

Bachelor Thesis

My Bachelor Thesis

Submitted by

Marcel Pascal Stolin

32168

born at 03.04.1993 in Kamen

Written at

Fraunhofer-Institut für Produktionstechnik und Automatisierung IPA

and

Stuttgart Media University

First Examiner: Prof. Walter Kriha
Second Examiner: Prof. Dr.-Ing. Marco Huber
Supervisor: M.Sc. Christoph Hennebold

Eidesstattliche Erklärung

Ich versichere hiermit, dass ich, Marcel Pascal Stolin, die vorliegende Bachelor Thesis selbstständig angefertigt, keine anderen als die angegebenen Hilfsmittel benutzt und sowohl wörtliche als auch sinngemäß entlehnte Stellen als solche kenntlich gemacht habe. Die Arbeit hat in gleicher oder ähnlicher Form noch keiner anderen Prüfungsbehörde vorgelegen.

Ort, Datum: _____

Unterschrift: _____
Marcel Pascal Stolin

Zusammenfassung

Hier kommt eine deutschsprachige Zusammenfassung hin.

Abstract

Abstract in English.

List of Figures

2.1	Autonomic computing concept - Source: Authors own model, based on [JSAP04].	9
2.2	The control-loop concept - Source: Authors own model, based on [Mur04].	10
2.3	Managed resource - Source: Authors own model, based on [JSAP04].	10
2.4	Autonomic manager - Source: Authors own model, based on [JSAP04].	11
2.5	The monitoring process - Source: Authors own model.	13
2.6	Push-based monitoring approach - Source: Authors own model.	14
2.7	Pull-based monitoring approach - Source: Authors own model.	14
4.1	Docker architecture - Source: Authors own model, based on [Doc].	21
4.2	Docker basic container structure - Source: Authors own model, based on [BMDM20].	22
4.3	Optimization process of the Spark Catalyst - Source: Authors own model, based on [Luu18].	25
4.4	Overview of a Spark cluster architecture - Source: Authors own model, based on [Theb].	26
4.5	Spark's cluster mode - Source: Authors own model, based on [CZ18].	27
4.6	Spark's client mode - Source: Authors own model, based on [CZ18].	27
4.7	Catalyst optimization with RAPIDS accelerator for Apache Spark - Source: Authors own model, based on [McD20].	29
4.8	Prometheus high-level architecture - Source: Authors own model, based on [Thec, BP19].	30
4.9	Basic architecture of a GitLab CI/CD pipeline - Source: Authors own model, based on [Git].	33
5.1	Overall cluster architecture - Source: Authors own model.	38
5.2	Autonomic manager component design - Source: Authors own model.	40

5.3	UML activity model of the autonomic manager process -	
	Source: Authors own model.	41

List of Tables

6.1	Auto-Scaler configuration parameter	47
-----	---	----

Listings

2.1	Example of a dimensionless-metric	14
2.2	Example of a metric with dimensions	14
4.1	Usage of master launch script	27
4.2	Usage of worker launch script	28
4.3	Example usage of the spark-submit executable	28
6.1	Auto-Scaler configuration YAML file	46
6.2	Environment configuration for all worker nodes	50
6.3	Prometheus target configuration in YAML syntax	51
A.1	Apache Spark base image Dockerfile	59
A.2	Apache Spark master image Dockerfile	60
A.3	Apache Spark worker image Dockerfile	60
A.4	GPU discovery script - Source: https://github.com/apache/spark/blob/v3.0.1/examples/src/main/scripts/getGpusResources.sh (Accessed: 2021-01-03)	60

Contents

1	Introduction	1
1.1	Distributed Computing	1
1.2	Computing Acceleration with GPUs	2
1.3	Auto-Scaling	2
1.4	Automated Deployment Pipeline	3
1.5	Research Objective and Research Questions	3
1.6	Problem Statement	4
1.7	Thesis Structure	5
1.8	Research Methodologies	5
2	Theoretical Foundation	7
2.1	Scalability	7
2.1.1	Horizontal Scaling	7
2.1.2	Limitations of Scalability	7
2.2	Deployment Pipeline	8
2.2.1	Continuous Integration	8
2.2.2	Continuous Delivery	8
2.3	Autonomic Computing	8
2.3.1	Autonomic Computing Concept	9
2.3.2	Managed Resources	10
2.3.3	Autonomic Manager	11
2.4	System Performance	12
2.4.1	Performance Metrics	12
2.4.2	Time-Based Utilization	12
2.5	Monitoring	12
2.5.1	Database	13
2.5.2	Push and Pull	13
2.5.3	Multi-Dimensional Metrics	13
2.5.4	Query Language	15
2.5.5	Choosing a Monitoring Tool	15
3	Related Work	17
3.1	Elastic Environments	17
3.1.1	Architecture	17
3.1.2	Auto-Scaler	17

3.1.3	Auto-Scaling Algorithms	18
3.2	Heterogenous GPU aware Spark systems	18
4	Technical Background	21
4.1	Docker	21
4.1.1	Docker architecture	21
4.1.2	Docker Image	22
4.1.3	Docker Container	22
4.1.4	Docker Swarm Mode	23
4.2	Apache Spark	24
4.2.1	Spark programming model	24
4.2.2	Application Architecture	26
4.2.3	Standalone Cluster Deployment	27
4.3	RAPIDS accelerator for Apache Spark	28
4.3.1	Extension of the Spark programming model	29
4.3.2	Installation Requirements for Apache Spark Standalone Mode	29
4.4	Prometheus	30
4.4.1	Prometheus architecture	30
4.5	cAdvisor	31
4.6	GitLab CI/CD	31
4.6.1	CI/CD Pipeline	32
4.6.2	Job Execution	32
4.7	Scaling Heat	33
4.7.1	Recurrence Factor	33
4.7.2	Scaling Heat Algorithm Concept	34
4.8	Kubernetes Horizontal Pod Autoscaler	35
5	Conceptual Design	37
5.1	Design Restrictions	37
5.2	CI/CD	37
5.3	Identification of suitable Metrics for Scaling	37
5.3.1	CPU Performance	38
5.3.2	GPU Performance	38
5.4	Computing environment Architecture	38
5.4.1	Overall	38
5.4.2	Apache Spark Cluster	39
5.4.3	Autonomic Manager	39
5.5	Auto-Scaler	41
5.5.1	Configuration	41
5.5.2	Analyze	42
5.5.3	Plan	43
5.5.4	Execute	43
6	Implementation	45
6.1	Background	45

6.2	Auto-Scaler	45
6.2.1	Configuration	45
6.2.2	Scaling the Apache Worker Service	47
6.2.3	Docker Image	48
6.3	Deployment of a Docker Swarm	48
6.3.1	Hardware	48
6.3.2	Software info	48
6.3.3	Swarm	48
6.3.4	Build Script	48
6.3.5	Apache Spark Cluster with GPU Acceleration	48
6.3.6	Autonomic Manager	50
6.4	Automatic Deployment of Apache Spark Applications	51
7	Evaluation	53
7.1	Test Environment	53
7.2	Algorithms and Datasets	53
7.2.1	K-MEANS	53
7.2.2	Naive Bayes Classifier	53
8	Outlook	55
8.1	Optimizing Scaling	55
8.2	Reinforcement Learning for Auto-Scaling	55
9	Conclusion	57
9.1	Cluster architecture	57
	Anhang	57
A	Apache Spark Cluster Implementation	59

Notation

Konventionen

x	Skalar
\underline{x}	Spaltenvektor
$\mathbf{x}, \underline{\mathbf{x}}$	Zufallsvariable/-vektor
$\hat{x}, \hat{\underline{x}}$	Mittelwert/-vektor
x^*, \underline{x}^*	Optimaler Wert/Vektor
$x_{0:k}, \underline{x}_{0:k}$	Folge von Werten (x_0, x_1, \dots, x_k) / Vektoren $(\underline{x}_0, \underline{x}_1, \dots, \underline{x}_k)$
\mathbf{A}	Matrizen
\mathcal{A}	Mengen
\preceq, \prec	schwache/strenge Präferenzrelation
\mathbb{R}	Reelle Zahlen
\mathbb{N}	Natürliche Zahlen
■	Ende eines Beispiels
□	Ende eines Beweises

Operatoren

\mathbf{A}^T	Matrixtransposition
\mathbf{A}^{-1}	Matrixinversion
$ \mathbf{A} $	Determinante einer Matrix
$ \mathcal{A} $	Kardinalität der Menge \mathcal{A}
$\text{pot}(\mathcal{A})$	Potenzmenge von \mathcal{A}
$\mathbb{E}\{\cdot\}$	Erwartungswertoperator
$\mathcal{O}(g)$	O-Kalkül entsprechend der Landau-Notation bei welcher beispielsweise $f(x) \in \mathcal{O}(g(x))$ besagt, dass die Funktion $f(x)$ die Komplexität $\mathcal{O}(g(x))$ besitzt

Spezielle Funktionen

$\Pr(\mathcal{E})$	Wahrscheinlichkeitsmaß, welches die Wahrscheinlichkeit angibt, dass Ereignis \mathcal{E} eintritt
$p(\underline{x})$	(Wahrscheinlichkeits-)Dichtefunktion für kontinuierliche \underline{x}

	und Zähldichte für diskrete \underline{x}
$p(\underline{x} y)$	Bedingte Dichtefunktion
$P(\underline{x})$	(Wahrscheinlichkeits-)Verteilungsfunktion
$\operatorname{erf}(x)$	Gauß'sche Fehlerfunktion
$\exp(x)$	Exponentialfunktion e^x
$\mathcal{N}(\underline{x}; \hat{\underline{x}}, \mathbf{C}_x)$	Gaußdichte, d. h. Dichtefunktion eines normalverteilten Zufallsvektors \underline{x} mit Mittelwertvektor $\hat{\underline{x}}$ und Kovarianzmatrix \mathbf{C}_x

Introduction

1.1 Distributed Computing

Machine Learning and Big Data projects consist of a combination of extract-transform-load (ETL) pipelines and compute intensive algorithms to create meaningful informations from massive datasets [Vad18]. Performing ETL applications on a single machine, limits the scale and performance of the application. Given this limitation, the idea of distributed computing evolved. In the concept of distributed computing, multiple machines with commodity hardware form a cluster to utilize their resources to solve high complex problems [GOKB16]. Several companies utilized the idea of distributed computing to solve some of their business problems. Google developed the MapReduce[DG04] framework. MapReduce gives the opportunity to solve massive complex problems in parallel on a cluster of single machines. Yahoo published an ecosystem platform for distributed computing called Apache Hadoop¹. Apache Hadoop contributed to create a cluster to process massive amounts of data in parallel. Implementing data pipelines with MapReduce requires to chain multiple MapReduce jobs together. This causes a huge amount of writing and reading operation to the disk with bad impact on the overall performance. Another framework called Apache Spark was developed to simplify writing and executing parallel applications at scale while keeping the benefits of MapReduce's scalability and fault-tolerant data processing. Apache Spark provides a performance improve of 10x in iterative Machine Learning algorithms over MapReduce [ZCF⁺10] and has evolved as a replacement for MapReduce as the distributed computing framework of choice.

¹ Apache Hadoop - <https://hadoop.apache.org/> (Accessed: 2020-01-08)

1.2 Computing Acceleration with GPUs

Distributed computing frameworks like Apache Spark perform applications on a huge amount of CPU cores to enable parallelization. A CPU is build of multiple cores which are optimized for sequential serial processing. Performing computationally intensive applications on an Apache Spark cluster, consumes a huge amount of CPU cycles with bad impact on the overall performance [PYY15]. To handle the complexity of Big Data applications, from executing Machine Learning algorithms or training Deep Learning models, a remaining option of distributed computing clusters is to scale-up individual nodes. Scaling-up is limited by resource capacity and can be become uneconomically at a specific point. To perform computationally complex applications with better performance, Graphical Process Units (GPU) have become first class citizens in modern data centers. The architecture of a GPU consists of a huge amount of smaller and more efficient cores which are optimized for parallel data processing (handling multiple tasks simultaneously). In general, GPUs process data at a much faster rate than CPUs are capable. Leveraging GPUs for distributed computing frameworks like Apache Spark, can boost the overall performance of performing complex algorithms on large datasets.

1.3 Auto-Scaling

Adjusting the resources in a computing environment is not an easy task. To do it manually, a system administrator needs a deep knowledge about the environment and has to watch performance spikes regularly. This is a resource wasting process. In an optimal way, an automated process would watch the computing environment, analyse performance metrics and automatically add or remove resources to optimize the performance and cost. This process is called auto-scaling.

Hiring experts to manually watching an application and scaling an computing environment is a waste of resources and cost. An Auto-Scaler takes care of watching the environment by adding and removing resources to adapt to the computing needs. The Auto-Scaler can be configured to take care of optimal resource allocation while keeping the cost of running low.

