9 | | |-- prometheus.sh #<----- 2nd, install metrics env. on the K8s
   |   ↳ cluster
10 | | |-- promethues.yaml
11 | |-- minikube-setup.sh #<----- 1st, install minikube env. and KubeMQ
12 | |-- test_environment
13 | | |-- app-hpa.yaml
14 | | |-- app-service.yaml
15 | | |-- app-with-qadapter.yaml
16 | | |-- bck
17 | | | |-- hpa-conf-backup.yaml
18 | | |-- deploy-sample-app.sh #<----- 3rd, deploy a simple Kueuer
   |   ↳ infrastructure
19 | | |-- proxy-outside.yaml
20 | | |-- responder-hpa.yaml
21 | | |-- responder-service.yaml
22 | | |-- responder.yaml
23 | | |-- siege.log
24 | |-- tests #<----- 4th, try launching a K6 test.
25 | | |-- script-func.js
26 | | |-- script-stress.js

```

For deployment, please open the ".sh" files instead of executing them, they contain instruction and listing of commands to help installing the environment.

## 7.2.2 Threshold

In the automated testing tool K6, there is the possibility to set threshold. Threshold are the minimum requirements for a test to pass. In the tests described in the sections bellow, there is one threshold set and it is based on the metric "http request failed".

For validating the test, the rate of http requests failed must be under 0.01.

```

1 | thresholds: {
2 |   http_req_failed: [{ threshold: 'rate<0.01', abortOnFail: false }],
3 | },

```

<sup>26</sup><https://gitlab.forge.hefr.ch/ps6-2021-kubemq/deployment>



### 7.2.3 Functional testing

First of all, the functional testing: does everything works?, is there issues / bugs?

This first part of the tests is checking main goals of the project as:

- Queueing incoming requests
- Auto-scaling based on metrics from queues

and also find and report possible bugs / issues.

To do so, the automated test uses 3 stages:

1. Normal amount of incoming requests -> Nothing to scale up/down
2. Increasing amount of incoming requests -> Must scale up regarding queues metrics
3. Bellow normal amount of incoming requests -> Must scale down regarding queues metrics

In K6 words, these three steps are translated to:

```
1 | stages: [  
2 |   { duration: '1m', target: 5 }, // normal load  
3 |   { duration: '40s', target: 20 }, // autoscaler increase  
4 |   { duration: '1m', target: 5 }, // autoscaler down  
5 | ],
```

The complete configuration used for the functional test is available at chapter 11.4.1.

The goal is to show if the environment works by using it normally, stressing a bit and then decrease the number of incoming requests. Here, the "target" value do not signify requests/seconds but number of concurrent virtual users making calls to the service.

#### Results:


The tests have been run several times and always pass, the pods are scaled up correctly and scaled down too.

Here an example of a K6 test result (figure 19):

```

agost@BunsenLabs:~/workspace/PS6/deployment/kueuer/tests$ k6 run script-func.js

```



```

execution: local
script: script-func.js
output: -

scenarios: (100.00%) 1 scenario, 20 max VUs, 3m10s max duration (incl. graceful stop):
* default: Up to 20 looping VUs for 2m40s over 3 stages (gracefulRampDown: 30s, gracefulStop: 30s)

running (3m10.0s), 00/20 VUs, 12 complete and 18 interrupted iterations
default ✓ [=====] 05/20 VUs  2m40s

data_received.....: 6.8 kB 35 B/s
data_sent.....: 7.4 kB 38 B/s
http_req_blocked.....: avg=1.12ms  min=0s    med=718.34µs  max=4.67ms  p(90)=3.15ms  p(95)=3.48ms
http_req_connecting.....: avg=1.05ms  min=0s    med=628.75µs  max=4.56ms  p(90)=3.09ms  p(95)=3.42ms
http_req_duration.....: avg=12.45s  min=11.04ms med=82.1ms    max=1m45s   p(90)=34.84s  p(95)=1m6s
{ expected_response:true }...: avg=12.45s  min=11.04ms med=82.1ms    max=1m45s   p(90)=34.84s  p(95)=1m6s
✓ http_req_failed.....: 0.00% ✓ 0    ✗ 42
http_req_receiving.....: avg=111.66µs min=69.04µs med=109.1µs    max=180.59µs p(90)=147.68µs p(95)=155.58µs
http_req_sending.....: avg=77.7µs   min=17.72µs med=83.85µs    max=173.03µs p(90)=118.34µs p(95)=123.03µs
http_req_tls_handshaking.....: avg=0s       min=0s    med=0s         max=0s       p(90)=0s      p(95)=0s
http_req_waiting.....: avg=12.45s   min=10.92ms med=81.95ms    max=1m45s   p(90)=34.84s  p(95)=1m6s
http_reqs.....: 42 0.22105/s
iteration_duration.....: avg=44.49s   min=22.11s med=32.78s     max=1m46s   p(90)=1m16s   p(95)=1m30s
iterations.....: 12 0.063157/s
vus.....: 0 min=0 max=20
vus_max.....: 20 min=20 max=20

```

Figure 19: K6 final result for functional test

On the result showed above (figure 19), we can see that the average response time is 12.45s what is accurate and what we can expect for a service waiting between 0 and 30 seconds. Also, the maximum time is 1m45s, this happened when entering the stressing stage of the test, when they were not enough replicas to handle the request.

On the Prometheus server (figure 20), we can see that incoming messages instant rate metric reacts very well corresponding to the 3 defined stages and we can see each stages.

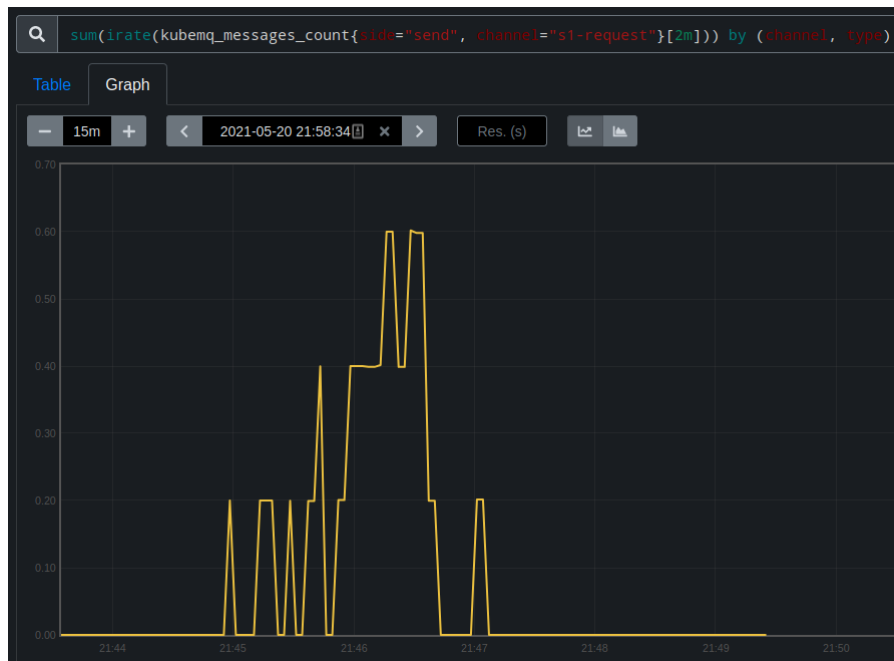


Figure 20: Prometheus metrics evolution for functional test

### 7.3 Stress test

The second kind of test made is a stress test. A stress test consists of testing the environment with a lot of incoming requests and check its behavior. For the stress test, a lower time consuming fake service is used due to resource limitation. It is possible to imagine scaling until 100 replicas of the same pod in big infrastructure with a lot of resource, but not with the environment used for these tests.

Here, the maximum amount of replicas allowed for the fake service is 30 and the threshold still the same. The goal of this stress test is to see if the queue can make a good buffer for the incoming requests rates while the auto-scaler increase the number of pods.

The test consists in 7 stages with a spike to 100 simultaneous virtual users trying to access the service. Below, the described stages in K6:

