

Implentation of an simplistic Interface for Big Data Workload Scheduler in Kubernetes

Implentation of an simplistic Interface for Big Data Workload Scheduler in Kuber- netes

Lukas Schwerdtfeger

A thesis submitted to the
Faculty of Electrical Engineering and Computer Science
of the
Technical University of Berlin
in partial fulfillment of the requirements for the degree
Bachelor Technische Informatik

Berlin, Germany
December 22, 2015



Main supervisor:

Prof. Dr. habil. Odej Kao, Technical University of Berlin

Hiermit erkläre ich, dass ich die vorliegende Arbeit selbstständig und eigenhändig sowie ohne unerlaubte fremde Hilfe und ausschließlich unter Verwendung der aufgeführten Quellen und Hilfsmittel angefertigt habe.

Berlin, den

Zusammenfassung

Kurze Zusammenfassung der Arbeit in 250 Wörtern.

Abstract

Short version of the thesis in 250 words.

Acknowledgements

This chapter is optional. First of all, I would like to...

Contents

List of Figures

List of Tables

1

Introduction

1.1 Motivation

The current population is producing more and more data. This creates an excellent opportunity for many businesses. Businesses willing to profit from collected data by using it to improve their sales strategies have to collect and store GigaBytes and upwards to ExaBytes of data. With storage costs becoming more affordable, companies are even less likely to toss away potential valuable data, creating so-called Data Lakes.

Collecting data is only the first step. It takes many stages of processing, through aggregation and filtering, to extract any meaningful information. Usually, the sheer mass of collected data makes it not very useful, to begin with.

Unfortunately, when working with ExaBytes of data, it is no longer feasible to work on a single machine. Especially when dealing with a stream of data produced by a production system and the information collected from yesterday's data is required the next day, or ideally immediately.

Scaling a single machine's resources to meet the demand is also not feasible. It is either very expensive or might just straight up not be possible. On the other hand, cheap commodity hardware allows the scaling of resources across multiple machines is much cheaper than investing in high-end hardware or even supercomputers.

The complexity of dealing with a distributed system can be reduced using the abstraction of a cluster. A **Cluster Resource Manager** is used, where a system of multiple machines forms a single coherent cluster that can be given tasks to.

Stream Processing of Data across such a cluster can carry out using Stream or Batch Processing Frameworks, such as Apache Spark or Apache Flink. These Frameworks already implement the quirks of dealing with distributed systems and thus hide the complexity.

The problem is that multiple Batch Jobs running on a single cluster need resources that need to be allocated across the cluster. While Resource Allocation is the Task of the Cluster Resource

Manager, the manager usually does not know how to allocate its resources optimal and often requires user (*TODO: 1*) to specify the resources that should be allocated per job. This usually leads to either too little resources being allocated per job, starving jobs and increasing the runtime, or more often over-committing resources and thus leaving resources in the cluster unused.

Another problem that arises is the fact that even though the reoccurring nature of Batch Jobs, not all Batch Jobs use the same amount of Resources. Some are more computationally intensive and require more time on the CPU, while others are more memory intensive and require more or faster access to the machine's memory. Others are heavy on the I/O usage and use most of the system's disk or network devices. This shows up in vastly different Job runtime (also total runtime) depending on the Scheduling of Batch Jobs across the Cluster.

Finding an intelligent Scheduling Algorithm that can identify reoccurring Jobs and estimate their resource usage based on collected Metrics and thus create optimal scheduling is not an easy task. It also requires a lot of setup when dealing with a Cluster Resource Manager.

1.1.0.1 TODO:

1. not just a user but the cluster user which is submitting the job

1.1.0.2 Open:

- ◇ How much detail is required here?

1.2 Problem Description

(TODO: Cluster Resource Manager, like YARN, were focused around Batch-Application because they existed because of Apache Hadoop/Map-Reduce ecosystem)

Cluster Resource Managers, like YARN, emerging from Apache Hadoop, were centered around Batch Frameworks.

With the rise of Cloud Computing and all the benefits that come with it, companies were quick to adopt new cloud computing concepts. The Concept of a Cluster Resource Manager introduced a notion of simplicity to those developing applications for the cloud. A Cluster Resource Manager now managed many aspects that used to be handled by dedicated Operations-Teams.

Kubernetes, a Cluster Resource Manager that was initially developed by Google, after years of internal use, provided an all-around approach to Cluster Resource Management for not just Batch-Application. The global adoption of Kubernetes by many leading companies, led to the growth of the ecosystem around it. Kubernetes has grown a lot since and has become the new industry standard, benefiting from a vast community.

Old Batch-Application-focused Cluster Resource Managers that used to be the industry standard are being pushed away by Kubernetes. Unfortunately, vastly different Interfaces or Scheduling Mechanism between other Cluster Resource Managers usually block the continuation of existing research done in the field of Batch Scheduling Algorithms.

Finding an efficient scheduling algorithm is a complex topic in itself. Usually, the setup required to further research existing scheduling algorithms is substantial. Dealing with different Cluster Resource Manager further complicates continuing on already existing work.

1.2.0.1 TODO:

1.2.0.2 Open:

- ◇ This sections contains a lot of text that may be better suited to the introduction section, but i don't really now what else to put in here
- ◇ Not happy with the ending of this chapter, like introduction it's really only one paragraph at the end that explains the intended contribution of this work

1.3 Goal of this Thesis

To aid further research in the topic of Batch-Scheduling-Algorithms, the goal of this thesis is to provide a simplistic interface for Batch-Scheduling on Kubernetes.

(TODO: Explain how already existing Schedulers like Mary and Hugo do not run on Kubernetes due to different interface/interactions)

Already existing Scheduling Algorithms, like Mary and Hugo, were initially developed for the Cluster Resource Manager YARN. Reusing existing Scheduling Algorithms on the nowadays broadly adopted Cluster Resource Manager Kubernetes is not a trivial task due to vastly different interfaces and interaction with the Cluster Manager.

(TODO: Explain why the Setup of Kubernetes has become easier: Cloud Providers, MiniKube)

Extending existing research to the more popular Resource Manager Kubernetes provides multiple benefits.

1. Research on Scheduling Algorithms for YARN has become less valuable due to less usage
2. The large ecosystem around Kubernetes allows for a better development environment due to debugging and diagnostic tooling
3. Initial setup of a Kubernetes cluster has become smaller due to applications like MiniKube, which allows a quick setup of a cluster in the local machine and Cloud Providers offering Kubernetes Clusters as a service.

(TODO: Describe the Interface here)

The interface should provide easy access to the Kubernetes Cluster, allowing an External-Scheduler to place enqueued Batch-Jobs in predefined slots inside the cluster.

For an External-Scheduler to form a scheduling decision, the interface should provide an overview of the current cluster situation containing:

1. Information about empty or in use slots in the cluster
2. Information about Jobs in the Queue
3. Information about the history of reoccurring Jobs, like runtime

It should be possible for an External-Scheduler to form a scheduling decision based on a queue of jobs and metrics collected from the cluster. The interface should accept the scheduling decision and translate it into Kubernetes concepts to establish the desired scheduling in the cluster.

(TODO: Explain shortcomings of Kubernetes)

Currently, the Kubernetes Cluster Resource Manager does not offer the concept of a Queue. Submitting jobs to the cluster would either allocate resources immediately or produce an error due to missing resources.

Kubernetes does not offer the concept of dedicated Slots for Applications either. While there are various mechanisms to influence the placement of specific applications on specific Nodes, these might become unreliable on a busy cluster and require a deep understanding of Kubernetes concepts, thus creating a barrier for future research.

1.3.0.1 TODO:

- ◇ Implementation of easy to use Interface that would allow already Batch Job Scheduling Algorithms likes Hugo and Mary to be run with small changes, on the popular Cluster Management Software Kubernetes

1.3.0.2 OPEN:

- ◇ Use of “should”. Okay? or Rather what it does?

1.4 Structure of this Thesis

The structure of this thesis allows the reader to read it in any order. To guide the reader through this thesis, the structure of this Thesis section will briefly explain which section contains which information.

The Background Chapter is supposed to give a brief overview of this thesis’s underlying concepts. This chapter introduces Big Data Streaming Processing, the Cluster Resource Manager Kubernetes, and Scheduling.

Following the Background chapter, the thesis provides an overview of the approach taken to tackle the problem described in the Problem Description Section. The Approach Section focuses on more profound concepts of Kubernetes and the Scheduling Cycle of the Kubernetes Scheduler. It summarizes the Kubernetes Operator pattern, which is commonly used to extend Kubernetes.

Implementation details will be given inside the Implementation Chapter, where an architectural overview and interaction between individual components are explained. The Implementation section also emphasizes the design Process for the Interface, which is exposed to an External-Scheduler. A significant part of the implementation is the Operator, which will be discussed extensively. The Implementation chapter shows how the points made inside the Approach Chapter are were implemented in the end. Finally, as the Goal of this Thesis section describes, changes that had to be made to already existing Scheduling Algorithm Implementations are disclosed and discussed.

(TODO: Hard to describe what is going to happen inside the Evaluation, if i don't have anything to evaluate yet)

An Evaluation of the research and contribution done by this thesis will be presented inside the evaluation chapter. Here its functionality is demonstrated. This section will also outline some of the limitations.

Before concluding the thesis, a comparison between State of the Art Technology for Kubernetes and Non-Kubernetes Scheduling Frameworks is made.

1.4.0.1 TODO:

- ◇ Thesis starts by giving a brief background to Big Data Streaming Processing, Cluster Management Systems (Kubernetes), and Scheduling
- ◇ Discuss the Approach this thesis takes on tackling the Problem Description, by explaining how scheduling in Kubernetes works and what it takes to Extend Kubernetes (using the Operator Pattern)
- ◇ Implementation Details that a worth mentioning:
 - An architectural Overview.
 - The Process of designing an Interface
 - The Operator that is used to extend Kubernetes
 - Changes that had to be made to existing Algorithms (and their tests)
- ◇ How the work of thesis is evaluated, by testing it's functionality, comparing results from previous work and finally outlining its limitations
- ◇ Comparing the Work that was done to current State of the Art Technology like the Batch Scheduling Framework Volcano and comparing to Scheduling approaches that are not available on Kubernetes
- ◇ A final Conclusion, with a note on future work, that is missing from the current implementation or requires rethinking.

2

Background

2.1 Distributed Dataflow Applications

Distributed Dataflow Processing aims to solve the problem of analyzing large quantities of data. In the last years, the amount of data that is being generated has exploded. This creates a Problem where single machines can no longer analyze the data in a meaningful time. While the Big Data Processing frameworks still work on single machines, computation is usually distributed across many processes running on hundreds of machines to analyze the data in an acceptable time.

Analyzing data on a single Machine is usually limited by the resources available on a single machine. Unfortunately, increasing the resources of a single machine is either not feasible from a cost standpoint or simply impossible. There is only a limited amount of Processor time, Memory, and IO available. Cheap commodity hardware allows a cluster to bypass the limitations of a single machine, scaling to a point where the cluster can keep up with the generated data and once again analyze data in a meaningful time frame.