Two different scaling approaches can be used to scale resources in a computing environment: Vertical-scaling and horizontal-scaling. Vertical scaling refers to adjusting the hardware resources of an individual node in the environment. Hardware adjustments can include adding (scale-up) or removing (scale-down) resources like memory or CPU cores [Wil12]. By adding more powerful resources to a node, a node can take more throughput and perform more specialized tasks [LT15]. Adjusting the nodes in a computing environment is referred as horizontal scaling [Wil12]. Increasing the number of nodes in an environment, increases the overall computing capacity and additionally, the workload can be distributed across all nodes [Wil12, LT15]. Both scaling approaches are not exclusive. A computing environment can

be designed to scale vertically, horizontally or both [Wil12]. Vertical scaling is limited by the maximum hardware capacity. Furthermore, a point can be reached where more powerful hardware resources become unaffordable or are not available [LT11]. Therefore, horizontal scaling is the preferred approach to enable auto-scaling.

1.4 Automated Deployment Pipeline

Building, testing and releasing software manually is a time-consuming and error-prone process. To overcome this issue, a pattern called deployment pipeline automates the build, test, deploy and release processes of an application development cycle. The concept of deployment pipelines is based on automation scripts which will be performed on every changes on an applications source code, environment, data or configuration [FH10]. A fully automated deployment pipeline has many improvements over deploying applications manually:

- Makes every process until release visible to all developers [FH10]
- Errors can be identified and resolved at an early stage [FH10]
- The ability to deploy and release any version of an application to any environment [FH10]
- A non automated deployment process is not repeatable and reliable [FH10]
- The automation scripts can serve as documentation [FH10]
- If an application has been deployed manually, there is no guarantee that the documentation has been followed [FH10]

The automated deployment pipeline is based on the Continuous Integration (CI) process. Furthermore, the deployment pipeline is the logical implementation of CI [FH10]. Nowadays,

1.5 Research Objective and Research Questions

The thesis work will be implemented at Center for Cyber Cognitive Intelligence at the Fraunhofer IPA². At the IPA, developers deploy ML applications using Apache Spark. Executing these applications should be accelerated using GPUs on the Apache Spark cluster. If the point is reached where the performance utilization of the Apache Spark cluster is too high, the environment should adapt to the computing needs and scale the Apache

2 Fraunhofer Institute for Manufacturing Engineering and Automation IPA - <https://www.ipa.fraunhofer.de/> (Accessed: 2021-01-07)

Spark cluster by adding more workers. In addition, developers should have the ability to submit an application to the Apache Spark cluster automatically which will be triggered by pushing changes to the code base of the application.

To address the goal of this thesis, the following three research question will be investigated:

- **RQ1:** Is it possible to scale the number of Apache Spark Worker in accordance to performance utilization?
- **RQ2:** How can Apache Spark be extended to accelerate application execution with GPU support?
- **RQ3:** Is it possible to automate the deployment process of applications to a running Apache Spark cluster?

The first research question searches for concepts to create a self-adapting computing environment. To answer this question, state-of-the-art computing architectures have to be investigated. Monitoring tools to collect performance metrics need to be evaluated. Additionally, tools which enable fast deployment of computing units. Furthermore, a suitable scaling approach has to be investigated.

The main goal of the second research question is to enable Apache Spark to perform algorithms with GPU acceleration included. Therefore, a concept needs to be investigated to extend Apache Spark to use GPUs for suitable algorithms in addition to the available CPUs.

The last research question has a more applied nature. Automating the development cycle of an application is a well investigated topic. The IPA is using a platform called GitLab (will be introduced in Section 4.6) which provides an API to build automated pipelines. To answer this research question, GitLab's functionality will be investigated to find a solution that fits the need of this project work.

1.6 Problem Statement

Given the previously introduced research questions and the research objective, this thesis work will provide a solution to the following three problem statements:

1. Developers at the Fraunhofer IPA perform ML applications on an Apache Spark cluster. On this cluster, the workload of applications is only distributed over a set of CPUs to enable parallelism. To improve the execution time of these applications, the cluster needs to be aware of distributing suitable operations on GPUs as well.

2. To enable GPU acceleration for Apache Spark alone is not sufficient to increase the performance. At some point, an Apache Spark Worker can reach the limit of its available computing resources. If this point is reached, the environment should automatically scale the number of Apache Spark worker to distribute the workload.
3. To perform an Apache Spark application to the cluster, developers have to submit the application manually. With an automated deployment pipeline, developers can submit an application by pushing changes to the code base. Additionally, a deployment pipeline will contribute to the reliability of executing applications and reduces the development time.

1.7 Thesis Structure

Chapter 2 provides the theoretical foundation about concepts which have been introduced in this chapter. Chapter 3 focuses on related work which provides solutions to solve the given problems of this thesis introduced in Section 1.6. In Chapter 4 all technologies which are being used to implement the objective of this thesis are being introduced. Afterwards in Chapter 5, a conceptual design of a dynamic computing environment and an automated deployment pipeline is being described. Chapter 6 describes the implementation of the computing environment and how the deployment pipeline is being used to automate the deployment of applications to the computing environment. In Chapter 7 the results of the implementation are being presented and analysed. Chapter 8 introduces further work, which has been discovered during the work of this thesis, as well as improvements for the implementation. Finally Chapter 9 ...

1.8 Research Methodologies

To reach the research objective and solve all problems of this thesis work, the first step will be to read state-of-the-art literature about the following topics:

1. Enabling GPU acceleration on Apache Spark
2. Self-adapting heterogeneous systems
3. Automated deployment pipelines

All literature has to be summarized to gain a deep understanding of how to solve the defined problem statements. During this step, suitable software tools which help to implement the research objective have to be identified. Afterwards, the summarized knowledge will be used to create a conceptual design of the implementation. The conceptual design will be used to implement an environment that fulfils the research objective and solves each

defined problem statement. After the implementation has been completed, the environment needs to be evaluated to ensure it fulfils the thesis research objective. Lastly, the results of the evaluation will be analysed.

Chapter 2

Theoretical Foundation

This chapter provides background information about ...

**TODO: Describe
Chapter**

2.1 Scalability

Scalability defines the ability of a computing system to handle an increasing amount of load [Far17]. The main reason for scaling is High-Availability TOOLKIT 2.1. To choose the right approach of scaling an environment is essential. To scale the computing capacity of a system, different approaches exist. A major scaling approach is horizontal scaling.

2.1.1 Horizontal Scaling

Horizontal scaling is accomplished by duplicating nodes in the computing environment [Wil12]. Duplication of the nodes in a computing environment increases the computational capacity of the environment. In addition, the workload can be distributed across all clones to handle and balance and increasing load of tasks [Wil12, LT15]. To increase the efficiency of horizontal scaling, all nodes should be homogeneous. Homogeneous nodes are able to perform the same work and response as other nodes [LT15].

2.1.2 Limitations of Scalability

The limit of scalability is reached, when a computing system is not able to serve the requests of its concurrent users [Wil12]. If the scalability of a computing environment is reached, an option is to add more powerful hardware resources to the system. This approach is called vertical scaling [QUL15]. By adding more powerful hardware resources, the point can be reached, where more powerful hardware becomes unaffordable or not available [Wil12]. To overcome the limits of hardware capacity, a computing system should be designed to scale horizontally in the first place [LT15].

2.2 Deployment Pipeline

A deployment pipeline is a implementation of the process for getting software from source code to production. It is based on the concept of Continuous Integration (CI) and Continuous Delivery (CD). The process involves building, testing and deploying software through automated scripts [FH10].

2.2.1 Continuous Integration

CI is the first part in a deployment pipeline. The CI process is responsible for building, testing and validating the software. This ensures that the software is in a deployable state at all time [Ros17]. A CI pipeline is based on several requirements:

1. **A single code base:** It is important that the source code of an application is in a single code base. This includes all automation script to build, test and validate the application. Each change on the code base will trigger the CI pipeline on the build server [Ros17].
2. **A build server:** The build server is responsible to monitor the code base for changes. If a change has been pushed, the build server automatically executes the CI scripts in order [Ros17].

The goal of CI pipeline is, to have a production ready artefact of the software after performing [Ros17].

2.2.2 Continuous Delivery

After the CI pipeline has succeeded, a CD pipeline is responsible to deploy the artefact to a production server [Ros17].

2.3 Autonomic Computing

Autonomic computing is the ability of an IT infrastructure to automatically manage itself in accordance to high level objectives defined by administrators [KC03]. Autonomic computing gives an IT infrastructure the flexibility to adapt dynamic requirements quickly and effectively to meet the challenges of modern business needs [Mur04]. Therefore, autonomic computing environments can reduce operating costs, lower failure rates, make systems more secure and quickly respond to business needs [JSAP04].

Computing systems need to obtain a detailed knowledge of it's environment and how to extend it's resources to be truly autonomic [Mur04]. An autonomic computing system is defined by four elements:

- **Self-configuring:** Self-configuring refers to the ability of an IT environment to adapt dynamically to system changes and to be able to deploy new components automatically. Therefore, the system needs

to understand and control the characteristics of a configurable item [Mur04, Sin06].

- **Self-optimizing:** To ensure given goals and objectives, a self-optimizing environment has the ability to efficiently maximize resource allocation and utilization [JSAP04]. To accomplish this requirement, the environment has to monitor all resources to determine if an action is needed [Mur04].
- **Self-healing:** Self-healing environments are able to detect problematic operations and then perform policy-based actions to ensure that the systems health is stable [Sin06, JSAP04]. The policies of the actions have to be defined and should be executed without disrupting the system [Sin06, JSAP04].
- **Self-protecting:** The environment must identify unauthorized access and threats to the system and automatically protect itself taking appropriate actions during its runtime [Sin06, JSAP04].

2.3.1 Autonomic Computing Concept

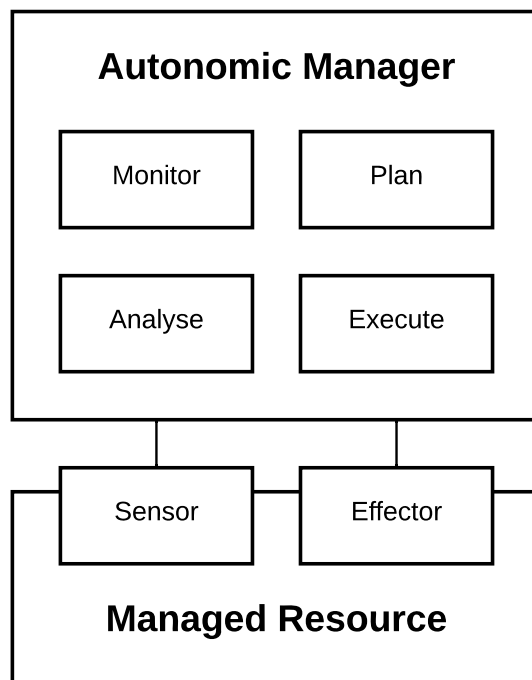


Figure 2.1: Autonomic computing concept - Source: Authors own model, based on [JSAP04].

Figure 2.1 demonstrates the main concept of an autonomic computing environment. The autonomic computing architecture relies on monitoring sensors and an adoption engine (autonomic manager) to manage resources

in the environment [GBR11]. In an autonomic computing environment, all components have to communicate to each other and can manage themselves. Appropriate decisions will be made by an autonomic manager that knows the given policies [JSAP04].

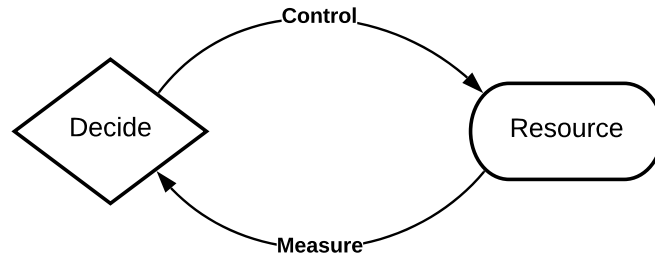


Figure 2.2: The control-loop concept - Source: Authors own model, based on [Mur04].

The core element of the autonomic architecture is the control-loop. Figure 2.2 illustrates the concept of a control-loop. The control-loop collects details about resources through monitoring and makes decisions based on analysis of the collected details to adjust the system if needed [Mur04].

2.3.2 Managed Resources

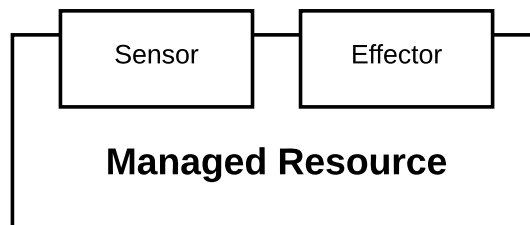


Figure 2.3: Managed resource - Source: Authors own model, based on [JSAP04].

A managed resource is a single component or a combination of components in the autonomic computing environment [Mur04, JSAP04]. A component can be a hardware or software component, e.g. a database, a server, an application or a different entity [Sin06]. They are controlled by their sensors and effectors, as illustrated in Figure 2.3. Sensors are used to collect information about the state of the resource and effectors can be used to change the state of the resource [JSAP04]. The combination of sensors and effectors is called a touchpoint, which provides an interface for communication with the autonomic manager [Sin06]. The ability to manage and control managed resources make them highly scalable [Mur04].