```
1 stages: [
2     { duration: '10s', target: 10 }, // below normal load
3     { duration: '1m', target: 10 },
4     { duration: '10s', target: 100 }, // spike to 100 users
5     { duration: '3m', target: 100 }, // stay at 100 for 3 minutes
6     { duration: '3m', target: 20 },
7     { duration: '7m', target: 10 }, // down scaling time
8     { duration: '40s', target: 50 }, // stress
9 ],
```

The complete configuration used for the stress test is available at chapter 11.4.2.

### Results:

On the figure 21 bellow, we can see the K6 result for this test:

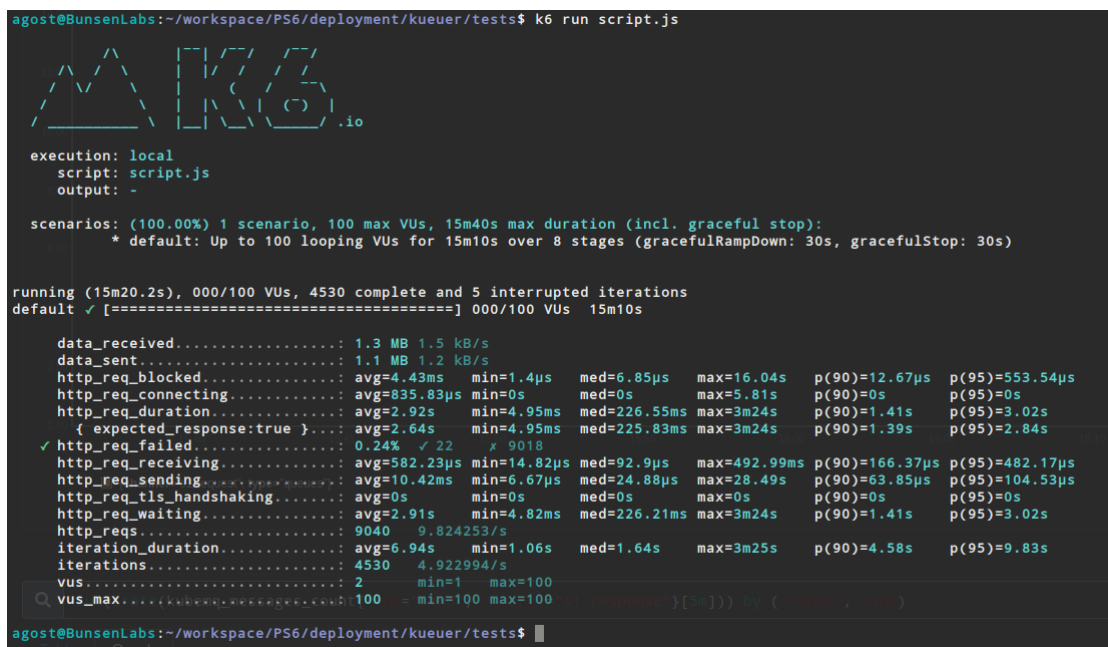


Figure 21: K6 final result for stress test

We can see that the system is capable of receiving a lot of incoming requests (regarding the environment's resources) and the threshold of requests failed rate passes. There is 0.24% of requests lost, due to different factors: environment limited resources, HPA scaling time, replicas creation time.

On the Prometheus server, the used metric is visible and pikes can be clearly distinguished (figure 22):



Figure 22: Prometheus metric evolution for final stress test

## 7.4 Bugs reported / fixed

The functional testing reveals 3 majors bugs, 2 on the Responder and one on the QAdapter:

1. Responder auto down scale closes pending requests
2. Responder stores too much events when there is multiple replicas
3. QAdapter auto down scale can cause loss of messages

Each of these three bugs are described in the next chapters.

### 7.4.1 Responder closes pending requests

During the stress test, a HPA has been setup for the Responder component. This HPA was also based (as the HPA for the service) on incoming messages to queue, but on the response queue (unlike the HPA for the service which is based on the request queue).

When the need of Responder component decrease, the HPA will scale down the number of replicas for the Responder. The issue arrives when there is connected request on a Responder and this is shut down. All the requests are stopped for the client.

No concrete solution has been found, but a little improvement can be set up by configuring the scale down behavior of the Responder HPA. This improvement consists of expanding the stabilization window to prevent the down scaling too close in time with a request pike. Also, limiting the number of replicas that must be scaled down at the same time and setting a bigger time window between two down scale improved the tests results.

To do what described just above, these lines were added to the Responder HPA:

```

1 behavior:
2   scaleDown:
3     stabilizationWindowSeconds: 600 # 10 minutes
4     policies:
5     - type: Pods
6       value: 1
7     periodSeconds: 600 # wait 10 minutes before scaling down another replica
8     selectPolicy: Min

```

#### 7.4.2 Responder stores too much events

As a reminder: The Responder component is based on a KubeMQ Events Store with Pub/Sub protocol on the response queue. The Pub/Sub protocol implies that the subscriber receives all the events published to the response queue. This means that if there are 3 subscribers to the response queue and an event "A" is published, each of the 3 subscribers will receive the event "A". This is not like a standard queue with publisher/consumer where only one of the consumers receive a message published in the queue.

The Responder must use the Pub/Sub protocol because when the Proxy component redirects to a Responder, it does not redirect to a Responder directly, but to its service/load balancer and it is impossible to predict to which Responder replica the load balancer will redirect the request. So, every Responder must receive all the events published in the response queue they are attached to.

As described in the chapter 6.4.2 (Responder implementation), the Responder implements a storage mechanism of events for all events not linked to an existing observer. The issue is that a responder does not receive all corresponding requests for all the events it receives because the requests are distributed across the Responder replicas but the events are sent to all Responders.

This implies that the store of the Responder makes 3 times the number of requests it receives (if considering an equitable distribution done by the load balancer) when there are 3 replicas (same logic for other number of replica). This behavior can cause performance loss and high consumption of disk resource.

The solution found is to setup a cron job that will clean all outdated instances of events. To know which is outdated, it has to save the time they were added. Are considered outdated, all events which were added more time ago than the server timeout. This means that if there is a connection which requests an event, this connection will be closed before the event will be cleared from the event store.

To implement this logic, two classes were added: ResponseStore and StoreCleanerCron.

The StoreCleanerCron class implements a unique method that is launched every 10 minutes in another thread than the main one. Below, the implementation of the single method:

```

1 @Scheduled(fixedRate = 1000*60*10)
2 public void cleanResponseStore() {
3     log.warn("Clearing ResponseStore - Cron");
4     ResponseStore.clearOldEvents(
5         ↪ environment.getProperty("spring.mvc.async.request-timeout", Duration.class));
6 }

```

The response store is a synchronized implementation of a response store. It must be synchronized because the StoreCleanerCron runs in another thread than the map which represents the store

is a shared resource.

On the figure 23 below, the class diagram of the implemented solution. 22):

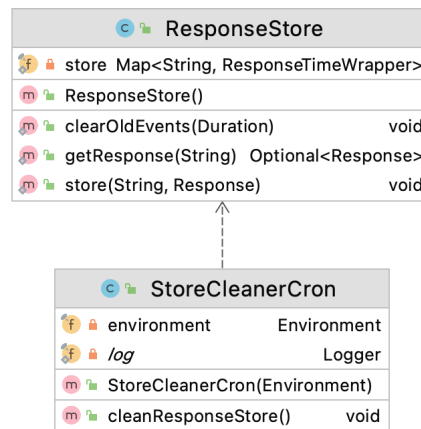


Figure 23: Class diagram of implemented solution for responder storing too much events

### 7.4.3 QAdapter auto-scale loss of messages

When the service HPA shut down replicas for down scaling, requests on which the terminating replicas were working on (at the time of the shut down) are lost.

The same issue appears when shutting down simple Rest service which has pending requests.

This issue is not compatible with the project constraints where all messages must be processed.

The solution found to prevent loss of messages when a replica of a service is stopped, is to implement the use of Message visibility feature of KubeMQ. This feature (described in the chapter 4.9.1) allows to set a visibility to a message in queue without removing it from the queue. Visibility can be "visible" or "not visible". If a message is visible, every consumer see it like a normal message, but if it is not visible, it acts like there is no message. The visibility has a timer and if there is no ack for the message within the allotted time, the message visibility changes to "visible" and be accessed by other consumers.

When retrieving a message in the queue and setting visibility to it, a initial visibility duration must be set. The consumer who got the message and set the visibility time has the possibility to expend this time if he needs more time.

The implementation of the visibility logic requires some changes in the QAdapter component. Instead of blocking on the response of the service the QAdapter is linked to (as described in chapter 6.3.2), it will send the request asynchronously and during the time between the request sending and the response receiving, it will check for the response (polling). If the response is not present and no issues appeared, it extends the visibility time to ask for more time. Time between checks (polling) is constant.

Sending visibility extension to KubeMQ can overload traffic the QAdapter must send the less visibility extension as possible.

So, a visibility extension is only sent when the current time + time between response check appears later than the end of the current visibility time. To illustrate this strategy, if we take the figure 24, if the current time is  $t_0 + 1$  and there is still no response from the service, we won't make a visibility extension because we know that at the next time the response will be checked

$(t_0 + 2)$ , the visibility still valid. On the other hand, if the current time is  $t_0 + 4$ , we now that at our next response check ( $t_0 + 5$ ) the visibility will not be valid anymore. This means that we need to extend the visibility before the next response check.

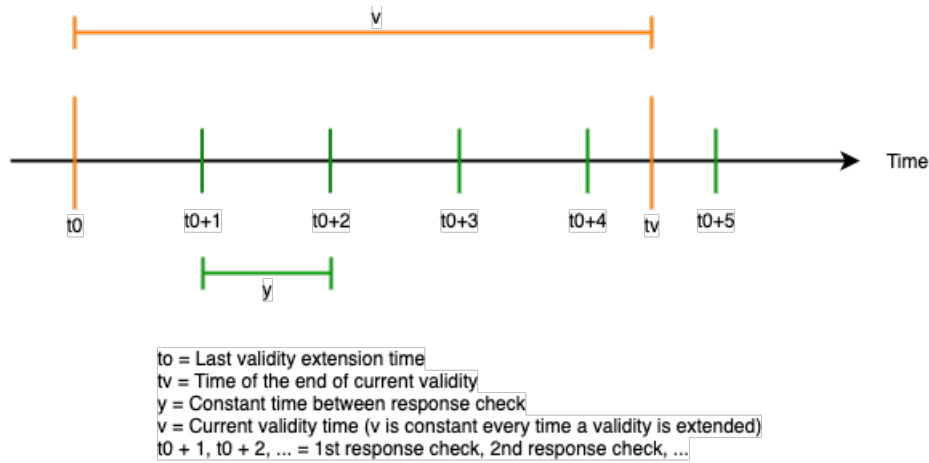


Figure 24: Graph representing validity check on QAdapter

Implementation of this strategy is done by 2 main changes: Adding a new class named "TransactionContextHolder" that is used as buffer class for response polling. When the response is received, it will push the response into an instance of this class and the polling loop will make the checks for response on the same instance. This class is synchronized to avoid conflict between the thread of the WebClient and the main train doing polling. This class also implement the possibility to save the last visibility extension time and is capable of return the duration between the last extension time and the current time. Adding a new method named "maintain-ConnectionUntilResponse" capable of polling and extension of the current visibility. This new method replaces the use of the method ".block()" on "Mono<?>" object used to send the request to the service.

## 7.5 Tests synthesis

The implemented solution works. It is not in production form but more like a proof of concept. There is still some little bugs related to KubeMQ Transaction that appears some times.

When comparing with a simple REST service without the Kueuer system we can find some advantages and some disadvantages of the Kueuer solution:

Disadvantages:

- Resource consumption. The need of deploying new objects in the Kubernetes environment, like the Proxy, the Responder, the Prometheus server and adapter and all the QAdapters makes the Kueuer solution not suitable for K8s clusters limited in resources.
- Scaling reactivity. With the QAdapter container in each service, it increases the time for Kubernetes to create a new replica of the service which slows down the scaling reactivity.
- Request duration. The request duration may be increased due to the complexity the Kueuer system adds.

Advantages:

- Distributed works. The Kueuer system works with queues and this offers the possibility to easily attach new worker on it.
- Less rejected requests. With the queue working as a buffer, the Kueuer system offers less rejections of requests, they will be processed and not returned to the client as rejected because of overloading.
- No message lost, less client re-call. When downscaling the amount of replicas of a service, no message will be lost, it will be automatically re-injected into the request queue.

## 8 Conclusion

This section discusses the conclusion of the project. The conclusion is divided in four different parts. The first part is the objectives of the project, then the second part, the state of the project: what was achieved. Next the third part is evolution of the project and the next steps and finally, the last part, a personal conclusion on the project by the writer (Romain Agostinelli).

### 8.1 Objectives and achievement

The goal of the project is to set a system capable of queuing and scaling (based on queue metrics) in a Kubernetes environment.

Multiple proposals have been made to solve this problem and one has been chosen. The name of the chosen solution is "Kueuer" (for Kubernetes and Queue). This solution is implemented and tested.

The main objective of the project has been achieved. It remains little bugs and improvements to do but a working solution can be easily deployed and tested.

Moreover, the developed system can be used in different environment and is easily configurable (do not depends on a specific Kubernetes environment).

### 8.2 State of the project

A system named "Kueuer" is developed, documented and tested. This system is actually based on KubeMQ and workers coded in Java.

The developed system is not a production system. It is more a "proof-of-concept" and must not be use as-is in a production environment. A complete deployment directory is available on the project's GitLab (see ch. 9).

### 8.3 Evolution of the project

There is two things to distinguish, the project and the developed solution.

For improvement for the developed solution, please refer to the chapter 6.7.

For the project, a good evolution would be to test other forms of solution, other system.

After testing the Kueuer system, the main problem detected is the consumption of resources. Maybe an alternative solution would be to decrease the number of new component in the architecture. Implementing a single proxy which does everything may be solution. It could improve reactivity by limiting the resources needed for launching new replica. I think it is a good way for next researches based on the same subject.



## 8.4 Personal conclusion

For me, it was the first time working with Kubernetes. It requires a bit of time to become familiar with Kubernetes environment but when I began to understand how everything works, how things communicate and what is possible, I immediately loved working with it. Today, I think I have a good comprehension of how Kubernetes works and it may help me developing project for Kubernetes in the future.

The subject is a real problem case and searching, experimenting and developing solutions was very exciting. During the whole project, I acquired a lot of knowledge, both in computer architecture and in development.

The only regret I have, is to not have more time to work on this project. I think I will continue to work on solutions on my side, testing things. It may be a good project for me to learn the Go language to move the created solution into a Go program.

I became really attached and interested in this topic, queuing and scaling. The project led me to something I never did: developing for IT infrastructure. I always developed solution for client requirement, not for architecture requirement and today, after this project, I think I found what I really want to do later: working on this kind of projects.

I made some mistakes during the project that I wanted to list here as a reminder for future projects:

- Plan more time for testing and begin testing at the beginning of the project
- Implement automated tests and deployment
- Plan development pause to check what already exists, to take step back
- Always make your estimated task time  $\times \pi$
- Keep it simple

## 9 Resources

All the developments, code and documentation is available on the project's GitLab at this address:

<https://gitlab.forge.hefr.ch/ps6-2021-kubemq>

Also, already compiled on images of the "Kueuer" system are available on the Docker Hub page of the writer:

<https://hub.docker.com/u/agost/> and begin with the prefix: "agost/io.kueuer.\*"

## 10 Déclaration d'honneur - Declaration of honour

Je, soussigné, Romain Agostinelli, déclare sur l'honneur que le travail rendu est le fruit d'un travail personnel. Je certifie ne pas avoir eu recours au plagiat ou à toutes autres formes de fraudes. Toutes les sources d'information utilisées et les citations d'auteur ont été clairement mentionnées.

Romain Agostinelli:

A handwritten signature in black ink, appearing to read 'R. Agostinelli', with a small dot above the final 'i'.

## Glossary

**K8s** Synonym for "Kubernetes". This mean "K" then 8 letters and finally "s" to form the word "Kubernetes". . 2, 3

**OS** Operating System. Most known are Windows (Microsoft), macOS (Apple) and Linux (multiple OS based on Unix kernel). 4

## References

- [1] Kubernetes. Kubernetes manage resources container.
- [2] Kubernetes. Kubernetes api server aggregation.
- [3] Kubernetes. Kubernetes controllers.
- [4] Kubernetes. Kubernetes components.
- [5] Kubernetes. Kubernetes manage resources container.
- [6] Kubernetes. Kubernetes pod autoscal walkthrough.
- [7] Kubernetes. Jobs.
- [8] Kubernetes. Cron job.
- [9] Kubernetes. Fine parallel processing work queue.
- [10] Bhargav Bachina. Kubernetes — learn sidecar container pattern.
- [11] Baeldung. Java fork join pool.
- [12] Docker. Docker command run.
- [13] Kubernetes. Define environment variable.
- [14] NGINX. Advanced configuration with annotations.
- [15] F5 Manual. Preventing tcp connection requests from being dropped.
- [16] Nick Ramirez. Protect servers with haproxy connection limits and queues.
- [17] KubeMQ. Kubemq java sdk.
- [18] Baeldung. Spring 5 webclient.
- [19] VERSPRITE. Spring 5 webclient.
- [20] Prometheus Adapter. Configuration walkthrough.
- [21] KubeMQ. Kubemq dashboard.

## 11 Appendixes

This section contains all the different appendixes of the project.

### 11.1 Simple App with HPA based on CPU consumption

```
1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: s1
5    labels:
6      app: s1
7  spec:
8    selector:
9      matchLabels:
10       app: s1
11    replicas: 1 # tells deployment to run 2 pods matching the template
12    template:
13      metadata:
14        labels:
15          app: s1
16      spec:
17        containers:
18          - name: s1
19            image: agost/io.kueuer.timesim:latest
20            imagePullPolicy: Always
21            ports:
22              - containerPort: 3000
23            resources:
24              requests:
25                cpu: "250m"
26  ---
27  apiVersion: v1
28  kind: Service
29  metadata:
30    name: s1-service
31  spec:
32    selector:
33      app: s1
34    type: LoadBalancer
35    ports:
36      - port: 8080
37        protocol: TCP
38        targetPort: 3000
39        nodePort: 30000
40  ---
41  apiVersion: autoscaling/v2beta2
42  kind: HorizontalPodAutoscaler
43  metadata:
44    name: s1-hpa
45  spec:
```