Dealing with distributed systems is a complex topic in itself. Many assumptions that could be made in a single process context are no longer valid. Scaling to more machines increases the probability of failures. Distributed Systems need to be designed to be resilient against hardware-failures, network outages/partitions and are expected to recover from said failures. Having a single loss resulting in no or an invalid result will not scale to systems of hundreds of machines, where it is unlikely not to encounter a single failure during execution.

Scaling Computational Resources beyond the limitation of single machines is not a new problem. HPC (High-Performance-Computing), like Super Computers, are already commonly used for computational heavy use cases. Distributed Dataflow Applications initially used on static clusters, the advance of cloud computing, and Cluster Management Systems allowed them to scale more dynamically, depending on the users' needs. Distributed Dataflow Applications may be classically categorized as HPC Systems bringing high performance and efficiency due to sophisticated scheduling, where the Cloud offers Resiliency, elasticity, portability, and manageability.[?]

Distributed Dataflow Frameworks can be put into two categories, although many fall in both categories. Batch Processing and Stream Processing. In Batch Processing, data size is usually known in advance, whereas Stream Processing expects new data to be streamed in from different sources during the Runtime. Batch Processing Jobs will complete their calculation eventually, and Stream Processing, on the other hand, can run for an infinite time frame.

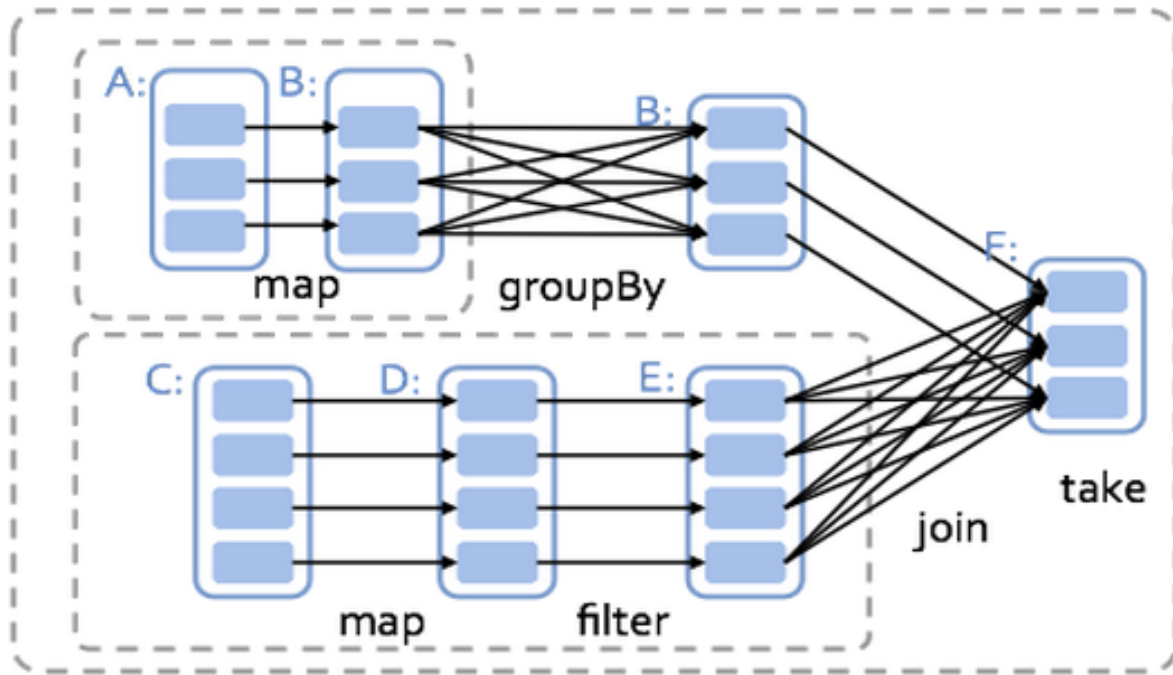


Figure 2.1: Example of Spark DAG [?]

Internally, Big Data Processing Frameworks build a directed acyclic graph (DAG) of Stages required for the analysis. Stages are critical for saving intermediate results, to resume after failure, and are usually steps during the analysis, where data moving across processes is required. A commonly used programming model is the MapReduce Programming model[?], where stages can be generalized in Map and Reduce Operations.

Map operations can be performed on many machines in parallel, without further knowledge about the complete data sets, like extracting the username from an Application-Log-Line. Reduce Operations require moving data around the cluster. These are usually used to aggregate data, like grouping by a key, summing, or counting.

Datasets that the Distributed Processing Frameworks analyze are usually in the range of Terabytes which is multiple magnitudes higher than the amount of Memory that each Machine has available. Partitioning of the Data is required due to the limitations of each single Machine, while persistent Storage, like Hard-Drives, might be closer to the extent of BigData, computation will quickly become limited by the amount of I/O a single machine can perform. Commonly a Distributed Dataflow Application is combined with a distributed file system (like HDFS[?]) that spans a shared Filesystem across many Nodes. In a distributed filesystem, the filesystem is usually broken down into smaller chunks, where every Node holds a few chunks of the complete filesystem as part of their local filesystem. If data from the distributed filesystem is not local to

the Nodes' filesystem, data will be moved across the network. Replicating chunks increases fault tolerance and thus prevents the loss of data but also requires synchronization across replicated chunks once they get updated. Data-Locality describes how close the data is, which the current task requires. If a Task has good Data-Locality on a Node, no data movement is required.

The user of Big Data Processing frameworks is usually not required to think about how an efficient partition of the data across many processes may be accomplished. Frameworks are designed in a way where they can efficiently distribute a large amount of work across many machines.

2.2 Scheduling

In general, scheduling is the process of assigning resources to a task. This includes the question:

1. Should any resources be allocated for the task at all?
2. At which point in time should resource be allocated?
3. How many resources should be allocated?
4. Which of the available resources should be allocated?

Scheduling is essential for Operating Systems that need to decide which process should get CPU time and which processes may need to wait to continue computation. In the case of multiple CPUs, a decision has to be made on which CPU should carry out the computation. The Operating System is not just concerned with CPU-Resources, but also I/O Device resources. Some devices may not work under concurrent usage and require synchronization.

In some cases, a simple FIFO scheduling that works on tasks in order there were submitted produces acceptable results. Scheduling depends on a goal. Some algorithms aim to find the optimal schedule to respect any given deadlines. Whereas some distinguish between Soft and Hard Deadlines, where ideally no deadlines would be missed at all, occasionally missing soft deadlines to guarantee Hard deadlines are met is acceptable. In general, finding a single best schedule that allows resources to be allocated optimally is not possible. Scheduling for a fast response time or throughput might prefer shorter tasks to be run, when possible, and might starve longer running tasks for a long time before progress can be made.

Common goals of scheduling are [?]:

- ◇ Minimizing the time between submitting a task and finishing the task. A Task should not stay in the queue for a long time and start running soon after submission.
- ◇ Maximizing resource utilization. Resources available to the scheduler should be used, even if that means skipping tasks waiting in the queue.
- ◇ Maximizing throughput. Finishing as many jobs as possible may starve longer running jobs in favor of multiple shorter running jobs.

To achieve their goals scheduling algorithm might allow preemption, where the currently active task could be preempted for another task to become active. Some scheduling algorithms account for the potential overhead of preempting the current task (like a context switch).

Scheduling can be applied at multiple levels of the Software Stack, like scheduling processes at the Operating System level, Virtual Machines at the Hypervisor level, or scheduling tasks across

many applications running on many machines at the cluster level. The higher up the stack, the more and more potential schedules become possible. It is a balancing act of complexity and performance. At the highest level, the scheduler has the most knowledge about the systems and could come up with optimal scheduling. However, it seems to be a wise choice for scheduling to be handled in their respective stack layer since a lot of complexity can be encapsulated inside each layer. Allowing the Top Level Scheduler to micromanage the complete stack creates an explosion of complexity.

For this work, the scope of scheduling is limited to Batch Scheduling of Jobs in a cluster.

The Question of Scheduling in a Distributed System is now the question of which machines resources should be used for which job. In Batch Scheduling algorithms need to pay attention to the characteristics of a Distributed System:

1. Potential heterogeneity of the system, with machines of different Hardware and different Operating Systems or Software
2. Spontaneously adding and removing resources of the Cluster
3. Interference between Applications residing on the same machine, same rack, same network switch, etc. (CO-Location)

While some of these factors can be controlled, different algorithms can be chosen for various use cases.

In Distributed Dataflow Applications, we deal with considerable different levels of scheduling. At the Framework level, applications build a DAG-based execution plan based on the job submitted. The initial DAG breaks down the job into their respective Map and Reduce Operations. These Operations will be broken down further into smaller Tasks based on the Partitioning of Data. Finally, Tasks are executed on an arbitrary number of processes across different machines. Optimizing the schedule of tasks to an executor process will be called DAG-Level scheduling and may now also include factors like Data-Locality.

Moving Up one Level Higher in the stack, we are concerned with running multiple Jobs inside the same cluster, and a decision needs to be made which job can spawn their executor on which Nodes. Executors are packaged on Containers. The containers are isolated so that they can not access each other.

A more profound introduction into the specifics of scheduling in Kubernetes is inside the Scheduling in Kubernetes chapter.

Cluster Resource Manager like Mesos[?] allows scheduling on multiple levels, with Resource Offers. The top-level Mesos scheduler finds multiple candidates and offers them to the DAG-Level scheduler framework. The potential for cross-level scheduling would allow the Top-level scheduler to respect data-locality constraints for the DAG-level scheduler.

2.3 Cluster Management Systems

In HPC Super Computers are commonly build with expansive hardwares, and tightly coupled Operating Systems, that allows parallel applications to take advantage of resources available on the super computer. Super Computers are carefully planned with software designed for specific

purposes, this usually meant that multiple applications running on a super computers, would have resources statically partitioned among them. Coarse grain partitioning introduces inefficient use of hardware, once part of the statically partitioned resources are no longer in use.

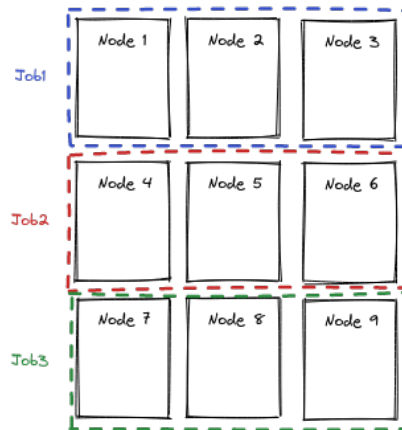


Figure 2.2: Static Partitioning

With the emerge of Distributed Dataflow Applications, like Hadoop MapReduce[?] the need for a more fine grain partitioning of cluster resources came along. The move away from Super Computers to cheap commodity hardware, meant that cluster can be scaled up or down easily, where previously careful planning was required. Managing a dynamic system of potential thousands of Nodes, can not be done in a manual fashion. A Cluster Resource Manager is required, to build the abstraction of a single cohesive cluster, that can be tasked with jobs.