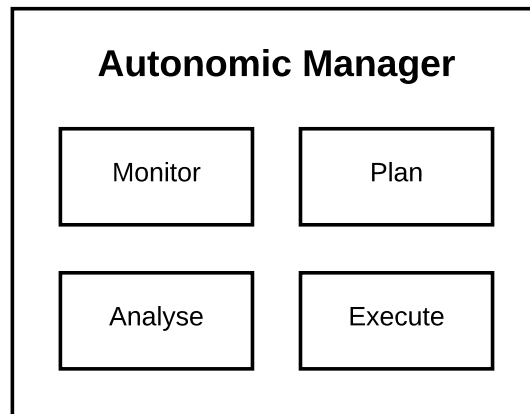


Figure 2.4: Autonomic manager - Source: Authors own model, based on [JSAP04].

2.3.3 Autonomic Manager

The autonomic manager implements the control-loop to collect, aggregate, filter and report system metrics from the managed resources. It can only make adjustments within its own scope and uses predefined policies to make decisions of what actions have to be executed to accommodate the goals and objectives [Mur04, Sin06]. In addition, the autonomic manager gains knowledge through analyzing the managed resources [Mur04]. The autonomic computing concept digests the MAPE-K model to implement an autonomic manager, as illustrated in Figure 2.4 [GBR11].

- **Monitor:** The monitor phase is responsible to collect the needed metrics from all managed resources and applies aggregation and filter operations to the collected data [Sin06].
- **Analyze:** The autonomic manager has to gain knowledge to determine if changes have to be made to the environment [Sin06]. To predict future situations, the autonomic manager can model complex situation given the collected knowledge [JSAP04].
- **Plan:** Plans have to be structured to achieve defined goals and objectives. A plan consists of policy-based actions [JSAP04, Sin06].
- **Execute:** The execute phase applies all necessary changes to the computing system [Sin06].

Multiple autonomic manager can exist in an autonomic computing environment to perform only certain parts. For example, there can be one autonomic manager which is responsible to monitor and analyse the system and another autonomic manager to plan and execute. To create a complete and closed control-loop, multiple autonomic manager can be composed together [Sin06].

2.4 System Performance

2.4.1 Performance Metrics

Performance metrics are statistics that describe the system performance. These statistics are generated by the system, applications or other tools [Gre20]. The following are examples of common types for performance metrics:

- **Throughput:** Volume of data or operations per second [Gre20].
- **Latency:** Time of operation [Gre20].
- **Utilization:** Usage of a component [Gre20].

Measuring performance metrics can cause an overhead. To gather and store performance metrics, CPU cycles must be spent. This can have a negative affect on the target performance [Gre20].

2.4.2 Time-Based Utilization

Utilization is a performance metrics that describes the usage of a device, e.g. CPU device usage. A time-based utilization describes the usage of a component during a time period where the component was actively performing work [Gre20]. When a resource approaches 100% utilization, the performance of that resource can degrade. If a component can process operations in parallel, the performance does not have to degrade much at 100% utilization and the component can process more work [Gre20].

2.5 Monitoring

Monitoring is a process, that aims to detect and take care of system faults. In a dynamic environment, becoming aware of the system is trivial [Lig12]. A monitoring system consists of a set of different tools. The tools are responsible to perform measurements on components in the computing environment and collect, store and interpret the monitored data [Lig12].

In the monitoring process, illustrated in Figure 2.5, data is continuously collected by agents. An agent is a process that continuously gathers metrics. Data can be device statistics, logs or system measurements. Agents will group these data into metrics and submit them to the monitoring system via a protocol. The monitoring system will store the metrics in its database [Lig12].

The requirements for a monitoring system, that is able to monitor a dynamic changing environment, are the following:

- An efficient database to store metrics
- A push or pull based way of gathering metrics [Far17]

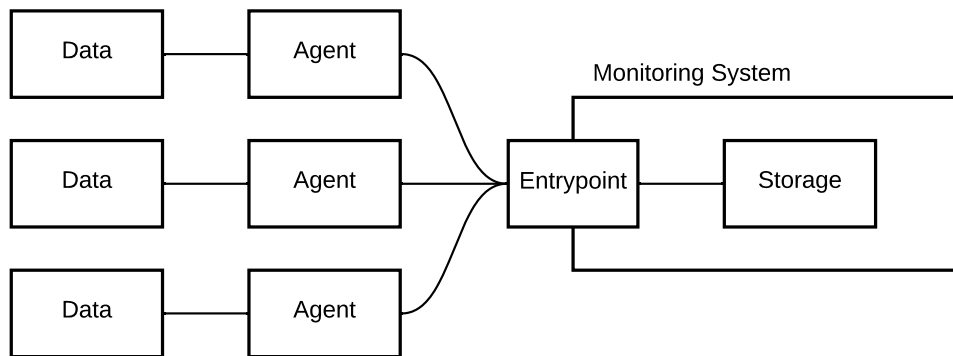


Figure 2.5: The monitoring process - Source: Authors own model.

- Multi-dimensional metrics [Far17]
- A powerful query language [Far17]

2.5.1 Database

Continuous data needs to be stored in the most efficient way. Time-series databases (TSDB) are optimized to store and retrieve time-series data. In a time-series database, metrics will be stored in a compact and optimized format. This allows the database to store a massive amount of time-series data on a single machine.

2.5.2 Push and Pull

The approach how the monitoring systems gathers metrics to store in the database plays a significant role. Push and pull based systems are the two primary approaches to gather metrics from services. Push based monitoring systems expect services to push metrics to their storage. On the other hand, pull based monitoring systems scrape metrics from all defined targets. Targets do not know about the existence of the monitoring system and only need to collect and expose metrics [Far17].

Service discovery is an important aspect to decide whenever to use a pull or push based monitoring system [Far17].

In a push-based environment, services only need to know the address of the monitoring service to push their data to the storage [Far17].

A pull-based monitoring tool needs to know the address of each target in the environment. The advantage of a pull-based monitoring systems is the simplicity to detect whenever a target has failed or is not available [Far17].

2.5.3 Multi-Dimensional Metrics

For query languages to be effective, metrics need to be dimensional. Metric without dimensions, are limited in their capabilities. In a dynamic environ-

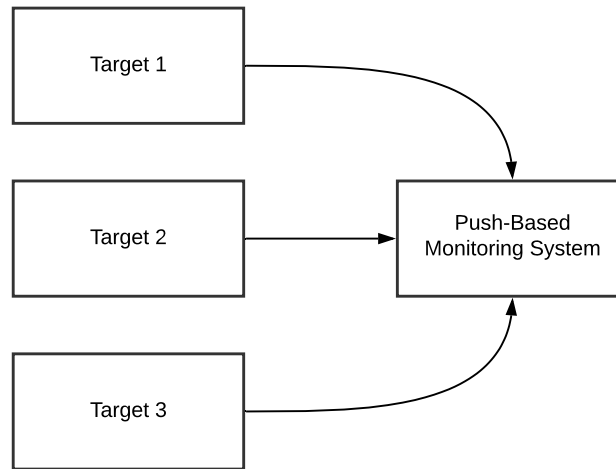


Figure 2.6: Push-based monitoring approach - Source: Authors own model.

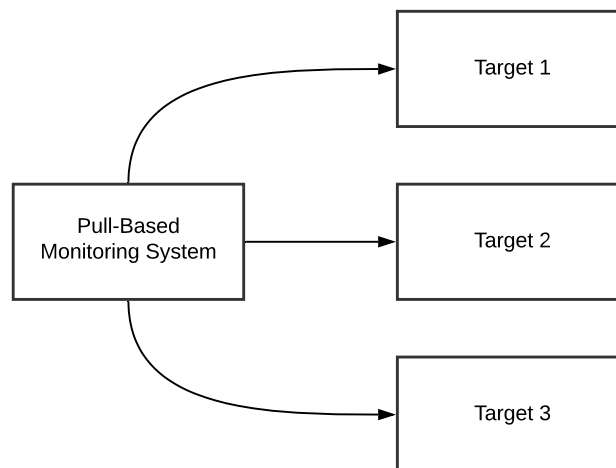


Figure 2.7: Pull-based monitoring approach - Source: Authors own model.

ment, services are dynamically added and removed. Therefore, a dynamic environment needs dynamic analytics where metrics represent all dimension in the environment [Far18].

```
1 container_cpu_user_seconds_total
```

Listing 2.1: Example of a dimensionless-metric

```
1 container_cpu_user_seconds_total{image="spark-worker:3.0.1-hadoop2.7"}
```

Listing 2.2: Example of a metric with dimensions

As the exasample Listing 2.1 and Listing 2.2 show, the metric with dimension provide more efficient querieing to gather informations about the environment.

2.5.4 Query Language

Another requirement is a powerful and flexible query language. It should have the ability to query a multi-dimensional data model. Furthermore, the query language must allow aggregations on time-series data.

Chapter 3

Related Work

This chapter provides background information about ...

**TODO: Describe
Chapter**

3.1 Elastic Environments

In recent years, container technologies have been used efficiently in complex IT environments. Dynamic scaling of containerized applications is an active area of research. The studied research can be divided in two parts.

3.1.1 Architecture

In the work by Lorido-Botrán et al. they reviewed state-of-the-art literatures about auto-scaling and proposed a process for auto-scaling homogeneous elastic applications. They mentioned three different problems, auto-scaler face while remaining the Quality of Service (QoS): Under-provisioning, over-provisioning and oscillation. Under-provisioning refers to, if not enough resources are available, over-provisioning means that more resources are available than needed and oscillation occurs when the environment gets scaled too quickly before the impact is clear. They mentioned the MAPE-Loop which consists of four different parts: Monitor, Analyze, Plan and Execute. The Auto-Scaler is part

3.1.2 Auto-Scaler

Srirama et al. [SAP20] designed a heuristic-based auto-scaling strategy for container-based microservices in a cloud environment. The purpose of the auto-scaling strategy was to balance the overall resource utilization across microservices in the environment. The proposed auto-scaling strategy performed better results than state-of-the-art algorithms in processing time, processing cost and resource utilization. The processing cost of microservices could be reduced by 12-20% and the CPU and memory utilization of cloud-servers have been maximized by 9-15% and 10-18%.

Lorido-Botrán et al. [LBMAL13] compared different representative auto-scaling techniques in a simulation in terms of cost and SLO violations. They compared load balancing with static threshold-based rules, reactive and proactive techniques based on CPU load. Load balancing is based on static rules defining the upper and lower thresholds of a specific load. For example *if CPU > 80% then scale-out; if CPU < 20% then scale-in*. The difficulty of this technique is to set the ideal rules. False rules can lead to bad performance. Proactive techniques try to predict the future values of performance metrics based on historical data. Reactive techniques are based on control theory to automate the systems management. In Addition, the authors proposed a new auto-scaling technique. To overcome the difficulties of static thresholds, the authors proposed a new auto-scaling technique using rules with dynamic thresholds. The results showed, that for auto-scaling techniques to scale well, it highly depends on parameter tuning. The best result was achieved with proactive results with a minimum threshold of 20% and a maximum threshold of 60%.

3.1.3 Auto-Scaling Algorithms

Barna et al. [BKFL17] proposed an autonomic scaling architecture approach for containerized microservices. Their approach focused on creating an autonomic management system, following the autonomic computing concept [KC03], using a self-tuning performance model. The demonstrated architecture frequently monitors the environment and gathers performance metrics from components. It has the ability to analyze the data and dynamically scale components. In addition, to determine if a scaling action is needed, they proposed the *Scaling Heat Algorithm*. The Scaling Heat Algorithm is used to prevent unnecessary scaling actions, which can throw the environment temporarily off.

Casalicchio et al. [CP17] focused on the difference of absolute and relative metrics for container-based auto-scaling algorithms. They analysed the mechanism of the *Kubernetes Horizontal Pod Auto-scaling* (KHPA) algorithm and proposed a new auto-scaling algorithm based on KHPA using absolute metrics called *KHPA-A*. The results showed, that KHPA-A can reduce response time between 0.5x and 0.66x compared to KHPA. In addition, their work proposed an architecture using cAdvisor for collecting container performance metrics, Prometheus for monitoring, alerting and storing time-series data and Grafana for visualizing metrics.

3.2 Heterogenous GPU aware Spark systems

Apache Spark is a computing framework that distributes tasks between CPU cores. Data and compute intensive applications profit from GPU acceleration. Therefore, various research projects took effort to bring GPU acceleration to Apache Spark.

Li et al. [PYYN15] developed a middleware framework called *HeteroSpark* to enable GPU acceleration on Apache Spark worker nodes. HeteroSpark listens for function calls in Spark applications and invokes the GPU kernel for acceleration. For communication between CPU and GPU, HeteroSpark uses the Java RMI¹ API to send data from the CPU JVM to the GPU JVM for execution. The design provides a plug-n-play approach and an API for the user to call functions with GPU support. Overall, HeteroSpark is able to achieve a 18x speed-up for various Machine Learning applications running on Apache Spark.

Klodjan et al. [HBK18] introduced HetSpark a heterogeneous modification of Apache Spark. HetSpark extends Apache Spark with two executors, a GPU accelerated executor and a commodity class. The GPU accelerated executor is based on VineTalk[MPK⁺17] for GPU acceleration. The authors observed, that for compute intensive tasks GPU accelerated executors are preferable while for linear tasks CPU-only accelerators should be used.

Yuan et al. [YSH⁺16] proposed SparkGPU to enable parallel processing with GPUs in Apache Spark and contributes to achieve high performance and high throughput in Apache Spark applications. SparkGPU extends Apache Sparks to determine the suitability of parallel-processing for a task to enable task scheduling between CPU and GPU. SparkGPU accomplished to improve the performance of machine learning algorithms up to 16.13x and SQL query execution performance up to 4.83x.