```

46  scaleTargetRef:
47    apiVersion: apps/v1
48    kind: Deployment
49    name: s1
50  minReplicas: 1
51  maxReplicas: 30
52  metrics:
53  - type: Resource
54    resource:
55      name: cpu
56      target:
57        type: Utilization
58        averageUtilization: 20
59  behavior:
60    scaleDown:
61      stabilizationWindowSeconds: 60 # 1 minutes
62      policies:
63      - type: Pods
64        value: 1
65      periodSeconds: 60 # wait 1 minutes before scaling down (1 by 1)
66      selectPolicy: Min

```

## 11.2 Environment deployment configuration files

### 11.2.1 App with QAdapter Deployment file

```

1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: s1
5    labels:
6      app: s1
7  spec:
8    selector:
9      matchLabels:
10     app: s1
11  replicas: 1 # tells deployment to run 2 pods matching the template
12  template:
13    metadata:
14      labels:
15        app: s1
16    spec:
17      containers:
18      - name: s1
19        image: agost/io.kueuer.timesim:latest
20        imagePullPolicy: Always
21      - name: qadapter
22        image: agost/io.kueuer.qadapter:latest
23        imagePullPolicy: Always

```

```

24     env:
25     - name: port
26       value: "3003"
27     - name: queue_addr
28       value: "kubemq-cluster-grpc.kubemq.svc.cluster.local:50000"
29     - name: request_queue
30       value: "s1-request"
31     - name: response_channel
32       value: "s1-response"
33     - name: service_port
34       value: "3000"

```

### 11.2.2 App - HPA Based on Queue metrics

```

1  apiVersion: autoscaling/v2beta2
2  kind: HorizontalPodAutoscaler
3  metadata:
4    name: s1-hpa
5  spec:
6    scaleTargetRef:
7      apiVersion: apps/v1
8      kind: Deployment
9      name: s1
10   minReplicas: 1
11   maxReplicas: 10
12   metrics:
13   - type: External
14     external:
15       metric:
16         name: queue-incoming-message-per-second
17         selector: {
18           matchLabels: {
19             channel: "s1-request",
20             type: "queues"
21           }
22         }
23     target:
24       type: Value
25       value: 0.05
26   behavior:
27     scaleDown:
28       stabilizationWindowSeconds: 30 # 30 secs
29     policies:
30     - type: Pods
31       value: 2
32       periodSeconds: 30 # wait 30 seconds before scaling down (2 by 2)
33     selectPolicy: Min

```

### 11.2.3 Proxy Deployment file

```

1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: s1-proxy
5    labels:
6      app: s1-proxy
7  spec:
8    selector:
9      matchLabels:
10       app: s1-proxy
11   replicas: 1 # set the number of replicas
12   template:
13     metadata:
14       labels:
15         app: s1-proxy
16     spec:
17       containers:
18         - name: s1-proxy
19           image: agost/io.kueuer.proxy:latest
20           imagePullPolicy: Always
21           ports:
22             - containerPort: 3001
23           env:
24             - name: port
25               value: "3001"
26             - name: queue_addr
27               value: "kubemq-cluster-grpc.kubemq.svc.cluster.local:50000"
28             - name: queue
29               value: "s1-request"
30             - name: responder_addr
31               value: "192.168.99.111:30001"
32               #value: "responder-service.default.svc.cluster.local:8080"
33   ---
34   apiVersion: v1
35   kind: Service
36   metadata:
37     name: s1-service
38   spec:
39     selector:
40       app: s1-proxy
41     type: NodePort
42     ports:
43       - port: 8080
44         protocol: TCP
45         targetPort: 3001
46         nodePort: 30002

```

## 11.2.4 Reponder deployment configuration

```

1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: responder
5    labels:
6      app: responder
7  spec:
8    selector:
9      matchLabels:
10       app: responder
11    replicas: 1 # set number of replicas here
12    template:
13      metadata:
14        labels:
15          app: responder
16      spec:
17        containers:
18          - name: responder
19            image: agost/io.kueuer.responder:latest
20            imagePullPolicy: Always
21            ports:
22              - containerPort: 8080 # Default app is 8080
23            env:
24              - name: queue_addr
25                value: "kubemq-cluster-grpc.kubemq.svc.cluster.local:50000"
26              - name: response_channels
27                value: "s1-response"
28    ---
29    apiVersion: v1
30    kind: Service
31    metadata:
32      name: responder-service
33    spec:
34      selector:
35        app: responder
36      type: NodePort
37      ports:
38        - port: 8080
39          protocol: TCP
40          targetPort: 8080
41          nodePort: 30001

```

## 11.3 Metrics configurations files

### 11.3.1 Prometheus Configuration



```

1  apiVersion: apps/v1
2  kind: Deployment
3  metadata:
4    name: prometheus-deployment
5    namespace: kubemq
6    labels:
7      app: prometheus-server
8  spec:
9    replicas: 1
10   selector:
11     matchLabels:
12       app: prometheus-server
13   template:
14     metadata:
15       labels:
16         app: prometheus-server
17     spec:
18       containers:
19         - name: prometheus
20           image: prom/prometheus
21           args:
22             - "--config.file=/etc/prometheus/prometheus.yml"
23             - "--storage.tsdb.path=/prometheus/"
24           ports:
25             - containerPort: 9090
26           volumeMounts:
27             - name: prometheus-config-volume
28               mountPath: /etc/prometheus/
29             - name: prometheus-storage-volume
30               mountPath: /prometheus/
31       volumes:
32         - name: prometheus-config-volume
33           configMap:
34             defaultMode: 420
35             name: prometheus-server-conf
36         - name: prometheus-storage-volume
37           emptyDir: {}
38 ---
39 apiVersion: rbac.authorization.k8s.io/v1beta1
40 kind: ClusterRole
41 metadata:
42   name: prometheus
43 rules:
44   - apiGroups: [""]
45     resources:
46       - nodes
47       - nodes/proxy
48       - services
49       - endpoints
50       - pods
51   verbs: ["get", "list", "watch"]

```