Hadoop is one of Many Open-Source MapReduce implementations. Cluster computing using commodity hardware was driven by the need to keep up with the explosion of data. The Initial Version of Hadoop was focused on Running MapReduce Jobs to process a Web Crawl [?]. Despite the initial focus, Hadoop was widely adopted evolved to a state where it was no longer used with its initial target in mind. Wide adoptions have shown some of the weaknesses in Hadoops architecture: - Tight Coupling between the MapReduce Programming model and Cluster Management - Centralized Handling of Jobs will prevent Hadoop from Scaling

The tight coupling leads Hadoop users with different applications to abuse the MapReduce Programming model to benefit from cluster management and be left with a suboptimal solution. A typical pattern was to submit ‘map-only’ jobs that act as arbitrary software running on top of the resource manager. [?] The other solution was to create new Frameworks. This caused the invention of many frameworks that aim to solve distributed computation on a cluster for many different kinds of applications [?]. Frameworks tend to be strongly tailored to simplify solving specific problems on a cluster of computers and thus speed up the exploration of data. It was expected for many more frameworks to be created, as none of them will offer an optimal solution for all applications.[?]

Initially, Frameworks like MapReduce created and managed the cluster, which only allowed a single Application across many machines—running only a single application across a cluster of Nodes led to underutilization of the cluster’s resources. The Next generation of Hadoop allowed it to build ad-hoc clusters, using Torque [?] and Maui, on a shared pool of hardware.

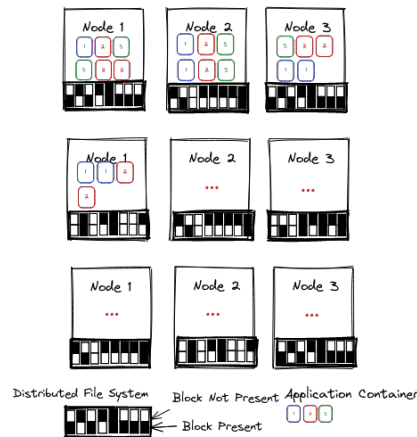


Figure 2.3: Dynamic Partitioning

Hadoop on Demand (HoD) allowed users to submit their jobs to the cluster, estimating how many machines are required to fulfill the job. Torque would enqueue the job and reserve enough machines once they become available. Torque/Maui would then start the Hadoop Master and the Slaves, subsequently spawn the Task and JobTracker that make up a MapReduce Application. Once all Tasks are finished, the acquired machines are released back to the shared resource pool. This can create potentially wasteful scenarios, where only a single Reduce Task is left, but many machines are still reserved for the cluster. Usually, Resources requirements are overestimated and thus leaves many clusters resources unused.

With the Introduction of Hadoop 3, the MapReduce Framework was split into the dedicated Resource Manager YARN and MapReduce itself. Now MapReduce was no longer running across a Cluster, but it was running on top of YARN, who manages the cluster beneath. This allows MapReduce to run multiple Jobs across the same YARN Cluster, but more importantly, it also allows other Frameworks to run on top of YARN. This moved a lot of complexity away from MapReduce and allowed the framework to only specialize on the MapReduce Programming Model rather than managing a cluster. This finally allows different Programming Models for Machine Learning tasks that tend to perform worse on the MapReduce Programming model [?] to run on top of the same cluster as other MapReduce Jobs.

Using the ResourceManager YARN allows for a more fine-grain partitioning of the cluster resources. Previously a static partitioning of the clusters resources was done to ensure that a specific application could use a particular number of machines. Many applications can significantly benefit from being scaled out across the cluster rather than being limited to only a few machines.

- Fault tolerance: Frameworks use replication to be resilient in the case of machine failure. Having many of the replicas across a small number of machines defeats the purpose
- Resource Utilization: Frameworks can scale dynamically and thus use resources that are not currently used by other Applications. This allows applications to scale out rapidly once new Nodes become available
- Data Locality: Usually, the data to work on is shared across a cluster. Many applications are likely to work on the same set of data. Having applications sitting on only a few Nodes, but the data to be shared across the complete cluster leads to a lousy data locality. Many unnecessary I/O needs to be performed to move data across the cluster.

Fine-grain partitioning can be achieved using containerization, where previously applications

were deployed in VMs, which would host a full Operating System. Many containers can be deployed on the same VM (or physical machine) and share the same Operating System Kernel. The hosting Operating makes sure Applications running inside containers are isolated by limiting their resources and access to the underlying Operating System.

Before Hadoop 3 with YARN was published, an Alternative Cluster Manager Mesos was publicized. Like YARN, Mesos allowed a more fine granular sharing of resources using containerization. The key difference between YARN and Mesos is how resources are scheduled to frameworks running inside the cluster. YARN offers a resource request mechanism, where applications can request the resources they want to use (*TODO: fact check*), and Mesos, on the other hand, offers frameworks resources that they can use. This allows frameworks to decide better which of the resources may be more beneficial. This enables Mesos to pass the resource-intensive task of online scheduling to the frameworks and improve its scalability.

Kubernetes[?] was initially developed by Google and released after multiple years of intern usage. Kubernetes was quickly adopted and has become the defacto standard for managing a cluster of machines. Kubernetes offers a descriptive way of managing resources in a cluster where manifests describing the cluster's desired state are stored inside a distributed key-value store[?] etcd. Controllers are running inside the cluster to monitor these manifests and do the required actions to bring the cluster into the desired state. Working with manifest abstracts away many of the problems that arise when deploying Applications to a cluster. Usually, an Operations Team was required to manage applications across the cluster. With Kubernetes offering the required building blocks and the mechanism of a control-loop, the operator pattern in combination with Custom Resource Definitions is commonly used to extend Kubernetes functionalities. Where Hadoop's YARN was focusing on Distributed Dataflow Applications, Kubernetes enables developers of all kind of Applications the scalability and resilience of the Cloud.

3

Approach

The Approach section gives a more profound overview of core concepts of Kubernetes, which intern are the building blocks of the work done in this thesis. Starting with the Operator Pattern, which is strongly connected with the method, Kubernetes achieves its cluster orchestration, the *control loop*. The second section details how scheduling works in Kubernetes and how it can be manipulated to be controlled by an External-Scheduler using the interface.

3.1 Extending Kubernetes using the Operator Pattern

The smallest unit of *deployment* in Kubernetes is a Pod[?]. While a Pod may consist of multiple containers, containers in a Pod are guaranteed to run on the same Node. Containers in Pods also share storage and network resources across the container boundary. In General, containers inside the same Pod are tightly coupled and commonly used in a sidecar pattern to extend the main container with common functionalities across the cluster, like the Kube-RBAC-Proxy[?] that is frequently used with Containers that interact with the Kubernetes API and require authorization.

Usually, in Kubernetes, Pods are not created by themselves but are managed by resources that build on top of them. Most commonly, Pods are used in Combination with Jobs, Deployments, or Statefulsets, which control the Lifecycle of the Pod.

The Operator pattern is based on the already existing design used by Kubernetes native Resources like Deployments the *Control-Loop*. Resources like Jobs, Deployments[?], and Statefulsets[?] describe Pods' desired state. Essentially Kubernetes has controllers that monitor changes to the Resources Manifest and the current cluster situation. The Control-Loop (or Reconciler-Loop) allows a Deployment consisting of a Pod template with a Replication-Factor to the exact amounts of Pods running inside the cluster. If any of the Pods fails, the Deployment Controller creates a new one. The controller ensures that new Pods are created first and become ready before old Pods are deleted. But on the Flip-Side, the controller also knows when the Resources Manifest is updated. For example, if the container is updated to a more recent version, all containers need to be replaced with the newer version. Deployment Resources have many policies that dictate

how these actions should happen. Usually, restarting all Pods at once would not be desired so that the deployment would allow for Rolling-Upgrades.

Using multiple controllers, that each intern only is concerned with a single resource creates a very loosely coupled and thus highly extendible interface. The Idea of a reconciliation controller loop improves the resiliency of a system by observing the current state, comparing it to the desired state, and taking actions for the observed state to converge into the desired state.[?]

The Operator Pattern is commonly used to extend Kubernetes Functionalities. Most of the time, however, just creating a new controller is not enough to extend Kubernetes. It usually also requires Custom Resource Definitions. Custom Resource Definitions (CRD) is the Kubernetes way of defining new resources that are allowed to exist in the cluster. Having multiple Controllers listing to the same Resource, like a Deployment, makes little sense or could even cause issues. Thus the Combination of a new Resource and a Controller that knows how to handle it creates the *Operator Pattern*. The term *Operator* is used as the controller is designed to replace previously manual work of configuring Kubernetes native Resources done by an Operator. An everyday use case for the Operator Pattern is to Control Applications at a Higher level, where previously Multiple Deployments and Services may have been required to operate a Database. The Operator Pattern could reduce that to just a single Manifest containing the meaningful configuration. Operators can thus be created by Experts operating the Software and be used by any Kubernetes Cluster.

Common best practices when developing Operators are that operators should only manage a single kind of resources, and multiple operators should be used if multiple resources need to be managed. Rather than micro-managing resources down to the creation of Pods, operators should delegate to existing operators where ever possible, e.g., using Kubernetes native resources, like Deployments. Multiple operators should not act on the same CRD, as this requires synchronization between operators and enforces tight coupling between components. Since all operators work with the Kubernetes API, many identical boiler code has been written over the years, and many best practices are contained inside opinionated frameworks. [?]

3.2 Scheduling in Kubernetes

In this section, the Scheduling Model of Kubernetes will be introduced as it is a vital part of the implementation for the Interface.

The scheduling Problem in Kubernetes is the problem of deciding which Pods are running on which Node. For some Pods, the question can be easily solved. For example, Pods controlled by a DaemonSet are by the specification of the DaemonSet running on every Node. Without further information, a feasible choice of scheduling Pods onto Nodes seems to be simple round-robin scheduling. Every Pod that requires scheduling gets scheduled onto the next Nodes until all Nodes have Pods then the cycle is repeated. However, both Pods and Nodes can influence the scheduling.

Pods can specify the resources they are going to use and may even set a Hard Requirement in the form of a request for resources they require to run. Scheduling needs to take the resources requests into account when scheduling a Pod across Nodes. Pods can also directly influence the Node they should be scheduled on, either through specifying a *nodeName* directly, a

NodeSelector, which identifies possible multiple Nodes, or a more general concept of *affinity*. Affinities provide the ability to set hard and soft requirements, where a Pod may become unschedulable if a hard requirement cannot be met, and a soft requirement is not preferred but an acceptable decision. With Affinities, even inter-pod Affinities can be specified where Pods can choose not to be scheduled on a Node where another Pod is already deployed.

