

Master of Science in Informatics at Grenoble  
Master Informatique  
Specialization MoSIG

---

# Simulation of a Kubernetes Cluster with Validation in Real Conditions

---

**LARUE Théo**

Defense Date, 2020

Research project performed at Laboratoire d'Informatique de Grenoble

Under the supervision of:

Michael Mercier

Defended before a jury composed of:

Head of the jury

Jury member 1

Jury member 2



## Abstract

TODO : remove the focus on HPC

The rise of containerized applications has provided web platforms with much more control over their resources than they had before with their physical servers. Soon enough, developers realized they could go even further by automating container management operations to allow for even more scalability. The Cloud Native Computing Foundation was founded in this context, and developed Kubernetes which is a piece of software capable of container orchestration, or in other words, container management. Now, as we observe a convergence between HPC (High Performance Computing) and the Big Data field where Kubernetes is already the standard for some applications such as Machine Learning, discussions about leveraging containers for HPC applications rose and interest in Kubernetes has grown in the HPC community. One of the many challenges the HPC world has to face is scheduling, which is the act of allocating tasks submitted by users on available resources. In order to properly evaluate and develop schedulers researchers have used simulators for decades to avoid running experiments in real conditions, which is costly both in time and resources. However, such simulators do not exist for Kubernetes or are not open to the public. While the default scheduler works great for most of the Cloud Native infrastructures Kubernetes was designed for, some teams of researchers would rather be able to experiment with different batch processing policies on Kubernetes as they do with traditional HPC. Our goal in this master thesis is to describe how we developed Batkube, which it is an interface between Kubernetes schedulers and Batsim, a general purpose infrastructure simulator based on the Simgrid framework and developed at the LIG.

## Acknowledgement

I would like to express my sincere gratitude to .. for his invaluable assistance and comments in reviewing this report... Good luck :)

## Résumé

Abstract mais en français

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgement</b>	<b>i</b>
<b>Résumé</b>	<b>i</b>
<b>1 Introduction</b>	<b>1</b>

---

<b>2</b>	<b>State of the art</b>	<b>3</b>
2.1	Infrastructure simulators . . . . .	3
2.2	Kubernetes schedulers . . . . .	3
<b>3</b>	<b>Implementation</b>	<b>5</b>
3.1	The scheduling problem . . . . .	5
3.2	Batsim . . . . .	6
3.3	Kubernetes . . . . .	6
3.4	Objectives . . . . .	8
3.5	Technical challenges . . . . .	9
3.5.1	Translation . . . . .	9
3.5.2	Time hijack . . . . .	9
3.5.3	Re-building the API . . . . .	9
<b>4</b>	<b>Evaluation</b>	<b>11</b>
<b>5</b>	<b>Conclusion</b>	<b>13</b>
	<b>Bibliography</b>	<b>15</b>

# Introduction

The need for scalable computing infrastructure has increased tremendously in the last decades. Nearly every field of computer science, from research to the service industry, now needs a proper infrastructure and by 2025, computation technology could reach a fourth of the global electricity spending[1]. Even the public sector is now in need for efficient distributed infrastructure as the concept of smart cities is developing.

Organizations generally know what type of infrastructure will meet their needs. It can take the form of Big Data centers to store and analyze data, High-Performance Computers for computing intensive tasks or GPU banks for machine learning or crypto-currency mining. However, studying those infrastructures extensively is much more challenging. As these computers reach scales in the order of warehouses[2], quantifying a system's performance under varying loads, applications, scheduling policies and system size quickly becomes undoable without expensive real world experiments. In fact, the nature of scheduling problems[7] alone make theoretical studies hard. This is an issue for organizations as they rely on those studies to determine the size of the required system or choose optimal scheduling policies.