```

52   - apiGroups:
53     - extensions
54   resources:
55     - ingresses
56   verbs: ["get", "list", "watch"]
57   - nonResourceURLs: ["/metrics"]
58     verbs: ["get"]
59 ---
60 apiVersion: rbac.authorization.k8s.io/v1beta1
61 kind: ClusterRoleBinding
62 metadata:
63   name: prometheus
64 roleRef:
65   apiGroup: rbac.authorization.k8s.io
66   kind: ClusterRole
67   name: prometheus
68 subjects:
69   - kind: ServiceAccount
70     name: default
71     namespace: kubemq
72 ---
73 apiVersion: v1
74 kind: ConfigMap
75 metadata:
76   name: prometheus-server-conf
77   labels:
78     name: prometheus-server-conf
79     namespace: kubemq
80 data:
81   prometheus.yml: |-
82     global:
83       scrape_interval: 5s
84       evaluation_interval: 5s
85     rule_files:
86     - ./kubemq.rule.yaml
87     scrape_configs:
88     - job_name: scrape_kubemq
89       kubernetes_sd_configs:
90       - role: pod
91         namespaces:
92           names:
93             - "kubemq"
94       relabel_configs:
95       - source_labels: [__meta_kubernetes_pod_container_port_number]
96         action: keep
97         regex: 8\d{3}
98   kubemq.rule.yaml: |-
99     groups:
100     - name: events
101       rules:
102       - record: kubemq:messages:count:events:send:sum
103         expr: sum(kubemq_messages_count{type="events",side="send"}) OR on()
           ↪ vector(0)

```

```

104 - record: kubemq:messages:count:events:receive:sum
105   expr: sum(kubemq_messages_count{type="events",side="receive"}) OR
      ↳ on() vector(0)
106 - record: kubemq:messages:count:events:sum
107   expr:
      ↳ kubemq:messages:count:events:send:sum+kubemq:messages:count:events:receive:sum
108 - record: kubemq:messages:volume:events:send:sum
109   expr: sum(kubemq_messages_volume{type="events",side="send"}) OR
      ↳ on() vector(0)
110 - record: kubemq:messages:volume:events:receive:sum
111   expr: sum(kubemq_messages_volume{type="events",side="receive"}) OR
      ↳ on() vector(0)
112 - record: kubemq:messages:volume:events:sum
113   expr:
      ↳ kubemq:messages:volume:events:send:sum+kubemq:messages:volume:events:receive:sum
114 - record: kubemq:messages:errors:events:sum
115   expr: sum(kubemq_errors_count{type="events"}) OR on() vector(0)
116 - name: events_store
117   rules:
118     - record: kubemq:messages:count:events_store:send:sum
119       expr: sum(kubemq_messages_count{type="events_store",side="send"})
          ↳ OR on() vector(0)
120     - record: kubemq:messages:count:events_store:receive:sum
121       expr:
          ↳ sum(kubemq_messages_count{type="events_store",side="receive"})
          ↳ OR on() vector(0)
122     - record: kubemq:messages:count:events_store:sum
123       expr:
          ↳ kubemq:messages:count:events_store:send:sum+kubemq:messages:count:events_store:receive:sum
124     - record: kubemq:messages:volume:events_store:send:sum
125       expr: sum(kubemq_messages_volume{type="events_store",side="send"})
          ↳ OR on() vector(0)
126     - record: kubemq:messages:volume:events_store:receive:sum
127       expr:
          ↳ sum(kubemq_messages_volume{type="events_store",side="receive"})
          ↳ OR on() vector(0)
128     - record: kubemq:messages:volume:events_store:sum
129       expr:
          ↳ kubemq:messages:volume:events_store:send:sum+kubemq:messages:volume:events_store:receive:sum
130     - record: kubemq:messages:errors:events_store:sum
131       expr: sum(kubemq_errors_count{type="events_store"}) OR on()
          ↳ vector(0)
132 - name: queues
133   rules:
134     - record: kubemq:messages:count:queues:send:sum
135       expr: sum(kubemq_messages_count{type="queues",side="send"}) OR on()
          ↳ vector(0)
136     - record: kubemq:messages:count:queues:receive:sum
137       expr: sum(kubemq_messages_count{type="queues",side="receive"}) OR
          ↳ on() vector(0)
138     - record: kubemq:messages:count:queues:sum
139       expr:
          ↳ kubemq:messages:count:queues:send:sum+kubemq:messages:count:queues:receive:sum

```

```

140 - record: kubemq:messages:volume:queues:send:sum
141   expr: sum(kubemq_messages_volume{type="queues",side="send"}) OR
      ↳ on() vector(0)
142 - record: kubemq:messages:volume:queues:receive:sum
143   expr: sum(kubemq_messages_volume{type="queues",side="receive"}) OR
      ↳ on() vector(0)
144 - record: kubemq:messages:volume:queues:sum
145   expr:
      ↳ kubemq:messages:volume:queues:send:sum+kubemq:messages:volume:queues:receive:sum
146 - record: kubemq:messages:errors:queues:sum
147   expr: sum(kubemq_errors_count{type="queues"}) OR on() vector(0)
148 - name: queries
149   rules:
150     - record: kubemq:messages:count:queries:send:sum
151       expr: sum(kubemq_messages_count{type="queries",side="send"}) OR
          ↳ on() vector(0)
152     - record: kubemq:messages:count:queries:receive:sum
153       expr: sum(kubemq_messages_count{type="queries",side="receive"}) OR
          ↳ on() vector(0)
154     - record: kubemq:messages:count:queries:sum
155       expr:
          ↳ kubemq:messages:count:queries:send:sum+kubemq:messages:count:queries:receive:sum
156     - record: kubemq:messages:volume:queries:send:sum
157       expr: sum(kubemq_messages_volume{type="queries",side="send"}) OR
          ↳ on() vector(0)
158     - record: kubemq:messages:volume:queries:receive:sum
159       expr: sum(kubemq_messages_volume{type="queries",side="receive"}) OR
          ↳ on() vector(0)
160     - record: kubemq:messages:volume:queries:sum
161       expr:
          ↳ kubemq:messages:volume:queries:send:sum+kubemq:messages:volume:queries:receive:sum
162     - record: kubemq:messages:errors:queries:sum
163       expr: sum(kubemq_errors_count{type="queries"}) OR on() vector(0)
164 - name: commands
165   rules:
166     - record: kubemq:messages:count:commands:send:sum
167       expr: sum(kubemq_messages_count{type="commands",side="send"}) OR
          ↳ on() vector(0)
168     - record: kubemq:messages:count:commands:receive:sum
169       expr: sum(kubemq_messages_count{type="commands",side="receive"}) OR
          ↳ on() vector(0)
170     - record: kubemq:messages:count:commands:sum
171       expr:
          ↳ kubemq:messages:count:commands:send:sum+kubemq:messages:count:commands:receive:sum
172     - record: kubemq:messages:volume:commands:send:sum
173       expr: sum(kubemq_messages_volume{type="commands",side="send"}) OR
          ↳ on() vector(0)
174     - record: kubemq:messages:volume:commands:receive:sum
175       expr: sum(kubemq_messages_volume{type="commands",side="receive"})
          ↳ OR on() vector(0)
176     - record: kubemq:messages:volume:commands:sum
177       expr:
          ↳ kubemq:messages:volume:commands:send:sum+kubemq:messages:volume:commands:receive:sum

```