On the other side, Nodes can also specify so-called Taints, where only Pods with the fitting Toleration can be scheduled onto the Node. Kubernetes also use the Taint mechanism to taint Nodes that are affected by Network problems, Memory-Pressure, or are not ready after a restart. The Taint specifies if just new Pods cannot be scheduled *NoSchedule* onto the Node, or if even already running Pods should be evicted *NoExecute*. Tainted Nodes with either type of Taints are not chosen during scheduling.

However, Pods can bypass the Taint Node if they have the fitting Toleration. For example, if a Node becomes unreachable due to a network outage, Nodes affected will be tainted with the “*node.kubernetes.io/unreachable*” Taint, which is set to *NoExecute*. Pods can neither be scheduled onto the unreachable Nodes nor keep running without the fitting Toleration. In a likely scenario where the network recovers and the machine becomes reachable again, it would be inefficient to restart all Pods instantly if a machine becomes unreachable for a short time. To prevent Pods from being evicted (and thus most likely restarted), Kubernetes by default creates the unreachable Toleration for all Pods. However, the Toleration set by Kubernetes also specifies a maximum duration where the taint is tolerated (by default, 5 minutes). If a node stays unreachable for more than 5 minutes, all Pods are evicted.

In the context of scheduling, Taints, Tolerations, Pods, and Nodes only refer to the objects stored inside the distributed key-value store. If a physical machine becomes unreachable, the objects manifest is updated. Whether the physical containers are still running or not cannot be determined, but for the duration of the Toleration, they do count as running and are not replaced with new Pods.

3.2.1 Scheduling Cycle

Kubernetes Scheduling is based on the Kubernetes Scheduling Framework. As the extensible nature of Kubernetes, the Scheduling framework is also extensible. Customization to the scheduling can be done at so-called Extension Points.

Each attempt of scheduling Pods to Nodes is split into two phases, the *Scheduling Cycle*, and the *Binding Cycle*.

The Scheduling Cycle is more interesting in the context of this work. In this part of the cycle, the decision is made on which Node to schedule which Pod. The Binding Cycle is the scheduling phase, where the Pod is being deployed to the physical machine. Multiple Scheduling Cycles will not run concurrently, as this would require synchronization between multiple Scheduling Cycles. Since the Scheduling Cycle is relatively fast, executing multiple Scheduling Cycles in parallel seems unnecessary. However, the Binding Cycle requiring the deployment of Pods is compared to the Scheduling Cycle rather long-living and thus may be executed in parallel. Both the Scheduling Cycle and the Binding Cycle may abort the scheduling of a Pod in case it may be unschedulable.

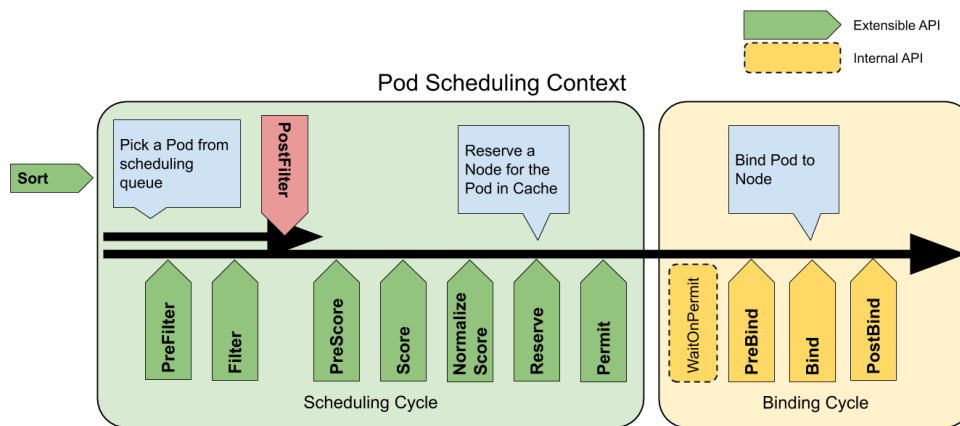


Figure 3.1: Scheduling Cycle

The Pod Scheduling Context can be modified through plugins that intern use the Scheduling Framework’s extension points. The graphic ?? shows the available extension points, which are the Pod Scheduling Cycle phases. Multiple Plugins can be registered for any of the different extension points and will run in order.

This section will give a brief breakdown of the different stages and highlight those critical in the implementation of the External-Scheduler-Interface.

Queue Sort: Sometimes, multiple Pods require a scheduling decision. The Queue Sort Extension point extends the Scheduling Cycle with a comparison function that allows the Queue of Pods to wait for a scheduling decision to be sorted. Usually, multiple Plugins can influence the scheduling decision at a time; however, having numerous Plugins sorting the queue of waiting Pods will not result in anything meaningful. Thus, only a single plugin can control the queue at a time. After sorting the Queue of Pods, the first Pod is chosen and passed onto the next stage.

Filter: The Scheduling decision is made by first using the PreFilter and Filter extension point. The Filter Extension point filters out Nodes that cannot run the Pod. Most extension points also have a pre-extension point used to prepare information about the Pod. As mentioned earlier, any extension point can return an error indicating that the Pod may be unschedulable.

The Filter extension point also has a **Post-Filter** Extension, which is only called if no node passes the Filter, and a Pod becomes *unschedulable*. The Post-Filter Extension can then be used to find a possible Scheduling using preemption. Preemption is very important, as it is required to model Pods having a higher priority than other Pods. In this work, preemption is required to guarantee that the scheduling from an External-Scheduler is correctly applied even if other Pods in the cluster exist that are unknown to the External-Scheduler.

Scoring: The filter plugin is concerned with hard requirements that prevent Pods from being scheduled. If previous filter plugins deem multiple Nodes suited for the Pod to be scheduled on, a decision needs to be made to find the best Node. The scoring plugin considers soft requirements and gives Nodes a lower or higher score that can or cannot fulfill the soft requirements. Finally, all Scores are normalized between two fixed values.

Reserve: The reserve phase is used to reserve resources on a Node for a given Pod. This is

required due to the asynchronous transition into the Binding Cycle. Pods that are supposed to be bound to a Node invoke reservation to prevent possible race conditions between future Scheduling Cycles. Reservation may fail, in which case the cycle moves to the unreserve phase, and all previously completed reservations are revoked. Usually, the Reservation extension point is used for applications that will not use the default containerization mechanism of Kubernetes and rely on different binding mechanisms.

Permit: The final stage of the Scheduling Cycle is the Permit Stage, which ultimately denies or delays the binding of a Pod in a case where binding might not be possible or still requires time.

3.2.2 Extending the Scheduler

Common ways of extending the Scheduler are either to implement a custom Scheduler and to replace the KubeScheduler in the cluster, or to use the *Extender* API that instructs the Scheduler to invoke an external API for its scheduling decision.

The pluggable architecture of the Scheduler allows plugins to extend the scheduling cycle at the extension points. Different Scheduler profiles can use different plugins. Profiles can be created or modified using a *KubeSchedulerConfiguration*. Only Pods that specify the Scheduler Profile using *.spec.schedulerName* are scheduled using the profiles.

The Extender Mechanism describes an external Application waiting for HTTP Requests issued during the Scheduling Cycle to influence the scheduling decision. Instructing the Scheduler to use an Extender, in the current Implementation of Kubernetes, is not limited to a specific Profile, but all Profile will use the configured Extender. An Extender can specify the Verbs it supports. Verbs in this context refer to Filtering, Scoring, Preempting, and Binding.

Note: This work does not replace the KubeScheduler, but deploys a Second Scheduler to the cluster with a Profile that uses the Extender API. This is done because of the limitation that every scheduling profile will use the Extender, which seems unnecessary in the context of a prototype.

Implementation

4.1 Designing the Interface

The External-Scheduler-Interface allows an External-Scheduler to control the resources Batch Applications like Apache Spark and Apache Flink will use for their TaskManager and Executor Pods. No Assumptions are made about the Driver and JobManager Pods. The Interface offers slots in the form of Testbeds to the External-Scheduler. The size (both in terms of the number of slots and resources) of testbeds are controlled via the Testbeds manifest. Managing the Testbed is not exposed through the Interface, as Testbeds are not expected to be changed by an External-Scheduler, except for using its slots.

The concept of a testbed allows the reservation of resources inside a cluster. It may be unreasonable to give an External-Scheduler all available resources in a multi-tenant cluster. The resources available to the External-Scheduler can be precisely controlled using Testbeds. Schedulings created by the External-Scheduler are always directed towards a Testbed. Currently, an active scheduling will claim its Testbed and prevent other Schedulings from using it.

The External-Scheduler-Interface offers endpoints for querying the current cluster situation in the form of the Testbeds slots occupation status and the status of Schedulings or Batch Jobs. The Interface provides a REST API that allows the creation and deletion of Schedulings and the ability to update information stored inside the Batch Job manifest. An additional Stomp [?] (WebSocket) server is available to not enforce any polling for updates.

The functionality is demonstrated as part of the evaluation, either through the Manual Scheduler, which acts as a testing tool and as visualization, and the Example Scheduler, which uses multiple Testbeds to profile Batch Jobs for its scheduling decisions.

4.2 Architecture

In this section, all components that make up the Interface are introduced. Here an Architectural overview is presented, and interactions between components are discussed.

The current implementation of the Scheduler Interface consists of 5 Components that will be introduced in this section but discussed in more detail in the Operator Section.