¹ Java Remote Method Invocation

Chapter 4

Technical Background

In this chapter bla bla

TODO: Describe Chapter

4.1 Docker

Docker is an open-source platform that enables containerization of applications. Containerization is a technology to package, ship and run applications and their environment in individual containers. Docker is not a container technology itself, it hides the complexity of working with container technologies directly and instead provides an abstraction and the tools to work with containers [NK19, BMDM20, PNKM20].

TODO: Mehr zum Thema DevOps

4.1.1 Docker architecture

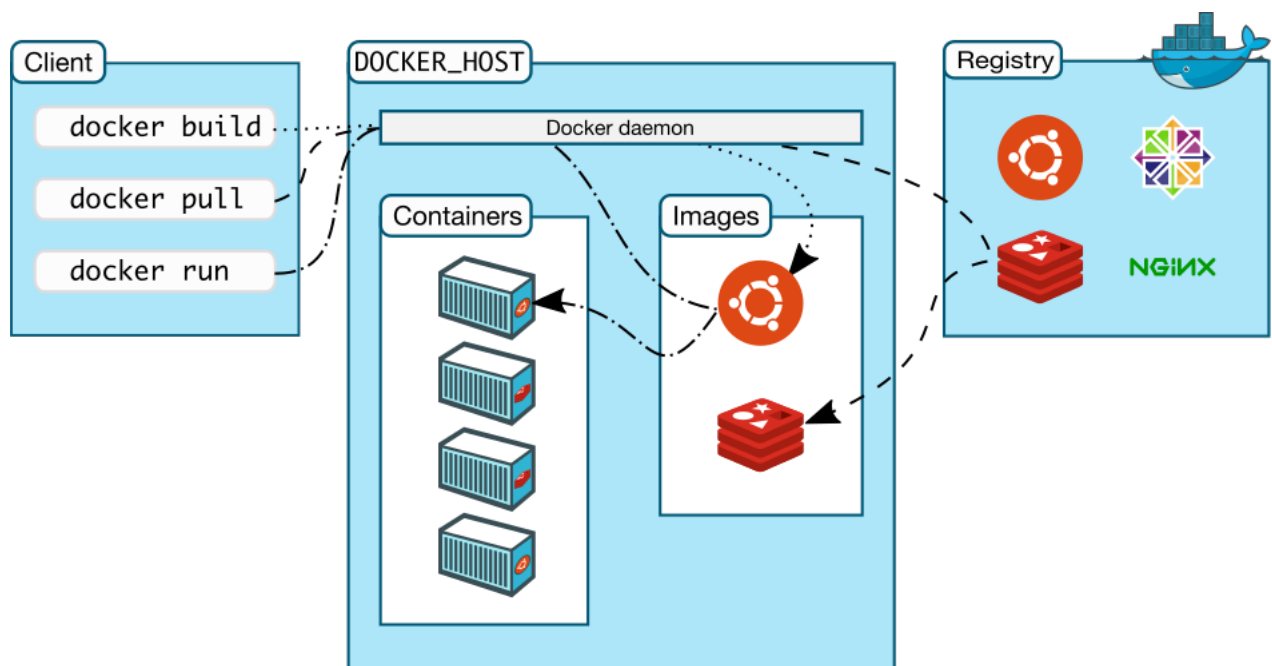


Figure 4.1: Docker architecture - Source: Authors own model, based on [Doc].

Figure 4.1 illustrates the client server architecture of Docker which consists of of a Docker client, the Docker daemon and a registry.

Docker client: The Docker client is an interface for the user to send commands to the Docker daemon [Doc].

Docker daemon: The Docker daemon manages all containers running on the host system and handles the containers resources, networks and volumes [BMDM20].

Docker Registry: A Docker registry stores images. Images can be pushed to a public or private registry and pulled from it to build a container [Doc].

4.1.2 Docker Image

An Image is a snapshot of the environment that is needed to run an application in a Docker container. The environment consists of all files, libraries and configurations that are needed for the application to run properly [Doc, NK19]. Images can be created from existing containers or from executing a build script called Dockerfile. A Dockerfile is a text file consisting of instructions for building an image. The Docker image builder executes the instructions of a Dockerfile from top to bottom [NK19].

4.1.3 Docker Container

A container is an execution environment running on the host-system kernel.

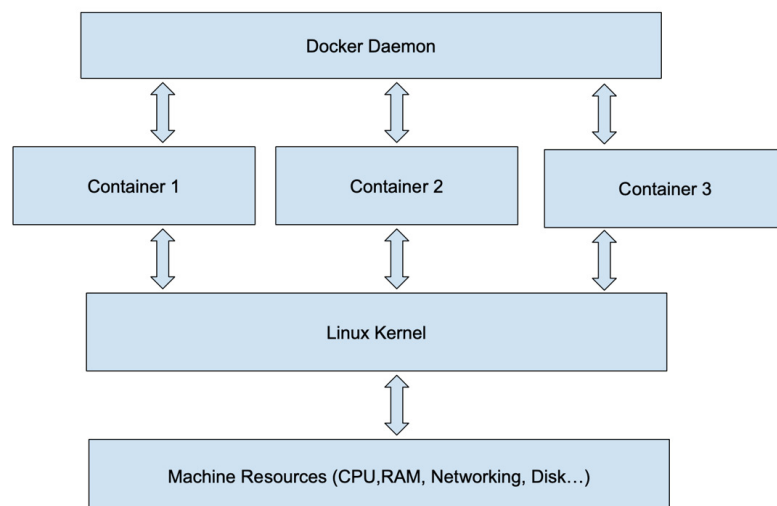


Figure 4.2: Docker basic container structure - Source: Authors own model, based on [BMDM20].

The advantage of a container is its lightweight nature. As illustrated in Figure 4.2, containers take advantage of OS-level virtualization instead of hardware-virtualization without the need of a hypervisor [Doc, NK19].

Containers share the resources of the host-system instead of using reserved resources [BMDM20]. Multiple containers can run on the host-system kernel and are by default isolated from each other [Doc]. In Docker, a container is a runnable unit of an image and is used for distributing and testing applications. A container can be configured to expose certain resources to the host system, e.g. network ports [BMDM20].

4.1.4 Docker Swarm Mode

Docker Swarm mode is the native cluster orchestration and management tool embedded in the Docker engine. In Docker Swarm mode, a cluster of multiple nodes is called a swarm. All nodes run in Swarm mode and act as managers or workers. In a swarm, multiple services can be deployed. The manager node is responsible to maintain the desired state of a service [Doc].

Docker Swarm mode will be used to create a self-healing and a self-adapting environment as described in Section ???. Many properties of Docker Swarm mode contribute the fact that it is an ideal candidate to create self-healing and self-adapting environment:

- **Desired state:** The manager node monitors the state of each service in the swarm and adapts the environment to maintain the desired state [Doc].
- **Cluster management and orchestration:** Docker Swarm mode is integrated with the Docker engine. A swarm can be created and managed using the Docker CLI [Doc].
- **Service model:** The Docker engine allows to define the desired state of a service. The manager node maintains the desired state of all services in the swarm [Doc].
- **Scaling:** The number of replicas can be defined for each service. The manager node will automatically adapt the number of replicas for a service to keep the desired state [Doc].
- **Multi-host networking:** A swarm runs all services in an overlay network. New services will automatically be added to the overlay network [Doc].

Nodes

A Docker engine participating in the swarm is called a node. Nodes can act as manager nodes, worker nodes or both [Doc].

The manager node is responsible for cluster orchestration and management. It maintains the desired state of all services and tasks in the swarm. In addition, the manager node dispatches tasks to worker nodes when service definitions will be submitted to the manager node [Doc].

Worker nodes are responsible to execute the tasks received by the manager node. While performing the tasks, the worker node notifies the manager node about the tasks state [Doc].

Services and Tasks

A service defines the desired state of a task. The state is defined by the number of replicas of a service and the configuration for the Docker container, e.g. Docker Image, resources, network, and more [Doc].

A task is a running Docker container. The task is defined by the corresponding service and will be managed by the manager node. A task can be performed on worker and manager nodes [Doc].

4.2 Apache Spark

Apache Spark is an open-source computing framework for parallel data processing on a large computer cluster. Spark manages the available resources and distributes computation tasks across a cluster to perform big-data processing operations at large scale [CZ18]. Before Spark was developed, Hadoop MapReduce [DG10] was the framework of choice for parallel operations on a computer cluster [ZCF⁺10]. Spark accomplished to outperform Hadoop by 10x for iterative Machine Learning [ZCF⁺10]. It is implemented in Scala¹, a JVM-based language and provides a programming interface for Scala, Java², Python³ and R⁴. In addition, Spark includes an interactive SQL shell and libraries to implement Machine Learning and streaming applications [CZ18]. It was developed in 2009 as the Spark research project at UC Berkeley and became an Apache Software Foundation project in 2013 [CZ18].

4.2.1 Spark programming model

Spark provides resilient distributed datasets (RDDs) as the main abstraction for parallel operations [ZCF⁺10]. Core types of Spark's higher-level structured API are built on top of RDDs [CZ18] and will automatically be optimized by Spark's Catalyst optimizer to run operations quick and efficient [Luu18].

Resilient distributed datasets

Resilient distributed datasets are fault-tolerant, parallel data structures to enable data sharing across cluster applications [ZCD⁺12]. They allow to express different cluster programming models like MapReduce, SQL and batched stream processing [ZCD⁺12]. RDDs have been implemented in Spark and serve as the underlying data structure for higher level APIs (Spark structured API) [ZCD⁺12]. RDD's are an immutable, partitioned collection of records and can only be initiated through transformations (e.g. map, filter) on data or other RDD's. An advantage of RDDs is, that they can

TODO:
Master-Slave
Architektur + Bild

¹ Scala programming language. <https://www.scala-lang.org/>

² Java programming language. <https://www.oracle.com/java/>

³ Python programming language. <https://www.python.org/>

⁴ R programming language. <https://www.r-project.org/>

be recovered through lineage. Lost partitions of an RDD can be recomputed from other RDDs in parallel on different nodes [ZCD⁺12]. RDDs are lower level APIs and should only be used in applications if custom data partitioning is needed [CZ18]. It is recommended to use Sparks structured API objects instead. Optimizations for RDDs have to be implemented manually while Spark automatically optimize the execution for structured API operations [CZ18].

Spark structured API

Spark provides high level structured APIs for manipulating all kinds of data. The three distributed core types are Datasets, DataFrames and SQL Tables and Views [CZ18]. Datasets and DataFrames are immutable, lazy evaluated collections that provide execution plans for operations [CZ18]. SQL Tables and Views work the same way as DataFrames, except that SQL is used as the interface instead of using the DataFrame programming interface [CZ18]. Datasets use JVM types and are therefore only available for JVM based languages. DataFrames are Datasets of type Row, which is the Spark internal optimized format for computations. This has advantages over JVM types which comes with garbage collection and object instantiation [CZ18].

Spark Catalyst

Spark also provides a query optimizer engine called Spark Catalyst. Figure 4.3 illustrates how the Spark Catalyst optimizer automatically optimizes Spark applications to run quickly and efficient. Before executing the user's code, the Catalyst optimizer translates the data-processing logic into a logical plan and optimizes the plan using heuristics [Luu18]. After that, the Catalyst optimizer converts the logical plan into a physical plan to create code that can be executed [Luu18].

Logical plans get created from a DataFrame or a SQL query. A logical plan represents the data-processing logic as a tree of operators and expressions where the Catalyst optimizer can apply sets of rule-based and cost-based optimizations [Luu18]. For example, the Catalyst can position a filter transformation in front of a join operation [Luu18].

From the logical plan, the Catalyst optimizer creates one ore more physical plans which consist of RDD operations [CZ18]. The cheapest physical will be generated into Java bytecode for execution across the cluster [Luu18].

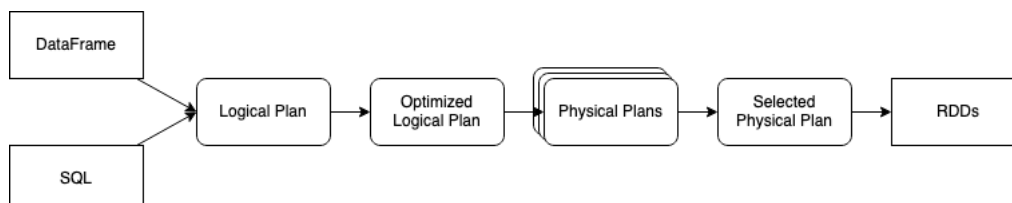


Figure 4.3: Optimization process of the Spark Catalyst - Source: Authors own model, based on [Luu18].