Simulation allows to tackle these issues by enabling users to draw conclusions empirically without the need to fire up real workloads. Indeed, running an entire experimental campaign on a real system is impossible for most due to large overhead and consequent costs both in time and money. With simulation, The gain in time and spent energy can be extreme : a HPC job spanning months on a real system can be resolved in a matter of minutes on any domestic computer. Another major point is that it also brings reproducibility to these experiments, that otherwise would have to be run on the exact same systems as their first iteration. With simulation, one can recreate the same conditions for any experiment anywhere they want, and expect the same results.

However, simulations need to be run with sound models for the results to be exploitable and in that regard, simulators usually fall under several pitfalls[8]. Very often simulators are implemented at the same time as new schedulers or Resource and Jobs Management Systems<sup>1</sup> in order to validate their algorithms. Thus, they are strongly coupled together and are not usable with any other software. They are either shipped with the software itself or worst, they are never released and discarded at the end of the development process. Moreover, still according to [8], strong coupling may lead to unrealistic models. In that case cluster resources can be accessed with ease by the scheduler, resulting in it having very precise information about the system state to take its decisions. This conflicts with the real world as a scheduler may not have access to all the information it wants, or may suffer from latency when getting it from the system.

To try and assess these issues a team of researchers at the LIG developed Batsim[5] which is a general purpose infrastructure simulator with modularity and separation of concerns in mind. Batsim is based on SimGrid[4] which is a framework for developing simulators for distributed computer systems. Simgrid

---

<sup>1</sup>The RJMS is the software at the core of the cluster. It is a synonym for a scheduler and manages resources, energy consumption, users' jobs life-cycle and implements scheduling policies.

is now a 20 years old framework that has been used in many projects<sup>2</sup>, making it a sound choice to run scalable and accurate models of the reality.

Batsim was designed to support algorithms written in any languages, as long as they support its communication protocol. It means that, while any scheduler found in the wild can potentially be run on a Batsim simulation, they still have to be adapted to make them compatible. This master's project is dedicated on developing an interface between Batsim and Kubernetes<sup>3</sup> schedulers in order to run Kubernetes clusters simulations. Kube<sup>4</sup> is an open source container management software originally developed by Google and then taken over by the Cloud Native Computing Foundation<sup>5</sup> (CNCF), widely exploited in the industry for its ease of use and wide range of capabilities. It has freed developers from the cumbersome task of setting up low level software infrastructure on their servers and automates maintenance, scaling and administration of their applications. For all these reasons it has become a de-facto solution for any organization that wishes to build new internet platforms from the ground up.

TODO : what we where able to do (summary of the simulator capabilities, experimentations, results)

---

<sup>2</sup><https://simgrid.org/usages.html>

<sup>3</sup><https://github.com/kubernetes/kubernetes/>

<sup>4</sup>Another term to designate Kubernetes. It is also sometimes called k8s.

<sup>5</sup><https://www.cncf.io/>

## State of the art

### 2.1 Infrastructure simulators

### 2.2 Kubernetes schedulers

#### **kube-scheduler**

Not exactly a batch scheduler, it is the default scheduler for Kubernetes made by the CNCF"" -> "hop là"

#### **kube-batch**

#### **Poseidon**

#### **bashScheduler by rothgar**

#### **random-scheduler by Banzaicloud**

#### **k8s-custom-scheduler by IBM**

#### **scheduler by kelseyhightower**





## Implementation

### 3.1 The scheduling problem

**schedule  $n$ .** : A plan for performing work or achieving an objective, specifying the order and allotted time for each part.

In a general way, scheduling is the concept of allocating available resources to a set of tasks, organizing them in time and space (the resource space). The resources can be of any nature, and the tasks independent from each others or linked together.

In computing the definition remains the same, but with automation in mind. Schedulers are algorithms that take as an input either a pre-defined workload, which is a set of jobs to be executed - the tasks are called jobs in this context -, or simple jobs submitted over time by users in an unpredictable manner. In the latter case, the jobs are added to a queue managed by the scheduler. Scheduling is also called batch scheduling or batch processing, as schedulers allocate batches of jobs at a time. Jobs are allocated on machines, virtual or physical, with the intent of minimizing the total execution time, equally distributing resources, minimizing wait time for the user or reducing energy costs. As these objectives often contradict themselves, schedulers have to implement compromises or focus on what the user really needs or requires.