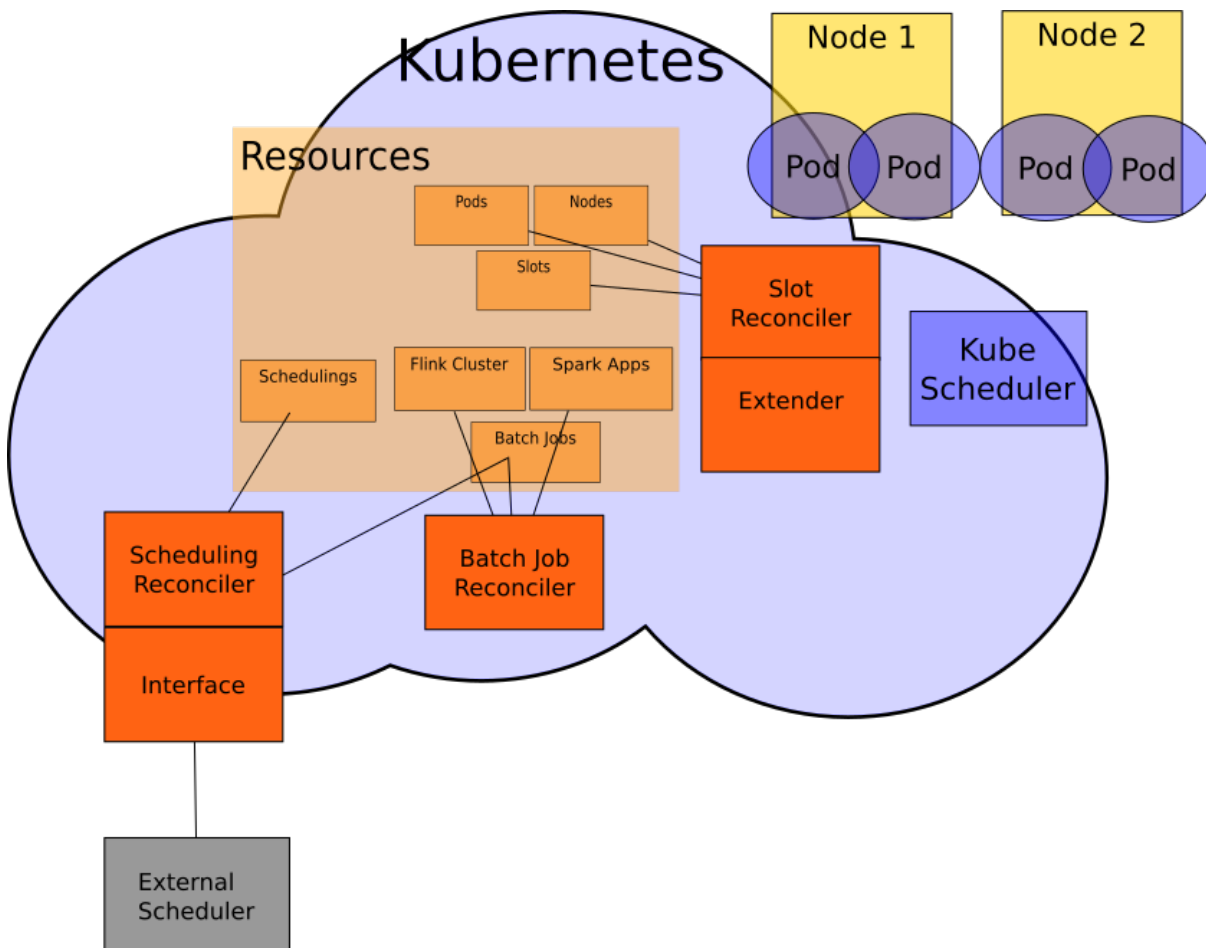


Figure 4.1: Components

The five components consist of three Reconciler or Control-Loops, the Batch-Job Reconciler, the Slots Reconciler, and the Scheduling Reconciler. The architecture also uses an Extender and, finally, the External-Scheduler facing Web-API.

The Interface that is visible to an External-Scheduler is supposed to be simple and only allows the querying of the current cluster situation, information about previous scheduling, and submission of new schedulings. Further, more concrete information, like Node metrics, can be queried from a Metric Provider that is commonly deployed along with the Cluster.

Additionally to the three Reconciler, three Custom Resource Definitions (CRDs) are created.

- ◇ *Batch Job* represents an Abstract Batch Job Application and can store information the External-Scheduler may want to remember for future invocations

- ◇ *Slots* represent a testbed of guaranteed resources available for a *Scheduling*. A *Slots* Custom Resource is a collection of slots across the clusters Node that are referenced in a *Scheduling*
- ◇ *Scheduling* represents the decision done by the External-Scheduler. A *Scheduling* maps multiple *Batch Jobs* to *Slots* available in the Cluster. The *Scheduling* also acts as a Queue and submits Jobs into the slots in order once *Slots* become available.

The Batch Job CR is used to Model a reoccurring Batch Job Application. In order to support both Flink and Spark Applications, an abstract Batch Job CR is chosen that maps the state of application-specific CR (link SparkApplications[?] and FlinkCluster[?]) to a common set of possible States. (Information about possible States and the corresponding State Machine are discussed in the Batch Job Operator Section). A Batch Job can be claimed by exactly one *Scheduling*, this is because the Batch Job CR models exactly the life cycle of a single application.

The Slots CR guarantees Resources in the Cluster by creating Ghost Pods, with a specific Resource Request representing the Size of an empty Slot. The Ghost Pods reserve resources by not allowing other Pods requiring scheduling to be scheduled onto the same Node. Using the extender and Preemption, the Slot Reconciler can reserve resources for Pods created by the Batch Job CR.

Finally, the Scheduling CR is passed to the external Interface by the External-Scheduler. Given a Set of Batch Jobs and the Slots and the Node they exist on, an External-Scheduler can compute a *Scheduling* that chooses Batch Jobs and the Slots they should run in.

Note: Reconciler and Control-Loop can be used interchangeably, but for less confusion with the Spring Boot Concept of a Controller, the principle of a Control-Loop will be called Reconciler**

4.3 Operator

4.3.1 Batch Job Operator

The Batch Job Operator comprises the Batch Job Reconciler and the Batch Job CRD. The Reconciler is listing for changes regarding the Batch Jobs CRs and Applications CRs, which are managed by the Spark Operator and the Flink Operator.

The Batch Job Operator knows how to construct the corresponding application given the Batch Job CRs specification. With the Spark and Flink Operator, reusing existing software allows the Batch Job CR to be only a thin wrapper around either a Spark or a Flink specification. In addition to the Spark and Flink CR, it may contain additional information that previous invocations of the External-Scheduler have stored.

Currently, the Batch Job CR only contains a partial application-specific specification. The Batch Job CR requires a user to specify only the Required Components, like Application Image containing the actual application and its arguments like a dataset or where to find it (e.g., using HDFS). Although it would not matter, if a user would submit a fully specified Flink or Spark Application, the Operator would overwrite most of the Driver/Executor Pod specific configuration and replication configuration.

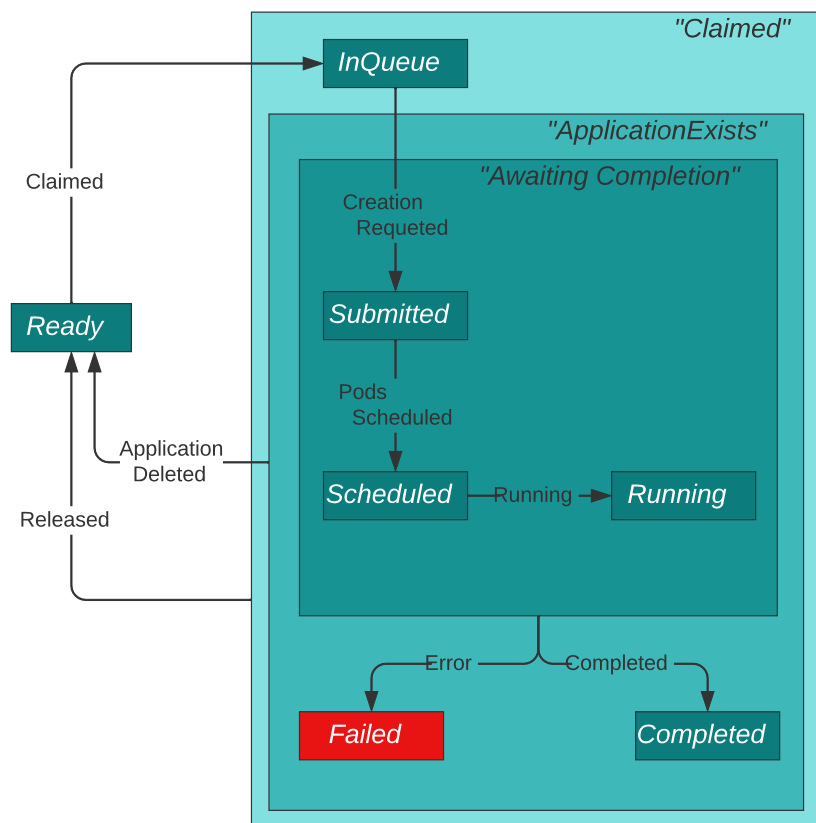


Figure 4.2: StateMachine

The Batch Job Reconciler is implemented as a nested state machine with anonymous sub-states. The Approach was chosen as it creates more comprehensible software, which can be split easier into Components and handle edge cases by design.

Initially, Batch Jobs submitted to the cluster remain in the Ready State. While in the Ready State, the Batch Job Reconciler will not do anything. A *Scheduling* can acquire a Batch Job. The Batch Job CR will move into the InQueueState until the *Scheduling* instructs the Batch Job Operator to create the application and track its lifecycle or releases it.

Communication between the Batch Job Reconciler and the Scheduling Reconciler is done via the *Batch Jobs* spec (*.spec.activeScheduling* and *.spec.creationRequest*). If *scheduling* wants to claim a Job, it updates the active scheduling spec. This mechanism ensures that only one *Scheduling* at a time can use the *Batch Job*. On the flip side, a *Scheduling* can claim multiple *Batch Jobs*. Suppose the active *scheduling* releases the Job; by removing the activeScheduling spec, the *Batch Job* moves back into the Ready State. Releasing a Job could happen at any time and may even cause any created application to be removed.

Once a Batch Job is in the InQueue State, the Reconciler waits for the creation request issued by the *Scheduling* Reconciler. The request is again done using the *Batch Jobs* spec and specifies desired replication and the *TestBed* and slots the application should use.

When configuring the application to be created by the corresponding Operator, there are two

types of configuration. Configuration can either be:

- ◇ Persisted inside the Batch Job CR, which is used on every invocation of the application. This includes the Applications Image and arguments, like the data set
- ◇ Scheduling dependent. These configurations can not be stored inside the CR and must be supplied with the creation request.

After a Batch Job was requested to create the application, application-specific logic is executed. In any case, the actual steps for deploying the applications to the cluster are done by the Applications Operator (Flink Operator[?] or Spark Operator[?]). The Batch Job Reconciler only instructs the Application Operators with configurations for the Executor/TaskManager Pods, so they are identifiable to the Extender.

When creating the application, the following aspects are configured for the Executor/TaskManager Pods:

- ◇ **Resource Requests:** The Container resources are specified by the Testbeds slot size. For the Pods to fit inside a slot, they need the correct Resource Request. (Currently only CPU and Memory)
- ◇ **Slot IDs:** The Scheduling (or the External-Scheduler) decides which slots are used by which Job. For the Executor/TaskManager Pods to be placed into the correct slot (technically the correct Node), Pods need to be made identifiable by the Scheduler Extender.
- ◇ **replication:** The Number of Executor/TaskManager Pods depends on the Number of Slots that will be used for the application.
- ◇ **Priority Class:** Application Pods need a *PriorityClass* otherwise, preemption will not be triggered by the Kubernetes Scheduler.
- ◇ **Scheduler Name:** Application Pods need a *SchedulerName* otherwise, the default KubeScheduler will handle the Scheduling and thus ignore the Scheduling Extender.

Configuration of **Resource Requests, Replication** is straightforward, as both the Spark and Flink Operator expose these via their respective CRDs. The Spark Operator actually exposes the complete PodSpec for both driver and executor Pods, whereas the Flink Operator only exposes a few PodSpec attributes. The Flink Operator had to be extended with the missing configurations. This way **Resource Requests, Replication, Priority Class, Scheduler Name** are configured.