4.2.2 Application Architecture

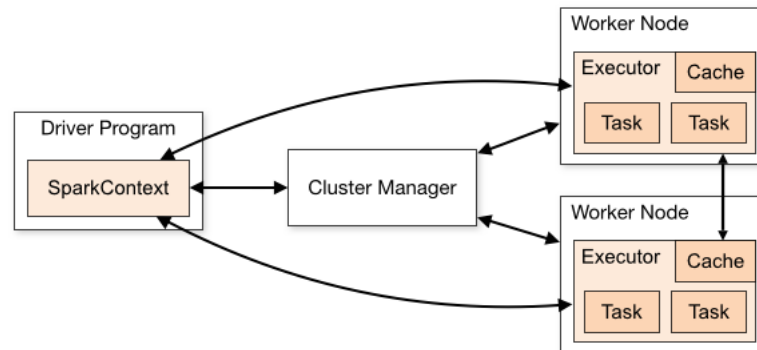


Figure 4.4: Overview of a Spark cluster architecture - Source: Authors own model, based on [Theb].

Figure 4.4 illustrates the main architecture of a Spark cluster. The architecture follows the master-worker model where the Spark driver is the master and the Spark executors are the worker [Luu18].

TODO: Nicht ganz richtig, Master und Worker sind machines und driver und executor sind prozesse

Spark Driver

The Spark driver is a JVM process on a physical machine and responsible to maintain the execution of a Spark application [CZ18]. It coordinates the application tasks onto each available executor [Luu18]. To get launch executors and get physical resources, the Spark driver interacts with the cluster manager [CZ18, Luu18].

Spark Executor

A Spark executor performs the tasks given by the Spark driver [CZ18]. It runs as a JVM process and runs until the Spark application finishes [Luu18]. After the executor finishes, it reports back to the Spark driver [CZ18]. Each task will be performed on a separate CPU core to enable parallel processing [Luu18].

Cluster Manager

The cluster manager is an external service that orchestrates the work between the machines in the cluster [Luu18, Theb]. The cluster manager knows about the resources of each worker and decides on which machine the Spark driver process and the executor processes run [Luu18, CZ18]. Spark supports different services that can run as the cluster manager: Standalone mode (introduced in Section 4.2.3), Apache Mesos⁵, Hadoop YARN[VMD⁺13] and Kubernetes⁶ [Theb]. The cluster manager provides three different deploy modes for acquiring resources in the cluster.

TODO: Explain standalone

⁵ Apache Mesos - <https://mesos.apache.org/> (Accessed: 2021-01-02)

⁶ Kubernetes - <https://kubernetes.io/> (Accessed: 2021-01-02)

- Cluster mode
- Client mode
- Local mode

To run an application in cluster mode, the user has to submit a precompiled JAR, python script or R script to the cluster manager [CZ18]. After that, the cluster manager starts the driver process and executor processes exclusively for the Spark application on machines inside the cluster [CZ18, Luu18].

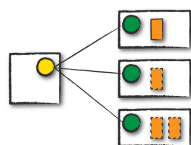


Figure 4.5: Spark's cluster mode - Source: Authors own model, based on [CZ18].

The difference between the client mode and the cluster mode is that, the driver process runs on the client machine outside of the Spark cluster [CZ18].

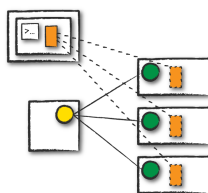


Figure 4.6: Spark's client mode - Source: Authors own model, based on [CZ18].

The local mode starts a Spark application on a single computer [CZ18]. It is not recommended to use the local mode in production, instead it should be used for testing Spark applications or learning the Spark framework [CZ18].

4.2.3 Standalone Cluster Deployment

The standalone mode is a basic cluster-manager build specifically for Apache Spark. It is developed to only run Apache Spark but supports workloads at large scale [CZ18].

Starting Master and Worker Nodes

Apache Spark provides executable launch scripts to start master and worker nodes in standalone mode. The executables can be found at *sbin/start-master.sh* to start a master node and at *sbin/start-slave.sh* to start a worker node. The worker launch executable requires the master node URI as parameter [Theb].

```
1 $ ./sbin/start-master.sh
```

Listing 4.1: Usage of master launch script

```
1 $ ./sbin/start-slave.sh spark://spark-master:7077
```

Listing 4.2: Usage of worker launch script

Listing 4.1 and Listing 4.2 provide an example of how to use both executables to start a master and a worker node. The URI *spark://spark-master:7077* in Listing 4.2 is an example of a master node URI. The master node launch script will print out the master URI after being executed successfully [Theb].

Resource Allocation

In standalone mode, worker need a set of resources configured. Therefore, a worker can assign resources to executors. To specify how a worker discovers resources, a discovery script has to be available [Theb].

Submitting Applications with spark-submit

To submit an application to a standalone cluster, Apache Spark provides the *spark-submit* executable. The executable file is available at *bin/spark-submit* in the installation folder of Apache Spark. In cluster mode (introduced in Section 4.2.2) the driver of a Spark application (see Section 4.2.2) will be launched from one of the worker processes inside the cluster. The submit process will finish after it has submitted the application. It does not wait for the submitted application to finish [Theb].

```
1 $ bin/spark-submit --master spark://spark-master:7077  
    application.py
```

Listing 4.3: Example usage of the spark-submit executable

Listing 4.3 shows how the *spark-submit* executable can be used to submit a Python application to a standalone Apache Spark cluster. *Spark-submit* requires the master node URI and the path to the desired Spark application file. With the *spark-submit* executable it is possible to submit Python, Java and R applications [Theb].

4.3 RAPIDS accelerator for Apache Spark

RAPIDS accelerator for Apache Spark is a plugin suite that aims to accelerate computing operation for Apache Spark using GPUs. It is available for Apache Spark 3 [NVI]. The plugin uses the RAPIDS⁷ libraries to extend the Apache Spark programming model, introduced in Section 4.2.1 [NVI, McD20].

⁷ Open GPU Data Science - <https://rapids.ai/> (Accessed: 01-01-2021)

4.3.1 Extension of the Spark programming model

The plugin suite extends the Spark programming model with a new DataFrame based on Apache Arrow⁸ data structures. In addition, the Catalyst optimizer (described in Section 4.2.1) is extended to generate GPU-aware query plans [McD20]. Apache arrow is a data platform to build high performance applications that work with large dataset's and to improve analytic algorithms. A component of Apache Arrow is the Arrow Columnar Format, an in-memory data structure specification for efficient analytic operations on GPUs and CPUs [Thea].

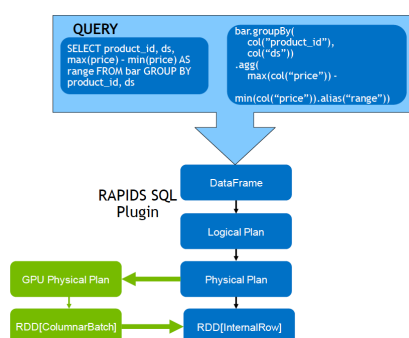


Figure 4.7: Catalyst optimization with RAPIDS accelerator for Apache Spark - Source: Authors own model, based on [McD20].

Figure 4.7 illustrates how the RAPIDS plugin suite extends the Catalyst optimization process illustrated previously in Figure 4.3. The Spark Catalyst optimizer identifies operator in a query plan that are supported by the RAPIDS APIs. To execute the query plan, these operators can be scheduled on a GPU within the Spark cluster [McD20]. If operators are not supported by the RAPIDS APIs, a physical plan for CPUs will be generated by the Catalyst optimizer to execute RDD operations [McD20].

4.3.2 Installation Requirements for Apache Spark Standalone Mode

The RAPIDS accelerator for Apache Spark is available for standalone Apache Spark cluster. To operate effectively, the following requirements need to be installed on all Apache Spark nodes in the cluster [NVI]:

- Java Runtime Environment (JRE)
- NVIDIA GPU driver
- CUDA Toolkit⁹

⁸ Arrow. A cross-language development platform for in-memory data - <https://arrow.apache.org/> (Accessed: 2020-12-03)

⁹ CUDA Toolkit - <https://developer.nvidia.com/cuda-toolkit> (Accessed: 2021-01-01)

- RAPIDS accelerator for Apache Spark Java library
- cudf Java library which is supported by the RAPIDS accelerator Java library and the installed CUDA toolkit
- GPU resource discovery script

4.4 Prometheus

Prometheus is an open-source monitoring and alerting system [Thec]. To collect and store data, Prometheus supports a multi-dimensional key-value pair based data model, according to Section 2.5.3, which can be analyzed in real-time using the PromQL query language [SP20]. Prometheus follows the pull-based approach, as described in detail in Section 2.5.2, to scrape metrics from hosts and services [BP19].

4.4.1 Prometheus architecture

TODO: Bilder für
components

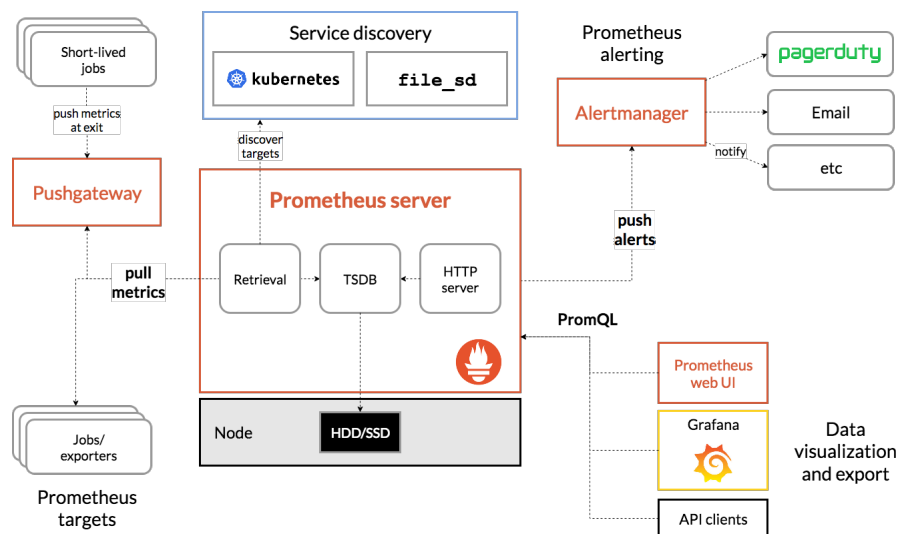


Figure 4.8: Prometheus high-level architecture - Source: Authors own model, based on [Thec, BP19].

Figure 4.8 illustrates the high-level architecture of Prometheus. The Prometheus ecosystem provides multiple components. Components can be optional, depending on the monitoring needs of the environment [BP19]. The main components of a Prometheus system are Prometheus server, Alertmanager, service discovery, exporters, Pushgateway and visualization tools [Thec].

Prometheus server: The Prometheus server is the main component of a Prometheus system. It is responsible to collect metrics as time-series data from targets and stores the collected data in the built-in TSDB [BP19]. Prometheus uses the concept of scraping to collect metrics from a target.

A target host has to expose an endpoint to make metrics available in the Prometheus data format [SP20]. Additionally, the Prometheus server triggers alerts to the Alertmanager if a configured condition becomes true [Thec].

Alertmanager: If an alerting rule becomes true, the Prometheus server generates an alert and pushes it to the Alertmanager. The Alertmanager generates notifications from the received alerts. A notification can take multiple forms like emails or chat messages. Webhooks can be implemented to trigger custom notifications [BP19].

Service discovery: As mentioned before, Prometheus follows a pull-based approach to fetch metrics from a target. To know about all targets, Prometheus needs a list of the corresponding hosts. The service discovery manages the complexity of maintaining a list of hosts manually in an changing infrastructure [BP19].

Exporters: If an application does not support an endpoint for Prometheus, an exporter can be used to fetch metrics and make them available to the Prometheus server. An exporter is a monitoring agent running on a target host that fetches metric from the host and exports them to the Prometheus server [SP20].

Pushgateway: If a target is not designed to be scraped, metrics can be pushed against the Pushgateway[Thec]. The Pushgateway converts the data into the Prometheus data format and passes them to the Prometheus server [SP20].

Visualization: Prometheus supports various tools for virtualization of the scraped data. Grafana¹⁰ is one of the widely used tools for this occasion.

4.5 cAdvisor

Container Advisor (cAdvisor) is a running daemon that collects, aggregates, analyses and exposes performance metrics from running containers. It has native support for Docker container and is deployed as a Docker container. cAdvisor collects metrics from the container daemon and Linux cgroups. Collected metrics will be exposed in the Prometheus file format [BP19, Goo].

4.6 GitLab CI/CD

GitLab CI/CD is a tool integrated into the GitLab platform that enables Continuous Integration (CI), Continuous Delivery (CD) and Continuous

¹⁰ Grafana: The open observability platform - <https://grafana.com/>

Deployment (CD) for software development. The GitLab platform integrates many development features like Git repository management and CI/CD. By pushing code changes to the codebase, GitLab CI/CD executes a pipeline of scripts to automate CI and CD processes of the software development cycle. A CI pipeline will consist of scripts that builds, tests and validate the updated codebase. A CD pipeline is responsible to deploy the application for production after the CI pipeline has executed successfully. Adding CI/CD pipelines to the software development cycle of an application, allows to catch bugs and errors early. This ensures that an application deployed to production will conform to established standards [Git].

4.6.1 CI/CD Pipeline

The fundamental component of GitLab CI/CD is called a pipeline. Pipelines will perform based on conditions. A conditions might be a push to the main branch of the repository [Git]. A pipeline comprises two components:

- **Stages:** A stage consists of one or multiple jobs that run in parallel. Furthermore, a stage defines how jobs will be executed. For example, a build stage only performs after a test stage has performed successfully [Git].
- **Jobs:** Jobs are responsible to perform the scripts defined by administrators. The scripts define necessary actions. For example compiling the source code or performing tests [Git].