The scheduler has many factors to keep in mind while trying to be as efficient as possible, such as :

- Resource availability and jobs resource requirements
- Link between jobs (some are executed in parallel and need synchronization, some are independent)
- Latency between compute resources
- Compute resources failures
- Jobs priority
- Machine shutdowns and restarts
- Data locality

All these elements make scheduling a very intricate problem that is at best polynomial in complexity, and often NP-hard[7]. Moreover, with the growing complexity of modern RJMS and the wide variety in infrastructures, scheduler performances are rather unpredictable.

## 3.2 Batsim

## 3.3 Kubernetes

### Kubernetes overview

In the early stages of application development, organizations used to run their services on physical servers. With this direct approach came many challenges that needed to be coped with manually like resources allocation, maintainability or scalability. In an attempt to automate this process developers started using virtual machines which enabled them to run their services regardless of physical infrastructure while having a better control over resources allocation. This led to the concept of containers which takes the idea of encapsulated applications further.

Containers can be thought of as lightweight virtual machines. Unlike the latter, containers share the same kernel with the host machine but still allow for a very controlled environment to run applications. There are many benefits to this : separating the development from deployment, portability, easy resource allocation, breaking large services into smaller micro-services or support of continuous integration tools (containers greatly facilitate integration tests).

The CNCF<sup>1</sup> (Cloud Native Computing Foundation) was founded in the intent of leveraging the container technology for an overall better web. In a general way, we now speak of these containerized and modular applications as cloud native computing :

*“Cloud native technologies empower organizations to build and run scalable applications in modern, dynamic environments such as public, private, and hybrid clouds. Containers, service meshes, microservices, immutable infrastructure, and declarative APIs exemplify this approach.*

*These techniques enable loosely coupled systems that are resilient, manageable, and observable. Combined with robust automation, they allow engineers to make high-impact changes frequently and predictably with minimal toil.”<sup>2</sup>*

Kubernetes<sup>3</sup> is the implementation of this general idea and was announced at the same time as the CNCF. It aims at automating of the process of deploying, maintaining and scaling containerized applications. It is industry grade and is now the de-facto solution for container orchestration.

The basic processing unit of Kubernetes is called a **pod** which is composed of one or several containers and volumes<sup>4</sup>. In the cloud native context a pod most often hosts a service or micro-service.

Pods are bundled together in **nodes** (figure 3.1) which are either physical or virtual machines. They represent another barrier to pass through to access the outside world which can be useful to add layers of security or facilitate communication between pods. Nodes take the idea of containerisation further by encapsulating the already encapsulated services. Each node runs at least one pod and also one **kubelet** which is a process responsible for communicating with the rest of Kubernetes (or more precisely, with the master node which in turns communicates with the api server). A set of nodes is called a **cluster**. Each Kubernetes instance is responsible for running a cluster.

Kubernetes revolves its API server which is its central component (figure 3.2). The majority of operations between components go through this REST API like user interactions through kubectl or scheduling operations.

<sup>1</sup><https://www.cncf.io/>

<sup>2</sup><https://github.com/cncf/toc/blob/master/DEFINITION.md>

<sup>3</sup><https://kubernetes.io/>

<sup>4</sup>A volume is some storage space on the host machine that can be linked to containers, so they can read persistent information or store data in the long term