The difference between any of the mentioned above configurations and the **Slot IDs** is that the Application Operators only allow (rightfully so) to specify a single Pod spec. This is because the Executor/TaskManager Pods are controlled by a Stateful Set, which scales up to the desired replication. However, the configurations mentioned above are valid for all Pods, but *Slot IDs* need to be different.

This issue can be circumvented by leaving the final decision of which Pod goes into which slot to the Extender—submitting only a list of all SlotIDs to the Extender. The Extender needs to decide which Pod goes into which slot. Pods are configured with an affinity of the combined set of Nodes where the slots reside on.

Once the application was created, the Job moves into the Submission State. It resides there until

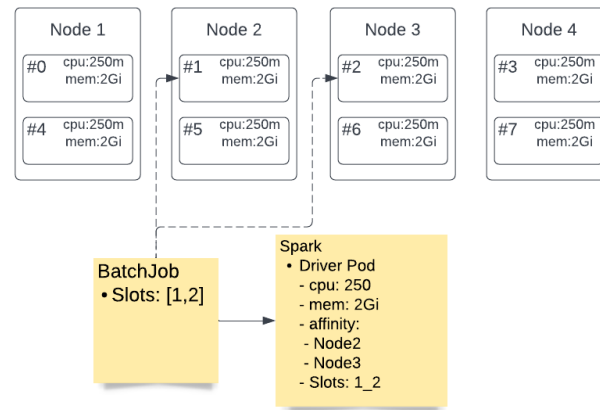


Figure 4.3: Affinities

all Pods where scheduled, at which point it moves on into the Running State. The underlying Applications state is monitored until it moves into the application-specific completed state (Spark: *Completed* and Flink: *Stopped*). During the implementation, scenarios were encountered in which the Batch Job reconciler was not running. Once restarted, found Applications in a completed state without passing the Scheduling, submission, or running state. To prevent any tight coupling, none of these transitions are required to be considered a successful execution.

The Batch Job reconciler tracks the time an application ran by creating timestamps once it started running and its completion.

4.3.2 TestBed Operator

The *TestBed* Operator comprises the Reconciler Loop and the *TestBeds* CRD. The *TestBeds* CR is supposed to model a Collection of slots located in a Cluster of Machines. Slots can have specified Resources. While no application is running inside a slot, it is considered *free*. To reserve resources in the cluster and thus guarantee applications supposed to be deployed inside a *free* slot actually to get the resources, the *TestBed* Operator needs to:

- ◇ **Reserve Resources** by using so-called Ghost Pods inside the cluster that specify a resource request and thus reserve the resources
- ◇ **Preempt** Ghost Pods for Pods that wants to be deployed inside a slot

The TestBed Reconciler listens to changes to the TestBed CR and the current cluster situation. It ensures that the correct number of Pods with the specified resource requests are always deployed onto the cluster. The TestBed CR is composed of the following configurations:

- ◇ Label Name to Identify any Nodes that are part of the TestBed. Only the Label Name is specified, not a specific value. The value is later used to create a distinct order of slots in the cluster.
- ◇ Number of slots per Node
- ◇ Resource Request per slot

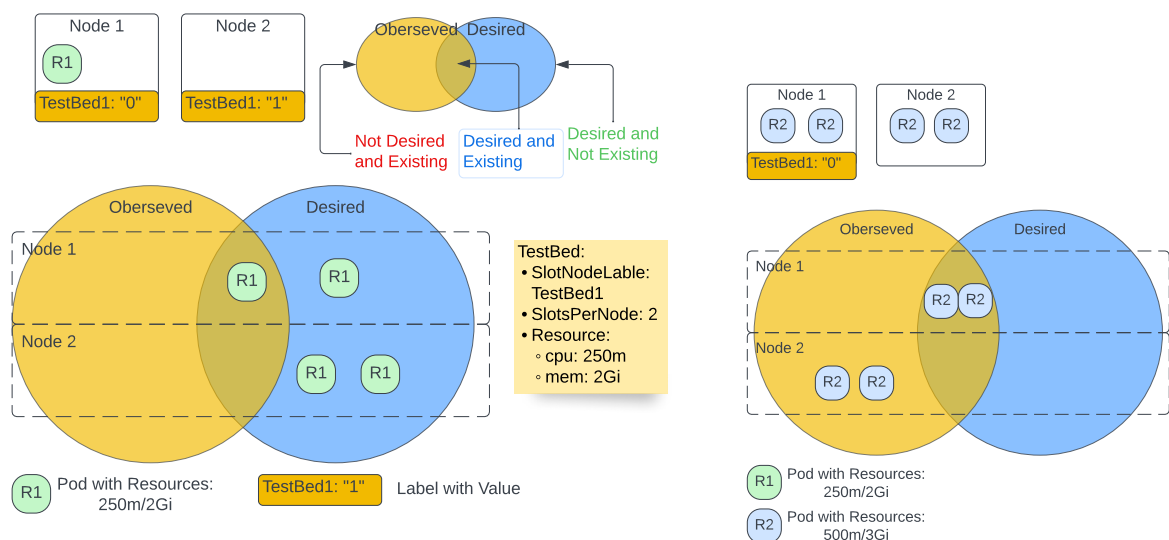
Given the Testbeds specification, the Reconciler listens for all changes to Nodes **with** the specified label. It also needs to listen for Nodes **without** any label if the label was removed and the Testbed needs to be resized. Further, it also listens to changes to any Pod part of the Test Bed.

The typical Reconciliation Loop works as follows:

- ◇ Fetch the current cluster situation
- ◇ Calculate the desired cluster situation
- ◇ Find the difference. Either delete undesired Pods or create desired Pods

Fetch the current cluster situation by fetching all Pods with the **SLOT** label. Pods are then grouped by their Node, thus creating a list of Pods per Node. The desired state is calculated by modeling Pods for every slot and grouping them by Nodes. When comparing Pods, we consider them equal if they reside on the same Node, have the same resource request, and have the same *SlotPositionOnNode*.

Note: The position of slots on a Node does not matter because slots on a Node are only a logical abstraction.



(a) New Pods need to be Created

(b) Node Change: Pods need to be Deleted

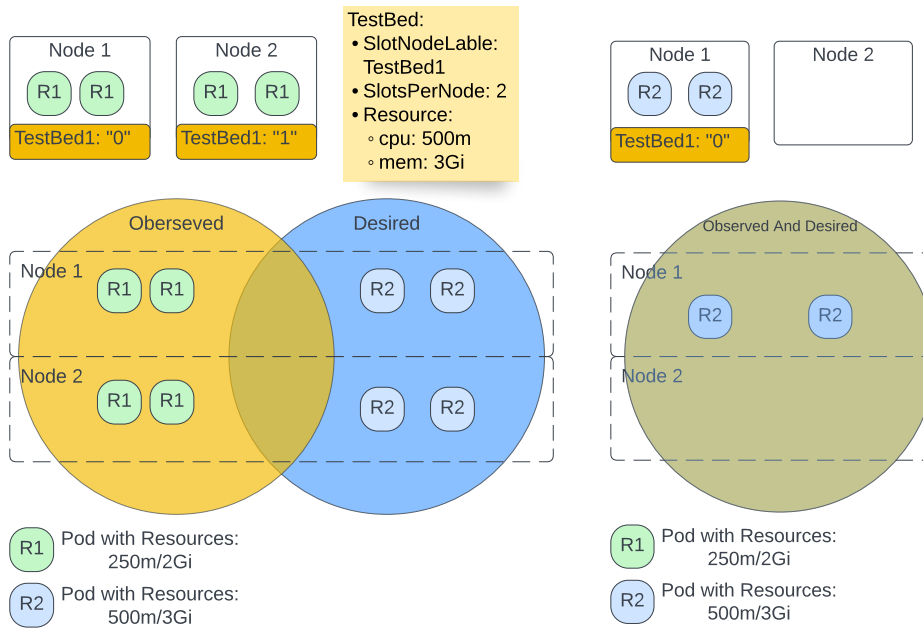


Figure 4.5: TestBed Observed and Desired

The Reconciler now builds a set of observed Pods and a set of desired Pods. ?? shows an example scenario where the control-loop realizes that Pods from the desired state are not in the current state, thus creating the missing Pods in the *desired and not existing* set. In a different scenario displayed by ?? the label on a Node was removed, thus reducing the number of slots inside the Testbed. Pods that are in the *existing and not desired* set will be removed. The final set is the *desired and existing* set, which contains Pods that already have the correct resources requirement and are placed on the correct Node.

Currently, the SlotOccupationStatus holds the following information:

- ◇ **NodeID** and **NodeName**: which is derived from the Test-Bed Selector Label on the Node
- ◇ **Position**: which is the SlotID,
- ◇ **slotPositionOnNode**: where the position does unique among the whole Test Bed, SlotPosition on Node is only unique per Node
- ◇ **PodName** and **PodUID**: The Name and the Unique Identifier of a Pod that is currently residing inside the slot
- ◇ **state**: is the current state of the slot, which can either be *free*, *reserved*, or *occupied*

4.3.3 Extender

Extender Component is integrated within the TestBed Reconciler. Suppose the reconciliation loop detects that the cluster is in progress. The loop is aborted to prevent changes from the Testbed Reconciler and the Extender to act concurrently on the TestBed CR. Currently, a cluster is considered in progress if any of the Pods require Scheduling (.spec.NodeName is not set) or are terminating (deletion timestamp is set).

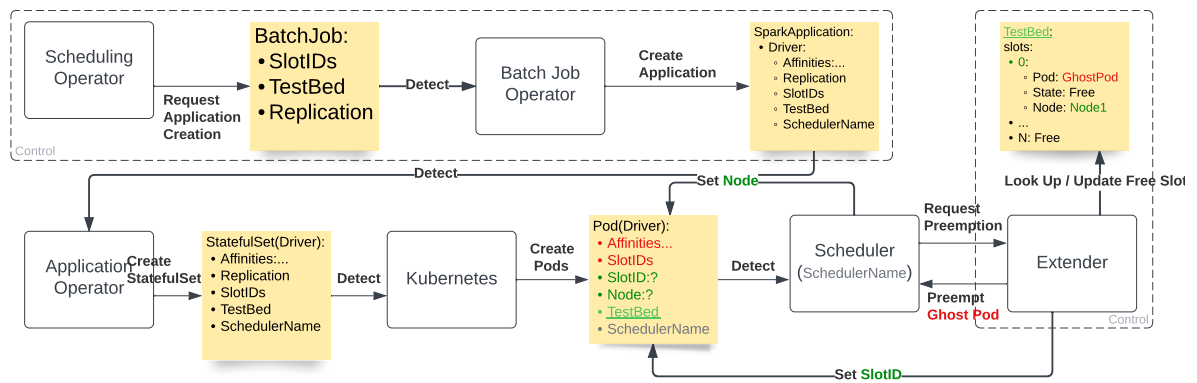


Figure 4.6: Components under control of the External-Interface System

The Extender is the component that directly interacts with the Kubernetes Scheduler. An additional scheduler, with an additional scheduling profile, is running concurrently to the default Kube-Scheduler. The custom scheduler (which will be referred to as Kube-Scheduler) is configured to use the Extender. To guarantee the Scheduling of Pods onto the TestBeds slots, the Extender extends the Filter and Preemption extension points of the Kubernetes scheduling cycle. The main problem the Extender can solve, is that the Batch Job Operator does not have full control of Pods created downstream by the Applications Operator.