Configuration

GitLab CI/CD is configured by a `.gitlab-ci.yml` file. It is necessary that this file is located in the repositories root directory. The configuration file will create a pipeline that performs after a push to the codebase [Git].

Basic Architecture

4.6.2 Job Execution

Jobs which are defined in the configuration file will be performed by GitLab runners. A GitLab runner is an agent that performs the jobs in its own environment and responds the result back to the GitLab instance. A runner is a lightweight and highly scalable application that runs on a server and performs one or multiple executors. An executor provides the environment where jobs will be executed in. GitLab runner provides multiple variants of executors. For example the Docker executor that connects to the underlying Docker engine. In addition, the Docker executor performs a job in a separate and isolated Docker container. GitLab runner can be set up only for specific projects or be available to all project on the GitLab platform [Git].

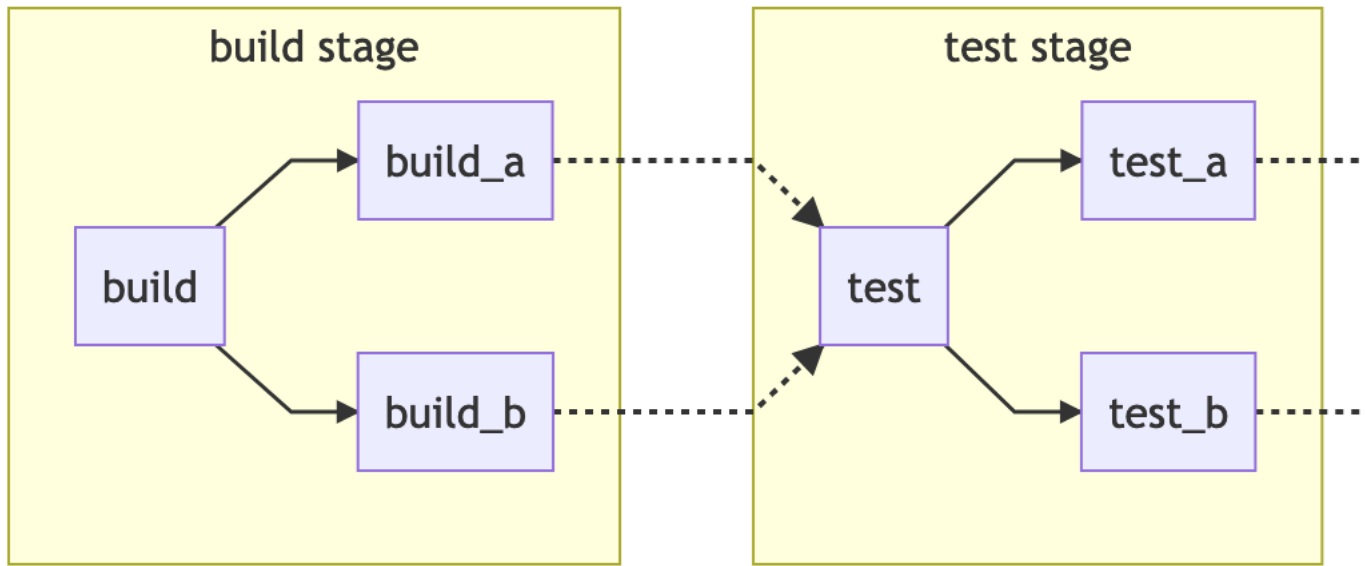


Figure 4.9: Basic architecture of a GitLab CI/CD pipeline - Source: Authors own model, based on [Git].

4.7 Scaling Heat

The Scaling Heat algorithm is a decision making algorithm to determine if a scaling action is necessary. It was introduced by Barna et al. [BKFL17] to overcome issues of traditional recurrence factor algorithms [BKFL17].

4.7.1 Recurrence Factor

In an Autonomic Computing environment, a scaling decision is made in each interval after data has been retrieved from the monitoring system (see Section 2.3 for the Autonomic Computing architecture). Sudden performance spikes can occur and can cause the decision algorithm to perform unnecessary scaling actions. These unnecessary scaling actions can have a negative impact on the overall computing performance. To overcome this issue, a recurrence factor needs to be introduced to the decision making algorithm. With a recurrence factor (n), a scaling action will only be performed until a performance threshold has been violated n times [BKFL17].

Traditional recurrence factor algorithm require violations to occur regularly. If a performance violation of the opposite direction occurs, the algorithm can get stuck in the violation process. Therefore, no scaling actions will be performed [BKFL17].

4.7.2 Scaling Heat Algorithm Concept

Algorithm 1: Scaling Heat decision making algorithm[BKFL17]

Input: *utilization* - The retrieved utilization for a performance metric

Input: *lower_threshold* and *upper_threshold* - Range limit of the performance metric

Input: *heat* - Current heat value of a performance metric. Indicating if a scaling action is necessary

Output: *heat* - New heat value for the next iteration

```

1 if utilization  $\geq$  upper_threshold then
    // cluster overload
2     if heat  $<$  0 then
        // reset heat for removal
3         heat  $\leftarrow$  0;
4         heat  $\leftarrow$  heat + 1;
5 else if utilization  $\leq$  lower_threshold then
    // cluster overload
6     if heat  $>$  0 then
        // reset heat for adding
7         heat  $\leftarrow$  0;
8         heat  $\leftarrow$  heat - 1;
9 else
    // utilization is within threshold range
    // move heat towards 0
10    if heat  $>$  0 then
11        heat  $\leftarrow$  heat - 1;
12    else if heat  $<$  0 then
13        heat  $\leftarrow$  heat + 1;
14 end
15 if heat = n then
16     Perform a scale-out action;
17     heat  $\leftarrow$  0;
18 else if heat =  $-n$  then
19     Perform a scale-in action;
20     heat  $\leftarrow$  0;
21 return heat

```

Algorithm 1 introduces the Scaling Heat algorithm. The algorithm is based in a concept called *heat*. The value of heat indicates if a scaling action of removing or adding components is necessary. If the given utilization of a performance metric violates the upper threshold, the heat value will increase. Violations of the lower threshold will cause a decrease respectively. When the heat reaches the recurrence factor n , positive for adding and negative for removing nodes, a scaling action will be executed. After executing a scaling

action, the heat value will be set to 0 [BKFL17].

4.8 Kubernetes Horizontal Pod Autoscaler

Kubernetes Horizontal Pod Autoscaler (KHPA) is a auto-scaling algorithm used in Kubernetes. KHPA scale the number of replicas per Pod based on the utilization of performance metrics. Kubernetes is an orchestration tool that allows to create and deploy units called Pods. A Pod is a running process on a cluster that encapsulates an application. The algorithm is based on a control loop. Each n seconds, the algorithm gathers performance metrics and computes the number of Pods to achieve the desired utilization of a performance metric.

Chapter 5

Conceptual Design

5.1 Design Restrictions

TODO: Describe
Chapter

The design of the computing environment will be restricted by several points. These factors are given because of technological choices and requirements from up there.

- **Running on a NVIDIA DGX A100¹:** The environment will run on a NVIDIA DGX A100 workstation. The workstation has 80 CPUs and 8 GPUs installed. For this computing environment, 2 GPUs will be available.
- **Apache Spark for distributed computing:** Apache Spark will be used as a distributed computing framework.
- **Python as the main programming language:** Python will be used as the main programming language for Apache Spark applications. Therefore, examples will use Python code. Configurations for the system are optimized for using Python in production.

5.2 CI/CD

5.3 Identification of suitable Metrics for Scaling

To analyze the Apache Spark cluster computing performance suitable metrics have to be defined. As mention in SECTION XY, Apache Spark distributes its workload across multiple CPU cores. In addition, one objective of this thesis is to accelerate the computing performance of the Apache Spark cluster with GPUs. Therefore suitable performance metrics are CPU utilization of all Apache Spark worker and the utilization of all available GPUs. These

¹ The Universal System for AI Infrastructure - <https://www.nvidia.com/en-us/data-center/dgx-a100/>

utilization metrics will be based on the time, when the Apache Spark cluster is actively performing computations.

5.3.1 CPU Performance

To adapt to business needs, the CPU percentage of each Spark Worker will be calculated. Prometheus provides several metrics to calculate the CPU percentage. The CPU percentage of all Worker can be calculated as follows:

Each Apache Spark worker has a number of CPU cores to perform work. To calculate the CPU utilization of an Apache Spark worker,

$$SparkWorkerCPUUtilization = \frac{\sum CPUCoreUtilization}{NumberOfCPUCores} \quad (5.1)$$

$$OverallCPUUtilization = \frac{\sum SparkWorkerCPUUtilization}{NumberOfSparkWorker} \quad (5.2)$$

5.3.2 GPU Performance

The system has a fixed number of GPUs to use.

$$GPUUtilization = \frac{\sum SparkWorkerGPUUtilization}{NumberOfSparkWorker} \quad (5.3)$$

5.4 Computing environment Architecture

5.4.1 Overall

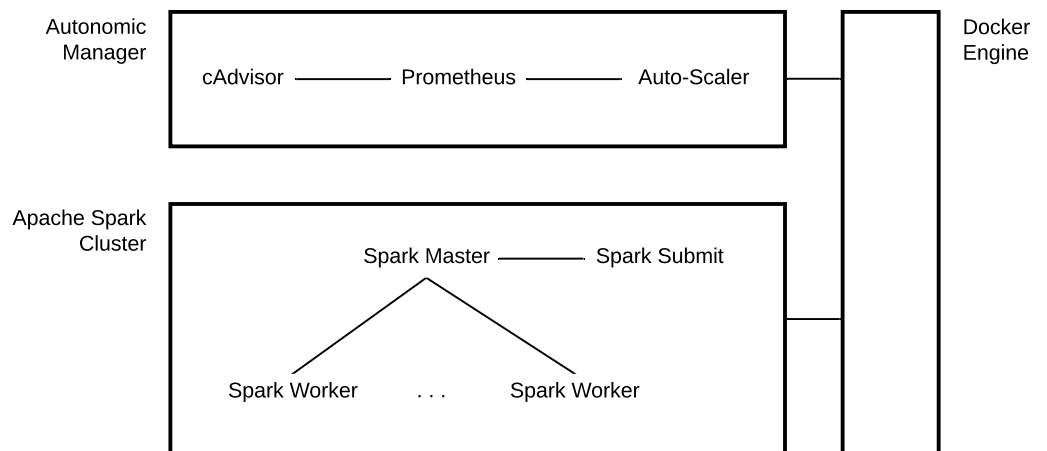


Figure 5.1: Overall cluster architecture - Source: Authors own model.

**TODO: Hier auf
Tabelle verweisen
(Anhang)
Metriken die von
Prometheus +
cAdvisor
bereitgestellt
werden**

Figure 5.1 illustrates the overall architecture with all components of the computing environment. The two main components in the environment are an Apache Spark cluster and an autonomic manager. The autonomic manager will be implemented according to the MAPE architecture (introduced in Section 2.3.3). It is responsible for monitoring and auto-scaling the Apache Spark cluster. To enable distributed computing, an Apache Spark cluster will be set up to execute machine learning applications. The autonomic manager and the Apache Spark cluster consist each of multiple nodes. Each node of a component will run as a Docker container. All containers will run as services in Docker Swarm mode. As described in Section 4.1.4, Docker Swarm mode will maintain a healthy state of all running containers. Therefore, the environment will enable self-healing according to the requirements of Autonomic Computing described in Section 2.3.

5.4.2 Apache Spark Cluster

Master and Worker

The Apache Spark cluster will consist of a Spark master node and a dynamic number of Spark worker nodes. The Spark master node is responsible to distribute the application workload across available Spark worker nodes. A Spark worker node will execute the workload given by the master node. Each Spark worker is homogeneous. Homogeneity is important to scale the number of worker nodes. To enable homogeneous nodes, each Spark worker node is a Docker container running the same Docker image. In addition, each worker is given the same computing resources. With homogeneous Spark worker nodes, each worker will respond as all other nodes. To enable GPU acceleration, each WORKER/MASTER will have the RAPIDS plugin installed. The cluster will be deployed in standalone mode. To be able to run Python applications

Spark Submit

Because the Apache Spark cluster will be executed in standalone mode, a node inside the cluster is required to run Spark applications. When a Spark application will be executed, a Spark Submit container will be deployed in the cluster. When the application has finished, the container will be automatically removed. Each app will be executed by a unique Spark Submit node. The Spark Submit node will be deployed via the CI pipeline (SECTION XY). The purpose of this

5.4.3 Autonomic Manager

The design of the autonomic manager needs to fulfill all requirements given by the MAPE architecture (described in Section 2.3.3). It will be responsible to monitor the performance of Apache Spark worker nodes in the environment, analyze the metrics and plan and execute scaling actions in

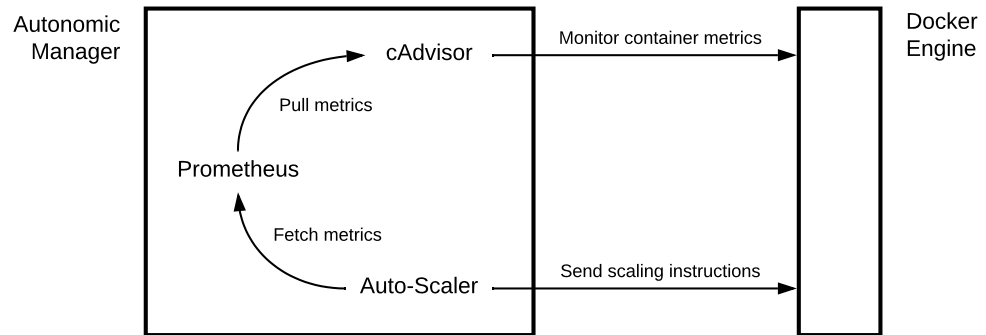


Figure 5.2: Autonomic manager component design - Source: Authors own model.

accordance to fulfill the performance goals. As illustrated in Figure 5.2, the autonomic manager consists of a cAdvisor, Prometheus and Auto-Scaler node. Together, all three nodes build a complete autonomic manager in accordance to the MAPE architecture. Each node will run as a Docker container. For monitoring, cAdvisor collects metrics from all available Docker container in the computing environment. Prometheus pulls metrics from cAdvisor. The Auto-Scaler will fetch metrics from Prometheus and send scaling instructions to the Docker engine.