```

178         - record: kubemq:messages:errors:commands:sum
179           expr: sum(kubemq_errors_count{type="commands"}) OR on() vector(0)
180     - name: responses
181       rules:
182         - record: kubemq:messages:count:responses:sum
183           expr: sum(kubemq_messages_count{type="responses"}) OR on()
184             ↪ vector(0)
185         - record: kubemq:messages:volume:responses:sum
186           expr: sum(kubemq_messages_volume{type="responses"}) OR on()
187             ↪ vector(0)
188         - record: kubemq:messages:errors:responses:sum
189           expr: sum(kubemq_errors_count{type="responses"}) OR on() vector(0)
190
191 ---
192 apiVersion: v1
193 kind: Service
194 metadata:
195   name: prometheus-service
196   namespace: kubemq
197   annotations:
198     prometheus.io/scrape: 'true'
199     prometheus.io/port: '9090'
200 spec:
201   selector:
202     app: prometheus-server
203   type: NodePort
204   ports:
205     - port: 8080
206       targetPort: 9090
207       nodePort: 30000

```

### 11.3.2 Prometheus adapter configuration

```

1 prometheus:
2   url: http://prometheus-service.kubemq.svc
3   port: 8080
4 rules:
5   external:
6     - seriesQuery: '{__name__=~"kubemq_messages_count.*"}'
7     resources:
8       overrides:
9         namespace:
10           resource: namespace
11         pod:
12           resource: pod
13   name:
14     matches: ^(.*)
15     as: "queue-incoming-message-per-second"
16 metricsQuery: sum(irate(kubemq_messages_count{<< range $key, $value :=
17   ↪ .LabelValuesByName >><< if ne $key "namespace" >><< $key >>="<< $value
18   ↪ >>",<< end >><< end >>side="send"}[5m])) by (channel, type)

```

```

17 ##### WARNING: queue-length is not guarantee
18 - seriesQuery: '{__name__=~"kubemq_messages_count.*"}'
19 resources:
20   overrides:
21     namespace:
22       resource: namespace
23     pod:
24       resource: pod
25   name:
26     matches: ^(.*)
27     as: "queue-length"
28   metricsQuery: sum(kubemq_messages_count{<< range $key, $value :=
    → .LabelValuesByName >><< if ne $key "namespace" >><< $key >>=<< $value
    → >>,<< end >><< end >>side="send"})-sum(kubemq_messages_count{<< range
    → $key, $value := .LabelValuesByName >><< if ne $key "namespace" >><<
    → $key >>=<< $value >>,<< end >><< end >>side="receive"})
29
30 # For the metric queries function, this is a workaround to an issue on external
    → metrics on prometheus-adaper:
31 # workaround proposition:
    → https://github.com/kubernetes-sigs/prometheus-adapter/issues/255#issuecomment-601955919
32 # Issue: https://github.com/kubernetes-sigs/prometheus-adapter/issues/282

```

## 11.4 K6 Tests configurations

### 11.4.1 Functional testing K6 configuration

```

1 import http from 'k6/http';
2 import { sleep } from 'k6';
3
4 export let options = {
5   stages: [
6     { duration: '1m', target: 5 }, // normal load
7     { duration: '40s', target: 20 }, // autoscaler increase
8     { duration: '1m', target: 5 }, // autoscaler down
9   ],
10  thresholds: {
11    http_req_failed: [{ threshold: 'rate<0.01', abortOnFail: false }],
12  },
13 }
14
15 export default function () {
16   let params = {
17     timeout: '300s', // chrome default timeout
18     tags: { k6test: 'yes' },
19   };
20   http.get('http://192.168.99.111:30002', params);
21   sleep(1);
22 }

```

### 11.4.2 K6 Stress test configuration