?? shows which of the Components and Resource Managed by them are under the control of the External-Interfaces System. The Batch Job Operator can only control the Application CR created. The Application CR only describes a single PodSpec, which will be later be replicated into multiple Pods by the Replication Controller, which is part of the StatefulSet. Thus it is not possible to set Pod specific configurations, like the SlotID, at the Batch Job Operator Level. However all Pods can be configured, with enough information, for the Extender to figure out which Pod belongs in which slot.

Note: PodSpec here only refers to the TaskManager/Executor PodSpec, as the External-Interface does not handle Scheduling of the JobManager/Driver Pods.

In order to influence the scheduler to schedule Pods onto Nodes with the correct slot, the number of possible Nodes is first limited by all Nodes containing any of the slots using affinities.

The Applications PodSpec created by the Batch Job Operator, limits the possible Nodes, the kube-scheduler can use affinities. During development multiple scenarios, of interaction between kube-scheduler and Extender were identified.

If the kube-scheduler, detects not enough available resources, it will trigger preemption. If no preemption is required the Filter endpoint of the Extender is queried, to further limit the possible Nodes. ?? describes how the set of possible Nodes is first shrunk by affinities and later the correct Node is chosen by the Extender, based on the first free slot. Once a Pod invokes the Extender a slot will be reserved. The Extender Stores Information inside the TestBed CRs Status. To prevent race conditions between the TestBed Reconciler and the Extender mutual access to the TestBed CR, is required, which is guaranteed since the TestBed Reconciler will abort its reconciliation if Scheduling is in progress.

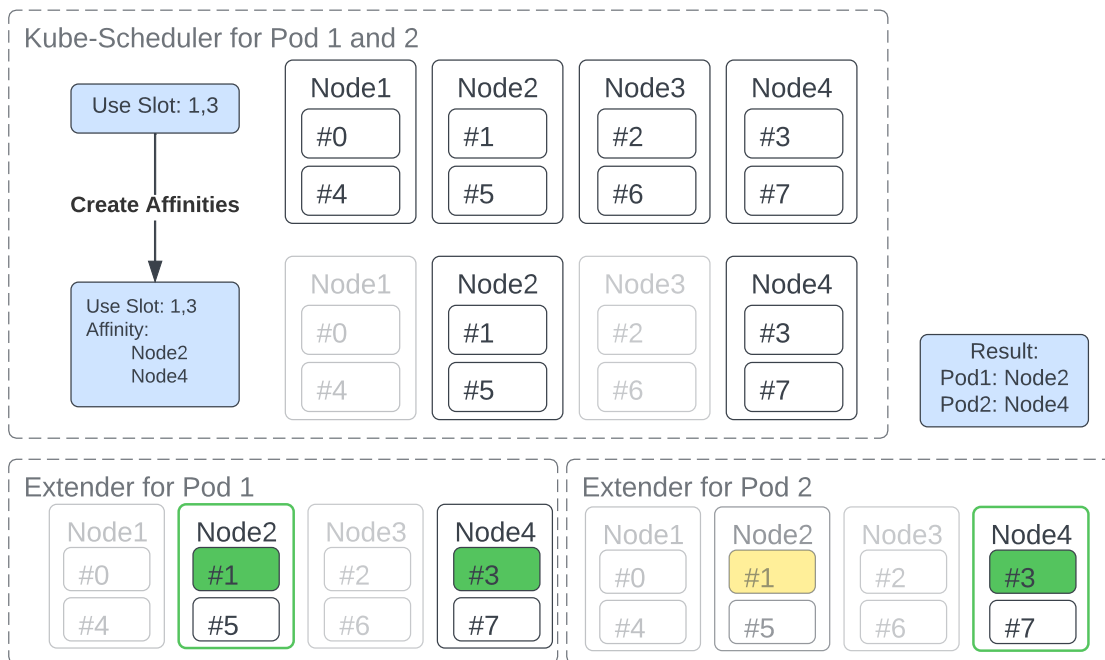


Figure 4.7: Kube-Scheduler limits Nodes, Extender selects Node with slot

In the Event no resources are available the kube-scheduler invokes the preemption endpoint of the Extender. Preemption is straightforward, since the TestBed already reserves resources using ghost Pods. The Extender chooses a slot based on the first free slot, and returns the currently residing ghost Pod to be the victim.

For the TestBed Reconciler to detect scheduled Pods the Extender will in addition to returning the target Node to the scheduler, also sets the **SlotID + NonGhostPod** of scheduled Pods. Once Scheduling has finished, the TestBed Reconciler will no longer abort execution, and preempt all ghost Pods, that were not previously preempted by the Kube-Scheduler. Further Information about the Extender can be found inside the source code, as it contains heavily commented code, to describe different scenarios.

4.3.4 Scheduling Operator

Once again the Scheduling Operator is composed of the Reconciler Loop and the Scheduling CRD. The Scheduling Reconciler like the Batch Job Reconciler implemented using a nested Statemachine. The Scheduling CR models, a collection of jobs, and a slot selection strategy. Once submitted the Reconciler first acquires all jobs, and starts running them. The Scheduling tracks the execution of all its jobs and submits new Jobs once old jobs have finished and slots become available again. Initially the Scheduling was planned to only support offline Scheduling, where an External-Scheduler plans the execution of multiple jobs in advance, however in theory updating the Scheduling spec would allow an online scheduling, but in the current state it is rather unreliable, as it only allows jobs to be added to the end of the queue.

Slot selection strategy do not aim to provide a full scheduling algorithm, they are just means for an External-Scheduler to describe which Job should use which slot.

- ◇ Implemented as a state machine
- ◇ Acquire State claims all Batch Jobs and the Test Bed
- ◇ Once all Jobs are in the InQueueState scheduling choses the first n runnable jobs
- ◇ Two Modes: SlotBased + QueueBased (Images)
- ◇ Once Creation was Requested, Reconciler waits until all jobs submitted were scheduled.
- ◇ This is required, because the Slot Reservation is not instantaneous. Wait for Batch Job Reconciler + Application Operator until the Extender marked slots as reserved.
- ◇ At this point the Scheduling waits until slots come available, different Modes require different Condition
- ◇ The Queue Based Scheduling, only requires a number of available slots
- ◇ Slot Base scheduling requires specific slots to come available
- ◇ Once the Queue is empty the Scheduling moves into the await completion state until all jobs have completed
- ◇ Note: Online Scheduling: is possible by updating the scheduling CR and extending the Queue

4.3.5 External-Scheduler-Interface

- ◇ Interaction with the Scheduling Interface is naturally done via the Kubernetes API, creating, updating, deleting CRs.
- ◇ If the external-scheduler chooses not to directly interact with kubernetes, a thin layer in form of web api is provided.
- ◇ The interface aims to abstract away some of the Kubernetes features like namespaces.
- ◇ The interface allows to create update and delete schedulings. Query for jobs inside the cluster. Query for slots inside the cluster.
- ◇ The interface contains a web socket server that broadcasts changes to jobs, schedulings, Testbed

4.4 Changes to existing Algorithm

5

Evaluation

5.1 Testing

To test the functionalities of the External-Scheduler-Interface, a combination of Unit Tests and Integration Tests are used.

Unit tests are highly specific towards smaller components of the System and thus will not fit in the scope of the written thesis, however they can be found inside the repository (Appendix).

Integration Tests aim to test the bigger picture. In a system composed of many distributed processes, testing all functionalities is not feasible, as setup would require recreating a cluster of software components. Testing the complete system with all its components, is only possible interacting with an established cluster. The Manual scheduler aims to ease the use of External-Scheduler-Interface.

For Integration Test to run in a timely manner, a common practice is to *Mock* Software components, which are either not under immediate control or are very expensive to setup. The Java Operator SDK and the Fabric8 Kubernetes client, can Mock the Kubernetes Cluster. Mocking the entirety of Kubernetes can not be accomplished, but the Kube-API-Server itself is enough to test the Operators. With the Kubernetes API-Server mocked, we can simulate the abstract state (declared state) and verify that changes to the abstract state trigger the correct actions by the Reconcilers.

Overview of Software Components used in the Integration Tests: - Kubernetes API: Mocked by Fabric8 - Batch Job Operator: Actual Reconciler + CRDs are used - Application Operator (Spark): Mocking changes to the SparkApplication - Scheduling Operator: Actual Reconciler + CRDs are used - (Scheduler): Out of Scope for this work - TestBed Operator: Actual Reconciler + CRDs are used. The actual Slot Occupation Status is mocked

Using integration tests, basic functionality of the Reconciler components, at least on the Abstract Resource level can be verified.

To test the External-Scheduler-Interface on an established Cluster with all components installed, the Manual Scheduler Frontend can be used.

5.1.1 Manual Scheduler

The Manual Scheduler is a web frontend that interacts with the Interface, and allows a user of the Interface to create a scheduling and test it.

The frontend both acts as a reference on how the Scheduler-Interface is supposed to be used, and to verify the functionality of the Cluster.

5.2 Example Scheduling Algorithm

The Interfaces usability is evaluated with an exemplary implementation of a non-trivial scheduling algorithm. Since the Interfaces can manage multiple TestBeds inside the same cluster, a Profiling-Scheduler approach is chosen to highlight some of the interface's features.

The example scheduling algorithm is used on a 5 Node cluster running inside the Google Cloud Platforms Kubernetes Engine (GKE). During development, smaller Nodes with a single vCPU and 8Gi of Memory were sufficient, but for the final evaluation, Nodes were doubled in capacity.