Monitoring System

The design introduced in Section 5.4.1 is a dynamic changing computing environment. To monitor this dynamic environment, a monitoring-system is needed that fulfills the requirements described in Section 2.5.

cAdvisors will be used as an agent to collect performance metrics from all running Docker containers in the environment. Prometheus pulls the collected performance metrics from cAdvisor and stores the data as time-series data in its database. In addition, Prometheus provides a powerful multi-dimensional query language to aggregate and analyze the stored data.

Workflow

The workflow of the autonomic manager is implemented as a loop.

Figure 5.3 illustrates all steps of each component of the autonomic manager process according to the MAPE architecture.

The workflow starts with collecting metrics in the monitor phase. cAdvisor is responsible to collect metrics from the Docker engine. After that, Prometheus stores the collected metrics as time-series data in its database. Next, in the analyze phase, the Auto-Scaler needs to determine if a scaling action is needed. If scaling the Apache Spark worker nodes is not needed, the process has finished and will repeat from the beginning in the next period. If the performance is over- or under-utilized, a scaling action is needed. Then, the Auto-Scaler needs to determine how many Apache Spark worker

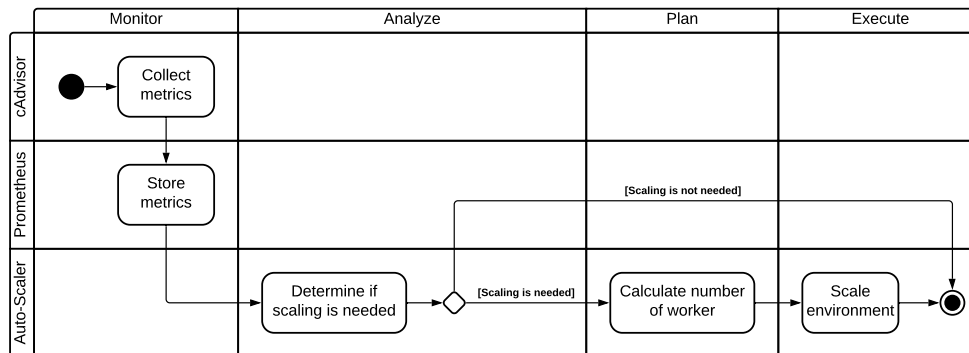


Figure 5.3: UML activity model of the autonomic manager process - Source: Authors own model.

nodes are needed to reach the performance goal. Lastly, the Auto-Scaler is responsible to send the scaling instructions to the Docker engine.

5.5 Auto-Scaler

The Auto-Scaler is a component of the the Autonomic Manager and is responsible for the analyze, plan and execution phase. Together with cAdvisor and Prometheus, the Auto-Scaler builds a complete Autonomic Manager according to Section 2.3.3. In addition, the Auto-Scaler implements the control-loop which is responsible to make adjustments in the environment according to the performance. To adjust the number of Apache Spark worker in the environment, the Auto-Scaler needs to perform queries on the performance metrics stored in Prometheus. After analyzing the performance metrics, the Auto-Scaler will send instructions to adjust the number of Apache Spark worker to the Docker engine. The Docker engine tells the Swarm manager to scale the replicas of the Apache Spark service in accordance to the Auto-Scaler instructions.

5.5.1 Configuration

The Auto-Scaler needs specific configuration properties to be able to collect the correct metrics from Prometheus and deploy new Apache Spark worker container in the environment. The following are properties that have to be defined to ensure that the Auto-Scaler is able to collect meaningful metrics and scale Apache Spark worker as expected.

General properties

- **Interval seconds:** The number of seconds when the loop has to repeat needs to be defined.
- **Cooldown period:** The duration in seconds, the Auto-Scaler has to wait after a scaling action was performed.

- **Recurrence factor:** To prevent too many scaling actions, the auto-nomic manager should only execute a scaling action, if the utilization thresholds is violated n times.
- **Prometheus URL:** The Auto-Scaler will fetch the configured metrics from the Prometheus REST API.

Metrics

To support to analyze multiple metrics, the user should be able to create a dynamic list of metrics. Each metric needs to have a variety of properties configured.

- **Target utilization:** The relative target utilization of a metrics needs to be defined to calculate the number of Spark worker to add or to remove to reach the defined goal.
- **Utilization thresholds:** To determine if a scaling action is needed, the scaling heat algorithm needs the minimum and maximum utilization defined by an administrator.
- **Query:** A PromQL query needs to be defined to collect the metric for all Spark Worker.

Apache Spark worker properties

- **Worker image:** To guarantee that each Spark worker is homogeneous, all worker container should be created with the same image.
- **Worker network:** To establish communication between all Spark worker and the Spark master, all new Spark worker container should be in the same network.
- **Worker thresholds:** The minimum and maximum number of concurrent Spark worker should be defined. To avoid the cold start effect, the minimum amount of worker should be 1.
- **Apache Spark master URI:** To distribute the workload across all Spark Worker, all Spark Worker need to communicate with the Spark master.

5.5.2 Analyze

In order to determine if a scaling-action is necessary, the Auto-Scaler has to process the collected metrics. During each period, the Auto-scaler queries the Prometheus time-series database with the configured queries to get all needed metrics. After the metrics are received, the Auto-Scaler determines if a scaling action is needed using the Scaling Heat algorithm (introduced in Section AB). If scaling is not necessary, the Auto-Scaler continues to collect metrics from Prometheus.

5.5.3 Plan

If a scaling-action is necessary, the Auto-Scaler is responsible to plan how to scale the number of Spark worker to satisfy the defined utilization goals. A scaling plan consists of instructions to add or remove Spark worker which will be send to the Docker engine. To calculate the number of Spark worker, needed to accomplish the defined target utilization, the Auto-Scaler uses the *Kubernetes Horizontal Pod Auto-Scaling* algorithm. In addition, the Auto-Scaler needs to check if the estimated number of Spark worker fall bellow the minimum threshold or exceed the maximum threshold of concurrent Spark worker.

5.5.4 Execute

After a scaling plan has been created, the Auto-Scaler needs to send the instructions to the Docker engine. After scaling the environment, it needs time for changes to take effect. Therefore a cooldown period will be activated after each scaling action. During the cooldown period, no scaling actions will be forwarded to the Docker engine.

Chapter 6

Implementation

6.1 Background

**TODO: Describe
Chapter**

Mentioned before, the goal is to create a self-healing and self-adapting computing environment. In addition, Apache Spark applications should be able to deploy automatically via a CI pipeline to the cluster. Therefore, to implement the conceptual design described in SECTION XY, the following requirements need to be implemented:

1. Implementing the Auto-Scaler Python module
2. Dockerizing the Auto-Scaler
3. Creating the computing environment in a Docker swarm
 - Create an Apache Spark cluster with GPU acceleration
 - Create an autonomic manager with a monitoring system and an Auto-Scaler
4. Implement a GitLab CI pipeline to deploy Apache Spark applications to the cluster

6.2 Auto-Scaler

The conceptual design of the Auto-Scaler is described in detail in Section 5.5. The Auto-Scaler will be implemented as a module in Python 3.8. To deploy the Auto-Scaler in a Docker swarm, the Auto-Scaler needs to be dockerized in a Docker Image.

6.2.1 Configuration

The configuration parameter for the Auto-Scaler have been introduced in Section 5.5.1. The configuration for the Auto-Scaler will be specified in a YAML file. Listing 6.1 describes an example of an Auto-Scaler configuration.

Overall, a configuration file is structured in three sections: General, metrics and worker. Table 6.1 lists all available configuration parameters. It describes the value type and the default value of each parameter. Some parameters are required to be defined by the administrator and have no default value.

General: The general section defines details about the scaling and heat algorithm and the Prometheus URL.

Metrics: Metrics is a list of performance metrics configuration parameters. The name of a metric (*cpu* and *gpu* in the example Listing 6.1) can be set by the administrator. A performance metric requires a query in the PromQL syntax. Additionally a target utilization is needed and the minimum and maximum utilization of the performance metric.

Worker: To scale the replicas of the Apache Spark worker service, the name of the Docker service needs to be set. In addition, the minimum and maximum number of concurrent worker nodes needs to be defined to prevent an overhead of running worker nodes.

```
1 general:
2   interval_seconds: 5
3   cooldown_period_seconds: 180
4   recurrence_factor: 3
5   prometheus_url: "http://localhost:9090"
6
7 metrics:
8   cpu:
9     query: 'sum(rate(container_cpu_user_seconds_total{image
10      = "spark-worker:3.0.1-hadoop2.7"}[30s]))'
11     target_utilization: 0.5
12     thresholds:
13       min: 0.2
14       max: 0.6
15
16   gpu:
17     query: 'sum(rate(container_cpu_user_seconds_total{image
18      = "spark-worker:3.0.1-hadoop2.7"}[30s]))'
19     target_utilization: 0.3
20     thresholds:
21       min: 0.2
22       max: 0.6
23
24 worker:
25   service_name: "computing_spark-worker"
26   thresholds:
27     min: 1
28     max: 30
```

Listing 6.1: Auto-Scaler configuration YAML file

Name	Type	Default
general		
interval_seconds	Integer	1
cooldown_period_seconds	Integer	180
recurrence_factor	Integer	1
prometheus_url	String	<i>Required</i>
metrics		
query	String	<i>Required</i>
target_utilization	Float	0.5
thresholds		
min	Float	0.5
max	Float	0.5
worker		
service_name	String	<i>Required</i>
thresholds		
min	Float	0.5
max	Float	0.5

Table 6.1: Auto-Scaler configuration parameter

6.2.2 Scaling the Apache Worker Service

docker service apache-worker scale replicas=4

Estimation of Necessary Scaling Actions

The Scaling Heat algorithm, introduced in Section 4.7, is being used to estimate if a scaling action is necessary. The algorithm is being used because it will prevent the Auto-Scaler to perform unnecessary scaling actions. During each interval, after performance metrics have been received from the monitoring system, a heat value will be calculated for each performance metric specified in the configuration under *metrics*. The algorithm uses a recurrence factor which has to be defined in design time. The Auto-Scaler configuration provides a parameter called *recurrence_factor* (see Table 6.1 for details).

To store and calculate the heat for each performance metric, a class called **HeatStore** was created. FIGURE XY describes the UML class diagram for the **HeatStore** Python class. The class can be used to retrieve, update and reset the heat for a list of performance metrics.

Calculating the Number of Needed Worker Nodes

The KHPA algorithm will be used to calculate how many worker are needed to reach the target utilization (see SECTION XY for algorithm details). In

this project, the calculation is done for the CPU and GPU utilization. The highest number of the desired worker node replicas is chosen.

Performing a Scaling Action

Docker provides a Python library for the Docker Engine¹. This library will be used to perform the swarm scaling action.

If worker need to be removed, it is necessary to check if the worker are running any applications. Removing a worker while an application is performing will cause the cancellation of the application. To check if applications

After a scaling action has been performed, a cooldown period will be applied. The cooldown period is needed because the number of desired worker nodes can keep fluctuating due to the dynamic nature of performance metrics.

6.2.3 Docker Image

- Dockerfile hier erklären - Wie den Auto-Scaler per konsole starten (mit config parameter)

6.3 Deployment of a Docker Swarm

6.3.1 Hardware

6.3.2 Software info

Hier tabelle mit versionen von eingesetzter software

6.3.3 Swarm

Vielleicht euch einfach das ganze kapitel Swarm nennen? - Dockerfile erläutern - GPUs (nur die 2 bestimmten)

6.3.4 Build Script

6.3.5 Apache Spark Cluster with GPU Acceleration

The Apache Spark cluster is created in standalone mode, see Section 4.2.3 for deployment details. The cluster consists of a single Apache Spark master node and a dynamic number of Apache Spark worker nodes. The master and worker container will run as a service in a swarm (see Section 4.1.4).

¹ Docker SDK for Python 4.4.1 Documentation - <https://docker-py.readthedocs.io/en/4.4.1/> (Accessed: 2021-01-05)

For submitting an application to the cluster, an individual container performing the `spark-submit` executable will be deployed. Each node runs in an independent Docker container.