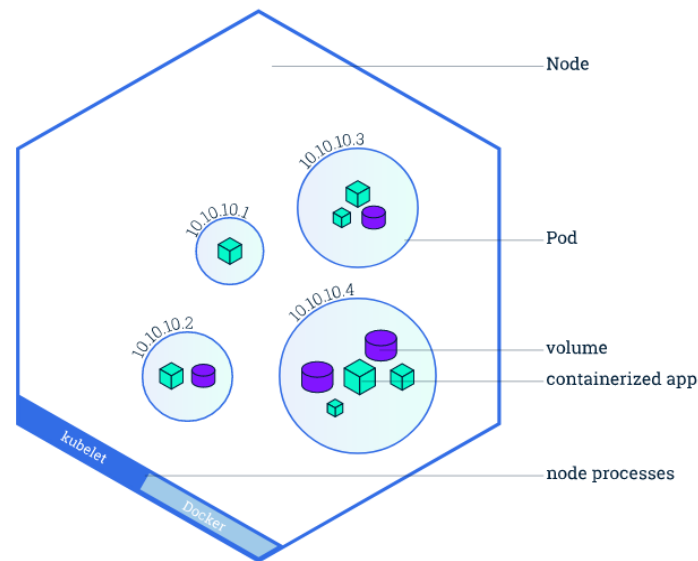


Figure 3.1: Node overview

**Source:** <https://kubernetes.io/docs/tutorials/kubernetes-basics/explore/explore-intro/>

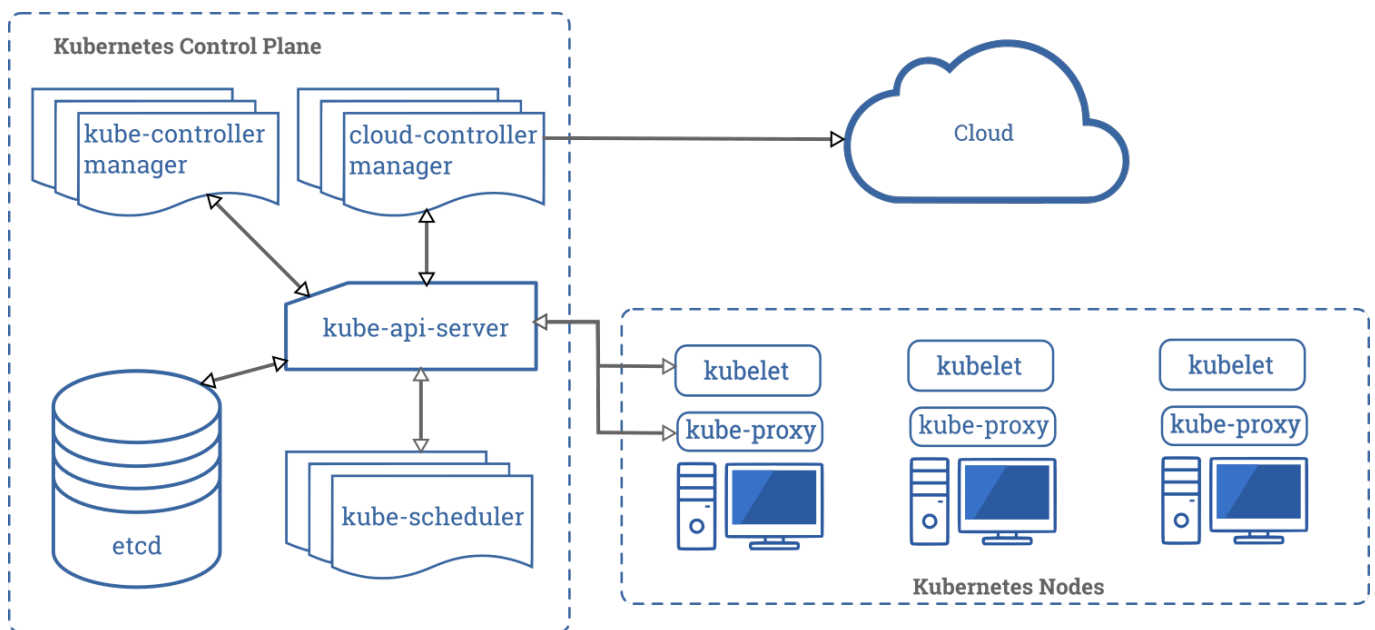


Figure 3.2: Components of Kubernetes

**Source:** <https://kubernetes.io/docs/concepts/overview/components/>

## HPC and Kubernetes

The difference between HPC and Cloud Native computing lies in the workloads they are intended to tackle. Kubernetes was designed for Cloud Native applications. Services or micro services are run in containers and are expected to be available at all times : they are replicated as many times as the user desires and restarted whenever a failure occurs. High availability is at the core of Kubernetes container management. On the other hand, depending on scheduling policies, HPC is focused on user wait time, maximizing resource usage, optimizing energy costs... For instance, in case of failure, it is sometimes not sufficient to restart the single job that failed : the entire submission must be re-run if it is part of several jobs computed in parallel.

Kubernetes is now the standard for AI and Machine Learning as shown by the many efforts at making this coupling an efficient environment[6][10][9], which brought an increasing interest for container driven

HPC aswell and Kubernetes for HPC in particular. Batch schedulers such as kube-batch<sup>5</sup> have been implemented for kube, and numerous HPC applications like slurm<sup>6</sup> now support containers as well.