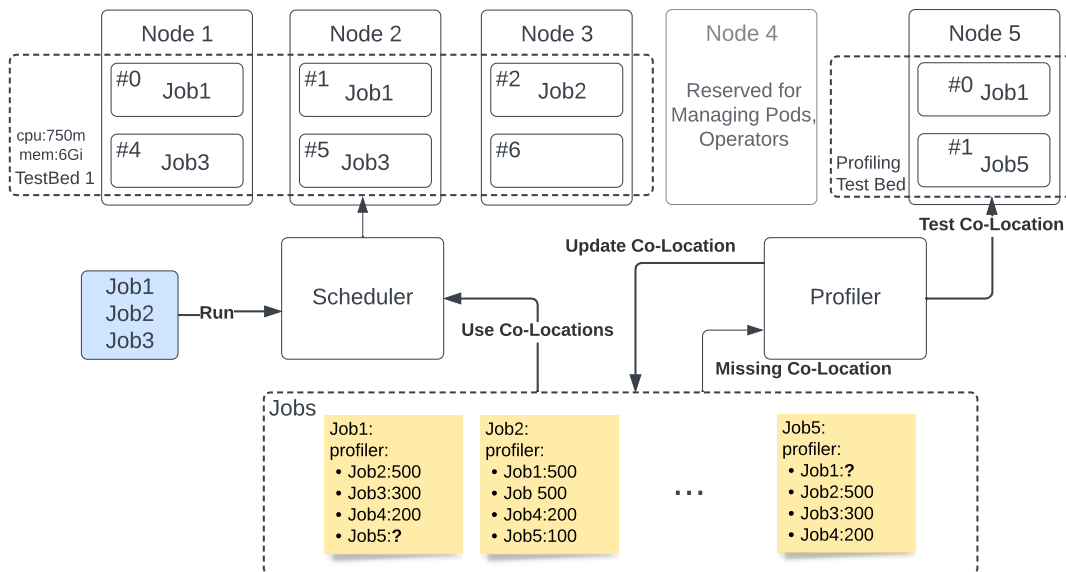


Figure 5.1: Architecture of the Example Profiler-Scheduler

Two Test-Beds are created using the Test-Bed CRD, the Profiling TestBed, and the TestBed for the actual execution of Jobs. The slot size was chosen depending on the resources available. For the final evaluation, a slot size of 750m CPU and 6Gi Memory leaves enough resources available for the cluster's control plane, managing Pods (Driver, JobManager), and the Operators running inside the cluster.

Contrary to the Scheduler Thread, which only creates a scheduling if requested (via stdin), the Profiler Thread runs at all times, updating and refining the Co-Location Matrix by choosing job pairings with the least data points.

Co-Locations are described as a simple runtime in seconds. By iterating over the available job pairings, the Profiler builds a Co-Location Matrix. The Profiler creates a cumulative moving average for each job pairing.

Jobs cannot be paired with themselves since the acquire/release mechanism only allows a single execution per job at a time. In theory, Jobs can be Co-Located with itself by deploying the application with a replication of two, but this could not be directly compared to Co-Location with a different job, as work done is split between both instances.

Note: The matrix is not symmetrical because the runtime of each job is used, not the runtime of both jobs (or the runtime of the complete scheduling). The runtime of the scheduling is the time after acquisition until all jobs have been completed. The Batch Job Operator only tracks the application's time inside the running state. This approach was chosen because applications may have vastly different startup times, which will become insignificant for long-running jobs.

The Scheduling Thread is the traditional scheduler. Given a list of Jobs, the Scheduling Thread tries to find optimal scheduling regarding total runtime. The scheduler takes a greedy approach choosing the co-located job based on the job with the shortest runtime to keep the evaluation simple. Replication of each job is selected based on the number of slots ($\text{Replication} = \text{NumberOfSlots} / \text{NumberOfJobs}$). Empty slots are again greedily filled with jobs suited best for co-location, not allowing a job to be chosen more than once.

Both the Profiler and the Scheduler run in parallel. If any of them cannot acquire their jobs, the scheduling will wait until they become available.

5.3 Limitations

The prototype implementation of the External-Scheduler-Interface comes along with limitations. This section will outline a few of them. The final Future Work Section will be a continuation of this section.

The intended use-case for the Interface is the development of new scheduling algorithms. Arbitrary-sized Testbeds can span across a single cluster controlled by Kubernetes. Kubernetes has been shown to scale with large clusters, and with Kubernetes v1.23 clusters with up to 5000 nodes are supported[?]. During the prototype development, a cluster with only a few nodes was used, thus leaving the scalability of the Extern-Scheduler-Interface in an uncertain state. It is unlikely for a single Testbed to overwhelm the operator due to limiting the access to *Testbeds* to a single *Scheduling*. Furthermore, Batch Jobs are rather long-living and do not require much maintenance (at least from the Interfaces perspective). However, there are currently no limitations on how many *Testbeds* may exist in the same cluster, and thus the Interface may be overwhelmed with too many concurrent *Testbeds*. Currently, the Interface has no leader election mechanism and therefore cannot scale horizontally. Using multiple operators in different namespaces does not require a leader election mechanism but would require refining the use of namespaces inside the Interface.

A Namespace in Kubernetes is the mechanism that isolates groups of resources within the same cluster. In a shared cluster, it is commonly used to not interfere with other users without revolving back to statically partitioned clusters and lose the benefit of resource sharing. The current implementation only supports limited use of namespaces since it was developed on a private cluster and could use the “default” namespace for all its components. The Integration Tests use a different namespace, and the Interface is not tied to a specific namespace. However, for now, all components run inside the same namespace. This leaves the question if actions done by the operator should be limited to a single namespace or if it should be able to interact with resources in every namespace. Furthermore, this question is not just limited to the operators, which are part of the prototype but also application-specific operators. The Spark-Operator currently allows both options, either watch all namespaces or limit the scope to only a single namespace.

The Interface currently does not expose the namespace to the External-Scheduler, because every resource must be inside the Interfaces namespace. This further promotes limiting the scope per instance of the Interface to a single namespace and thus limiting required privileges. However, some actions required by the Interface need access to resources across all namespaces. First of all, the *Testbed Operator* requires access to the nodes (not namespaced) and requires access to all pods running in every namespace to calculate the available resources before creating a *Testbed*. The issue with Nodes not being limited by a namespace is that multiple *Testbeds* may use the same node, which would undoubtedly cause trouble with the current implementation.

The last limitation that would prevent the Interface from being used in a *productionish* environment is the lack of Security. Both in terms of malicious Users (or just forgetful), who can reserve cluster resources with a Testbed, and accessing the cluster from outside via the API Endpoints, without proper authentication and authorization. The first issue is hard to circumvent because a malicious user could do that anyway, and Kubernetes Resource Quotas[?] can limit the resources per namespace, but they would also limit a non-malicious User. Automatically resizing a Testbed if it is not in use for a while could prevent unnecessary reservation of resources and defeat the purpose of reserving resources for an External-Scheduler. The Second issue requires API-Endpoints not to be publicly available. During development, the Interface was accessed using port-forwarding, which requires authorization to the Kubernetes cluster. For future work relying on port-forwarding seems plausible. However, creating an actual External-Scheduler down the line would require proper Ingress configuration to the cluster and a security model.

5.4 Discussion

6

State of the Art

6.1 Volcano

Volcano is a System Batch-Job Scheduler made for High-Performance Workloads on Kubernetes. Volcano extends Kubernetes with functionalities that Kubernetes do not natively support. Some of these functionalities are critical when working with High-Performance Workloads, like “PodGroups”. In a scenario where a Framework might want to create multiple Pods for its computation, the resources inside the cluster only allow for a few of them to be deployed. Applications could encounter deadlocks, requiring more Pods to be deployed to progress. The Concept of PodGroups prevents Pods from being scheduled unless all of them can be scheduled.

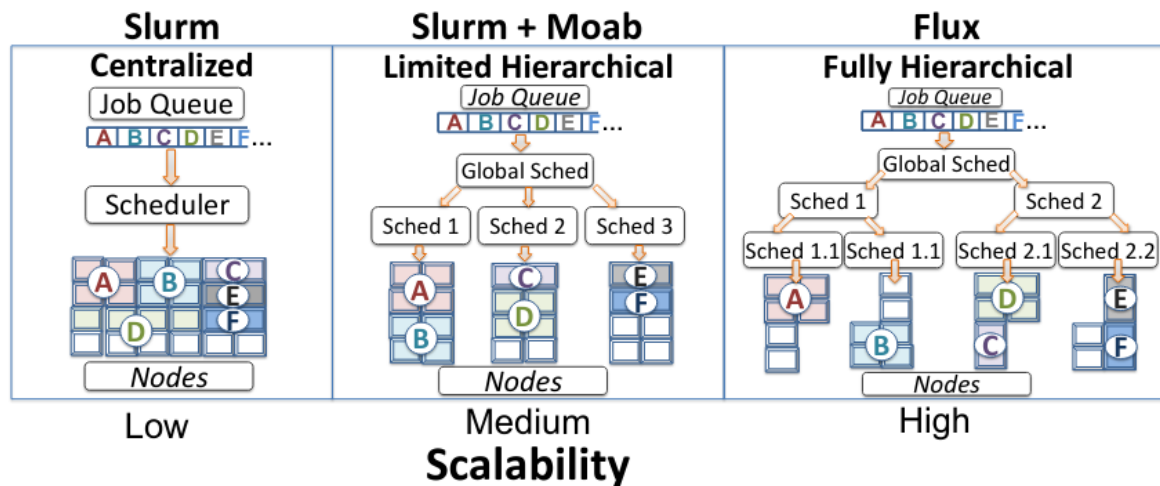
The Volcano scheduler is based on the Kubernetes Scheduling Framework, influencing the scheduling cycle at the extension points. This way, Volcano can implement many scheduling policies, which High-Performance Batch Applications commonly use. Volcano is more concerned with policies around the actual scheduling algorithm. In contrast, the External-Scheduling-Interface, introduced in this work, focuses on aiding the development of the actual scheduling algorithm. In theory, algorithms could be implemented using Volcano, but since Volcano provides just a thin layer above the Kubernetes Scheduling Framework, using Volcano for the development of new scheduling algorithms does not seem like a plausible choice.

With both Frameworks supporting Batch Scheduling in different ways, the External-Scheduling-Interface and Volcano could complement each other, but there has been no further investigation.

6.2 Flux

Flux aims to find a solution for Converged Computing, a Term used, when describing the move of traditional HPC computing to the Cloud Native Computing Model. HPC Systems traditionally bring high performance and efficiency due to sophisticated scheduling, where the Cloud offers Resiliency, elasticity, portability and manageability. Traditional HPC Batch Scheduler, will not keep up with the growth of systems enabled by the cloud. Established HPC Frameworks,

such as Slurm[?], use a centralized Scheduler. Flux identifies scalability issues, that the Scheduler is limited, in the maximum number of jobs it can handle. To Prevent the scheduler from overwhelming, job submission needs to be throttled, which will decrease job throughput. While not strictly related to scheduling, a centralized approach, will also fail at tracking the status of jobs running inside larger clusters. Sli



Flux introduces a new HPC Scheduling model to address the challenges, by using one common resource and job management framework at both system and application levels. Using an Hierarchical Scheduler applying the divide-and-conquer approach to scheduling in a large cluster. The hierarchical scheduling model, enables jobs to create their own schedulers, which is used to schedule its sub-jobs.

Another approach the Flux Scheduler takes, to scale with the increasing number of jobs in a cluster, is to aggregate jobs, that are similar and arrive within the same time-frame, into single larger job.

7

Conclusion and Future Work

7.1 Conclusion

7.2 Future Work