```

1  import http from 'k6/http';
2  import { sleep } from 'k6';
3
4  export let options = {
5      stages: [
6          { duration: '10s', target: 10 }, // below normal load
7          { duration: '1m', target: 10 },
8          { duration: '10s', target: 100 }, // spike to 100 users
9          { duration: '3m', target: 100 }, // stay at 100 for 3 minutes
10         { duration: '10s', target: 20 }, // scale down. Recovery stage.
11         { duration: '3m', target: 20 },
12         { duration: '7m', target: 10 }, // down scaling time
13         { duration: '40s', target: 50 }, // stress
14     ],
15     thresholds: {
16         http_req_failed: [{ threshold: 'rate<0.01', abortOnFail: false }],
17     },
18     discardResponseBodies: true,
19 }
20
21 export default function () {
22     let params = {
23         timeout: '300s', // chrome default timeout
24         tags: { k6test: 'yes' },
25     };
26     http.get('http://192.168.99.110:30002', params);
27     sleep(1);
28 }

```

## 12 Appendix - Cahier des charges - Specification

### 12.1 Introduction

These days, cloud services are receiving thousands of requests per second. During a normal development, it is often considered to use request queuing systems (RabbitMQ, Apache Kafka) to process requests in their order of arrival. These systems can create long waits if the number of requests per second is very high or if the service in question takes a long time to respond.

Another way to manage this request flow is to duplicate services so that several requests can be handled in parallel and thus reduce the wait.

Kubernetes is a popular state-of-the-art open-source computing workload orchestration system. It provides means to automate application deployments and scaling capabilities depending on the load of service requests.

The aim of the project is to combine hyperscaling capabilities of Kubernetes and queuing mechanisms to create an infrastructure capable of handling thousands of requests per second.

## 12.2 Context

The company Swisscom SA intends to migrate some of its services in a K8s environment and thus have a large Kubernetes micro-services infrastructure that exposes these different services.

Some micro-services take longer to respond than others, sometimes due to a service in general or sub-services that individually take a lot of time and slow down the processing of requests.

The goal is to be able to set up a system to manage requests and automatically increase the number of replicas of a service, whether it is a service in its entirety or only sub-services.

## 12.3 Objective

Implementation of a solution based on different techniques (such as hyperscaling and queuing) to handle numerous queries in order to improve request processing performance in a test environment simulating a large-scale Kubernetes infrastructure. The main constraint is the implementation of the *exactly-one semantics* concept in the solution.

## 12.4 Infrastructure topologies

In such infrastructure, multiple types of micro-services topologies can be imagined. This project considers two types: a first one, very simple, that will be used and setup first and a second one, more complex on which it will be tried to extend the solution found for the first one to make it works on a more complex infrastructure.

The project can then be divided in two parts, each one based on different micro-service topologies.

The both presented infrastructure topologies in the chapters below (figures 25 & 26) represent base infrastructure topologies. They are not optimized and the project aim is to improve these topologies regarding different constraints.

### 12.4.1 First part

For the first part of the project, a simple topology (figure 25) will be considered. The project will be focused on improving this topology and implementing a solution regarding the objectives and constraints.



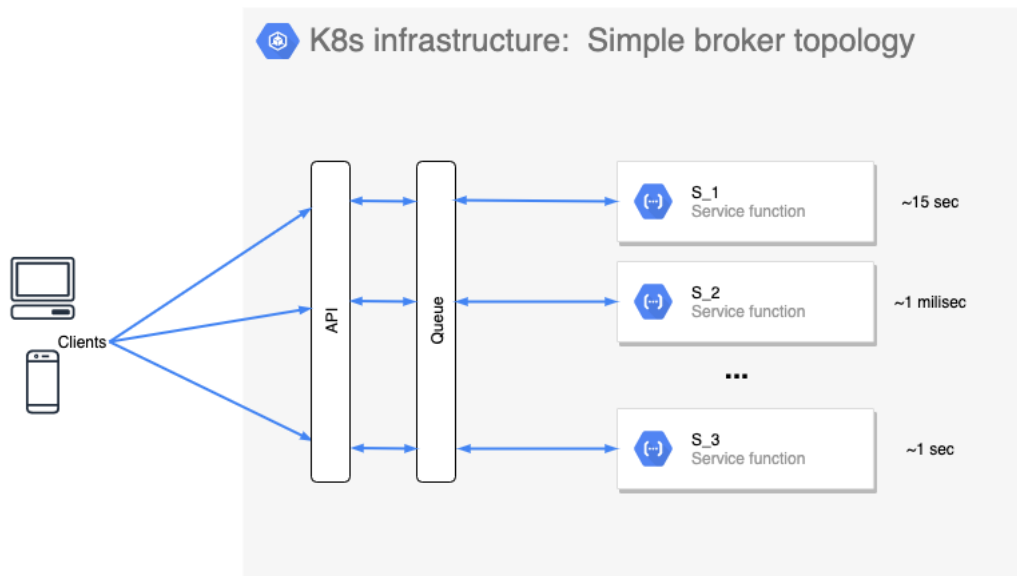


Figure 25: First simple topology

### 12.4.2 Second part

For the second part of the project, a more complex topology (figure 26) will be used and the solution found for the first part will be extended to fit this second topology and improve this one too.

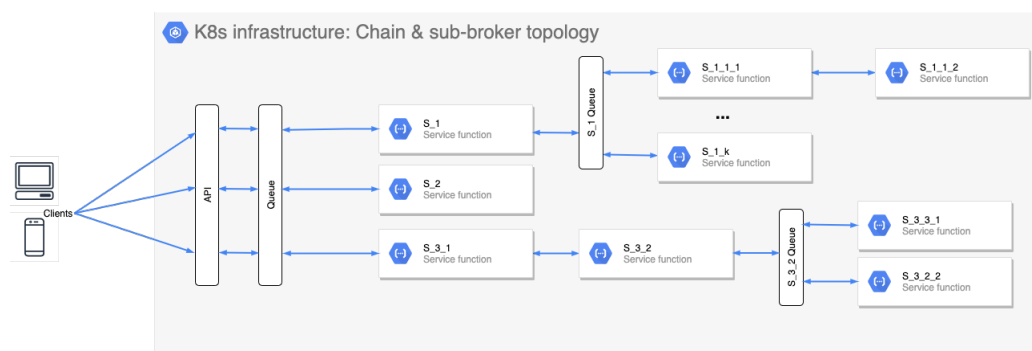


Figure 26: Second complex topology

## 12.5 Activities

In this section you will find the different activities present in the planning. Each sub-section represents a particular activity (step) of the project with a more precise definition of the deliverable as well as the tasks to be performed for the activity.

### 12.5.1 Kickoff and specs

**Deliverable :** Specification + planning

**Tasks :**

- Understand and getting familiar with K8s.
- Creation of the projects specifications and validate it.
- Creation of the planning using Agile method.
- Preliminary search:
  - Kubernetes
  - Overscaling
  - Queuing technique:
    - \* Apache Kafka
    - \* KubeMQ
    - \* RabbitMQ
  - Constraint:
    - \* Exactly-once semantics

### 12.5.2 Setting up a development environment

This second step of the project consists in setting up a development environment with Kubernetes using K3s or K8s solution.

**Deliverable :** Working development environment (K3|8s running).

**Tasks :**

- Compare and choose between K8s and K3s
- Setup K3|8 environment
- Tests on infrastructure, learning technology
- Research about technology update

### 12.5.3 First micro-service

The third step of the project consists in the creation and integration in the K3|8s development environment of multiple micro-services which simulate short and long processing request time. This is a simple micro-service environment, a more complex will come later.

It also contains the first research and experimentation of solutions allowing to solve the objective.

**Deliverable :** Micro-services which simulate request processing with various amount of time.

**Tasks :**

- Setup first micro-service topology infrastructure
- Research about exactly-once semantics
- Research of existing solution
- Comparison of solutions
- Proposal of one solution

#### 12.5.4 Solution

This fourth step consists of the implementation of the chosen solution in the previous step in the test environment. This contains multiple parts as analysis and conception of the chosen solution and then implementation.

**Deliverable :** A working implementation of the solution.

**Tasks :**

- Analysis and conception of the solution
- Implementation of the solution
- Testing of the solution

#### 12.5.5 Complex infrastructure

This last step consists of extending the solution implemented in the previous step for a more complex micro-service topology.

**Deliverable :** Extension of the solution for complex topology.

**Tasks :**

- Extension of the topology infrastructure
- Extension of the solution for the new infrastructure
- Testing of the solution

### 12.6 Planning

This project follows an Agile planning in order to be able to adapt as well as possible to the progress of research. This planning (figure 27 p.72) will allow us to focus on a single problem per sprint.

The decision to use *GitLab* and its *issues* system for planning was made. The flexibility of the platform better corresponds to the needs of Agile planning and will allow teachers and externals to follow the progress of the project live.



Figure 27: First Agile planning