Indeed, containers have many advantages that HPC users can benefit from. Here are some notable ones:

- First off, research has shown that Kuberenetes offer similar performance to more standard bare metal HPC[3].
- Users will get the same environment everywhere making up for a uniform and standardized work-place.
- Portability : users could seamlessly hop from one infrastructure to another based on their needs and criteria like price, performance, and capabilities rather than compatibility.
- Encapsulation : HPC applications often rely on complex dependencies that can be easily concealed into containers.

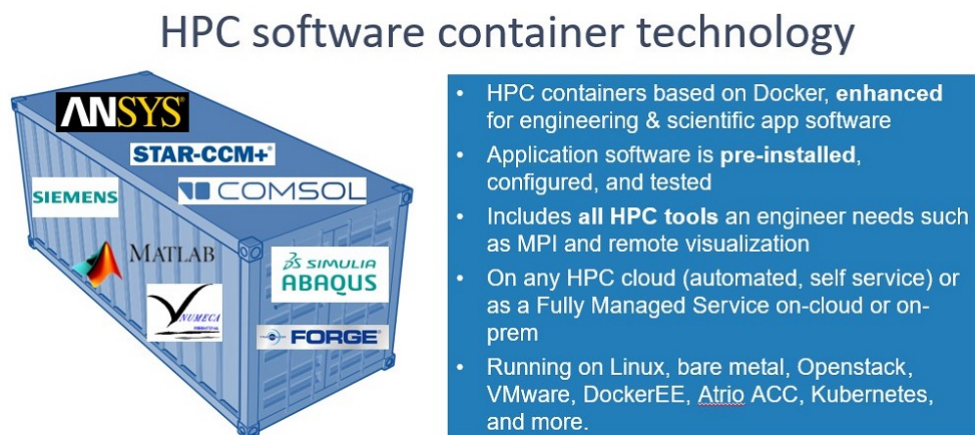


Figure 3.3: The container technology for HPC

**Source:** <https://www.hpcwire.com/2019/09/19/kubernetes-containers-and-hpc/>

Despite all those advantages, Kubernetes is not ready yet to be used in proper HPC environment because it lacks vital components like a proper batch job queuing system, and support for MPI applications. It cannot yet compete against the very well established HPC ecosystem, but that time may come soon as containers are becoming more and more integrated in modern infrastructures.

## 3.4 Objectives

The goal of this project is to design and implement Batkube, which will be an interface between Batsim and Kubernetes schedulers. With this interface, we want to compare Batsim results gainst data from a real Kubernetes cluster, given HPC workloads.

<sup>5</sup><https://github.com/kubernetes-sigs/kube-batch>

<sup>6</sup><https://slurm.schedmd.com/containers.html>

## 3.5 Technical challenges

### 3.5.1 Translation

### 3.5.2 Time hijack

TODO

---

#### Algorithm 1: Requester loop

---

```

Input: req: request channel, res: result channel map
1 while Batkube is not ready do
2   | wait
3 requests = []request
4 while req is not empty do
5   | m = <- req /* Non blocking receive */
6   | requests = append(requests, m)
7 sendToBatkube(requests) /* Only requests with duration > 0 are actually sent.
   Batkube will always answer. */
8 now = responseFromBatkube()
9 for m in range requests do
10  | res[m.id] <-now /* The caller continues execution upon reception */

```

---



---

#### Algorithm 2: Time request (time.now())

---

```

Result: Current simulation time
Input: d: timer duration, req: request channel, res: response channel map
Output: now : simulation time
1 if requester loop is not running then
2   | go runRequesterLoop() /* There can on ly be one loop runing at a time */
3 id = newUUID()
4 m = newRequestMessage(d, id) /* Requests are identified using uuids */
5 resChannel = newChannel()
6 res[id] = resChannel /* A channel is associated with each request */
7 req <- m /* The code blocks here until request is handled */
8 now = <-resChannel /* The code blocks here until response is sent by the
   requester loop */
9 return now

```

---

### 3.5.3 Re-building the API



— 4 —

## Evaluation





## Conclusion



# Bibliography

- [1] Anders Andrae. “Total consumer power consumption forecast”. In: *Nordic Digital Business Summit* 10 (2017).
- [2] Luiz André Barroso, Urs Hölzle, and Parthasarathy Ranganathan. “The datacenter as a computer: Designing warehouse-scale machines”. In: *Synthesis Lectures on Computer Architecture* 13.3 (2018), pp. i–189.
- [3] A. M. Beltre et al. “Enabling HPC Workloads on Cloud Infrastructure Using Kubernetes Container Orchestration Mechanisms”. In: *2019 IEEE/ACM International Workshop on Containers and New Orchestration Paradigms for Isolated Environments in HPC (CANOPIE-HPC)*. 2019, pp. 11–20.
- [4] Henri Casanova et al. “Versatile, Scalable, and Accurate Simulation of Distributed Applications and Platforms”. In: *Journal of Parallel and Distributed Computing* 74.10 (June 2014), pp. 2899–2917. URL: <http://hal.inria.fr/hal-01017319>.
- [5] Pierre-François Dutot et al. “Batsim: a Realistic Language-Independent Resources and Jobs Management Systems Simulator”. In: *20th Workshop on Job Scheduling Strategies for Parallel Processing*. Chicago, United States, May 2016. URL: <https://hal.archives-ouvertes.fr/hal-01333471>.
- [6] Mikyoung Lee, Sungho Shin, and Sa-Kwang Song. “Design on distributed deep learning platform with big data”. In: (2017).
- [7] Peter Brucker, Sigrid Kunst. *Complexity results for scheduling problems*. June 29, 2009. URL: <http://www2.informatik.uni-osnabrueck.de/knust/class/> (visited on 06/10/2020).
- [8] Millian Poquet. “Simulation approach for resource management”. Theses. Université Grenoble Alpes, Dec. 2017. URL: <https://tel.archives-ouvertes.fr/tel-01757245>.
- [9] Seetharami R. Seelam and Yubo Li. “Orchestrating Deep Learning Workloads on Distributed Infrastructure”. In: *Proceedings of the 1st Workshop on Distributed Infrastructures for Deep Learning. DIDL ’17*. Las Vegas, Nevada: Association for Computing Machinery, 2017, 9–10. ISBN: 9781450351690. DOI: 10.1145/3154842.3154845. URL: <https://doi.org/10.1145/3154842.3154845>.
- [10] Boris Tvaroska. “Deep Learning Lifecycle Management with Kubernetes, REST, and Python”. In: Santa Clara, CA: USENIX Association, May 2019.