Overall four Docker images are needed to create the Apache Spark cluster introduced in Section 5.4.2:

- Base image
- Master image
- Worker image
- Submit image

The master, worker and submit image require common packages to be installed and a set of common configuration. Therefore, an additional base image will serve as a base image.

Apache Spark Base Image Installation Details

The base image Dockerfile is available at Listing A.1. As parent, the base image uses the `nvidia/cuda:11.0-devel-ubuntu16.04` Docker image. The parent image runs Ubuntu² in version 16.04. Additionally the CUDA Toolkit and the NVIDIA GPU driver are already installed. Docker provides the ability to set build arguments. To be able to install a specific Apache Spark version, two arguments, can be set when building the Docker image. Apache Spark will be installed at `/opt/Spark`. This directory will be set as the working directory for the Docker image as well. Furthermore, required Ubuntu packages will be installed. This includes the Java Runtime Environment Version 8, which is a requirement for Apache Spark. To enable GPU acceleration on all Apache Spark nodes, the base image will install the compiled Java files for the RAPIDS plugin at `/opt/sparkRapidsPlugin` (for RAPIDS installation requirements, see Section 4.3.2). The `.jar` files can be downloaded in the maven repository. To enable Apache Spark to discover available GPUs, a GPU discovery script is needed (see Section 4.2.3 for details about resource allocation). This discovery script will be placed at `/opt/sparkRapidsPlugin` as well. The discovery script is introduced at Listing A.4.

Standalone Master and Worker Image

The master and worker image are build on top of the Apache Spark base image. Therefore, no additional installation steps are required. The master and worker nodes will be launched in standalone mode (see Section 4.2.3 for standalone mode details).

² Enterprise Open Source and Linux - <https://www.ubuntu.com/> (Accessed: 2021-01-03)

Master image: Implementation of the master node Dockerfile is available at Listing A.2. The master node image needs two ports configured: The Apache Spark service port and the port for the web user interface. The Apache Spark service port is set to 7077 and the web user interface port to 4040. To start the master node in standalone mode, the start-master.sh launchable will be set as image entrypoint which requires no additional arguments.

Worker image: Listing A.3 describes the implementation of the worker node Dockerfile. The port for the worker web interface will be exposed at 4041. To start the worker in standalone mode, the start-slave.sh executable will be set as entrypoint for the image. The launch script requires the master node URI as a parameter. To keep the configuration simple, the environment variable HIERDIEVAR will be set in the compose file. Listing 6.2 describes the configuration environment for the worker. As mentioned previously (in Section 5.1), for this project two GPUs are available on the DGX workstation. Furthermore, the worker needs to know where to find the GPU resource discovery script.

```
1 SPARK_WORKER_OPTS="-Dspark.worker.resource.gpu.amount=2 -
    Dspark.worker.resource.gpu.discoveryScript=/opt/
    sparkRapidsPlugin/getGpusResources.sh"
```

Listing 6.2: Environment configuration for all worker nodes

Submit Image

As mentioned in Section 5.1, a requirement is, that Apache Spark application will be implemented in Python. Accordingly, the pyspark Python module needs to be installed on all submit nodes. Apache Spark application will be placed at `/opt/spark-apps`. SECTION CI describes how an Apache Spark application will be copied to a submit node. As entrypoint, the image will perform a custom submit script (available at LISTING AB). This script performs the spark-submit executable (usage described in detail in Section 4.2.3).

6.3.6 Autonomic Manager

As mentioned in Section 5.4.3, the autonomic manager will consist of a monitoring system and the Auto-Scaler to create a complete control loop.

The monitoring system conceptual design was introduced in SECTION XY. It consists of a cAdvisor (SECTION XY) service and a Prometheus (SECTION XY) service. All modules will run as individual Docker services in the overall swarm. The following DOcker images will be used for the monitoring system:

- **cAdvisor:** google/cadvisor

- **Prometheus:** prom/prometheus

Prometheus Target Configuration

As mentioned in Section 4.4, Prometheus is a pull-based monitoring tool. It requires a list of targets to pull performance metrics from.

Listing 6.3 specifies the scrape configuration of the Prometheus system. The cAdvisor service is specified as a target. Prometheus will scrape every 5 seconds performance metrics from cAdvisor. All performance metrics will be labeled with the cAdvisor label. The cAdvisor service is available at cadvisor:8080.

```
1 scrape_configs:
2   - job_name: prometheus
3     scrape_interval: 5s
4     static_configs:
5       - targets: ["localhost:9090"]
6
7   - job_name: cadvisor
8     scrape_interval: 5s
9     static_configs:
10      - targets: ["cadvisor:8080"]
11        labels:
12          group: "cadvisor"
```

Listing 6.3: Prometheus target configuration in YAML syntax

6.4 Automatic Deployment of Apache Spark Applications

Vielleicht eher das Kapitel so nennen - gitlab-ci.yml erklären - Screenshot von webui output

7.1 Test Environment

TODO: Describe
Chapter

About DGX how powerful is the DGX ...

7.2 Algorithms and Datasets

Description of used datasets

7.2.1 K-MEANS

7.2.2 Naive Bayes Classifier

Chapter 8

Outlook

8.1 Optimizing Scaling

TODO: Describe
Chapter

Current approach: Wait until all applications have finished. Better approach: Blacklist worker for removing so no executor will be launched on the worker. Create an external shuffle service, so worker can be removed on runtime.

8.2 Reinforcement Learning for Auto-Scaling

bla bla

Chapter 9

Conclusion

9.1 Cluster architecture

TODO: Describe
Chapter

Appendix A

Apache Spark Cluster Implementation

```
1 FROM nvidia/cuda:11.0-devel-ubuntu16.04
2
3 LABEL maintainer="marcel.pascal.stolin@ipa.fraunhofer.de"
4
5 ARG SPARK_VERSION
6 ARG HADOOP_VERSION
7
8 # Install all important packages
9 RUN apt-get update -qy && \
10     apt-get install -y openjdk-8-jre-headless procps python3
11     python3-pip curl
12
13 # Install Apache Spark
14 RUN mkdir /usr/bin/spark/ && \
15     curl https://ftp-stud.hs-esslingen.de/pub/Mirrors/ftp.
16     apache.org/dist/spark/spark-${SPARK_VERSION}/spark-${
17     SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}.tgz -o spark.
18     tgz && \
19     tar -xf spark.tgz && \
20     mv spark-${SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}/*
21     /usr/bin/spark/ && \
22     rm -rf spark.tgz && \
23     rm -rf spark-${SPARK_VERSION}-bin-hadoop${HADOOP_VERSION
24     }/
25
26 # Add GPU discovery script
27 RUN mkdir /opt/sparkRapidsPlugin/
28 COPY getGpusResources.sh /opt/sparkRapidsPlugin/
29     getGpusResources.sh
30 ENV SPARK Rapids_DIR=/opt/sparkRapidsPlugin
31
32 # Install cuDF and RAPIDS
33 RUN curl -o ${SPARK Rapids_DIR}/cudf-0.15-cuda11.jar https
34     ://repo1.maven.org/maven2/ai/rapids/cudf/0.15/cudf-0.15-
35     cuda11.jar
36 RUN curl -o ${SPARK Rapids_DIR}/rapids-4-spark_2.12-0.2.0.
37     jar https://repo1.maven.org/maven2/com/nvidia/rapids-4-
38     spark_2.12/0.2.0/rapids-4-spark_2.12-0.2.0.jar
39 ENV SPARK Rapids_CUDF_JAR=${SPARK Rapids_DIR}/cudf-0.15-cuda11.jar
```

```

29 ENV SPARK Rapids_PLUGIN_JAR=${SPARK_RAPIDS_DIR}/rapids-4-
    spark_2.12-0.2.0.jar
30
31 # Set all environment variables
32 ENV JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
33 ENV SPARK_HOME /usr/bin/spark
34 ENV SPARK_NO_DAEMONIZE true
35 ENV PYSARK_DRIVER_PYTHON python3
36 ENV PYSARK_PYTHON python3
37 ENV PATH /usr/bin/spark/bin:/usr/bin/spark/sbin:$PATH
38
39 WORKDIR ${SPARK_HOME}

```

Listing A.1: Apache Spark base image Dockerfile

```

1 ARG SPARK_VERSION
2 ARG HADOOP_VERSION
3
4 FROM spark-base:$SPARK_VERSION-hadoop$HADOOP_VERSION
5
6 LABEL maintainer="marcel.pascal.stolin@ipa.fraunhofer.de"
7
8 # Set ports
9 ENV SPARK_MASTER_PORT 7077
10 ENV SPARK_MASTER_WEBUI_PORT 4040
11
12 EXPOSE 4040 7077
13
14 # Start master-node in standalone mode
15 ENTRYPOINT [ "sbin/start-master.sh" ]

```

Listing A.2: Apache Spark master image Dockerfile

```

1 ARG SPARK_VERSION
2 ARG HADOOP_VERSION
3
4 FROM spark-base:$SPARK_VERSION-hadoop$HADOOP_VERSION
5
6 LABEL maintainer="marcel.pascal.stolin@ipa.fraunhofer.de"
7
8 # Add spark-env
9 COPY spark-env.sh ${SPARK_HOME}/conf/spark-env.sh
10
11 # Set port
12 ENV SPARK_WORKER_WEBUI_PORT 4041
13
14 EXPOSE 4041
15
16 # Start worker-node
17 ENTRYPOINT ./sbin/start-slave.sh ${SPARK_MASTER_URI}

```

Listing A.3: Apache Spark worker image Dockerfile

```

1 #!/usr/bin/env bash
2
3 #

```

```

4 # Licensed to the Apache Software Foundation (ASF) under one
  # or more
5 # contributor license agreements. See the NOTICE file
  # distributed with
6 # this work for additional information regarding copyright
  # ownership.
7 # The ASF licenses this file to You under the Apache License
  # , Version 2.0
8 # (the "License"); you may not use this file except in
  # compliance with
9 # the License. You may obtain a copy of the License at
10 #
11 #     http://www.apache.org/licenses/LICENSE-2.0
12 #
13 # Unless required by applicable law or agreed to in writing,
  # software
14 # distributed under the License is distributed on an "AS IS"
  # BASIS,
15 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either
  # express or implied.
16 # See the License for the specific language governing
  # permissions and
17 # limitations under the License.
18 #
19
20 # This script is a basic example script to get resource
  # information about NVIDIA GPUs.
21 # It assumes the drivers are properly installed and the
  # nvidia-smi command is available.
22 # It is not guaranteed to work on all setups so please test
  # and customize as needed
23 # for your environment. It can be passed into SPARK via the
  # config
24 # spark.{driver/executor}.resource.gpu.discoveryScript to
  # allow the driver or executor to discover
25 # the GPUs it was allocated. It assumes you are running
  # within an isolated container where the
26 # GPUs are allocated exclusively to that driver or executor.
27 # It outputs a JSON formatted string that is expected by the
28 # spark.{driver/executor}.resource.gpu.discoveryScript
  # config.
29 #
30 # Example output: {"name": "gpu", "addresses
  #                 ": ["0", "1", "2", "3", "4", "5", "6", "7"]}
31
32 ADDR5=`nvidia-smi --query-gpu=index --format=csv,noheader |
  sed -e 'a' -e 'N' -e '$!ba' -e 's/\n/"/g'`
33 echo {"name": "gpu", "addresses": ["$ADDR5"]}

```

Listing A.4: GPU discovery script
 - Source: <https://github.com/apache/spark/blob/v3.0.1/examples/src/main/scripts/getGpusResources.sh> (Accessed: 2021-01-03)

Bibliography

- [BKFL17] BARNA, C. ; KHAZAEI, H. ; FOKAEFS, M. ; LITOIU, M.: Delivering Elastic Containerized Cloud Applications to Enable DevOps. In: *2017 IEEE/ACM 12th International Symposium on Software Engineering for Adaptive and Self-Managing Systems (SEAMS)*, 2017, S. 65–75
- [BMDM20] BULLINGTON-MCGUIRE, Richard ; DENNIS, Andrew K. ; MICHAEL, Schwartz: *Docker for Developers*. Packt Publishing, 2020. – ISBN 9781789536058
- [BP19] BASTOS, Joel ; PEDRO, Araujo: *Hands-On Infrastructure Monitoring with Prometheus*. Packt Publishing, 2019. – ISBN 9781789612349
- [CP17] CASALICCHIO, Emiliano ; PERCIBALLI, Vanessa: Auto-Scaling of Containers: The Impact of Relative and Absolute Metrics, 2017
- [CZ18] CHAMBERS, Bill ; ZAHARIA, Matei: *Spark: The Definitive Guide*. 1st. O'Reilly Media, 2018. – ISBN 9781491912218
- [DG04] DEAN, Jeffrey ; GHEMAWAT, Sanjay: MapReduce: Simplified Data Processing on Large Clusters. In: *OSDI'04: Sixth Symposium on Operating System Design and Implementation*. San Francisco, CA, 2004, S. 137–150
- [DG10] DEAN, Jeffrey ; GHEMAWAT, Sanjay: MapReduce: A Flexible Data Processing Tool. In: *Commun. ACM* 53 (2010), Januar, Nr. 1, 72–77. <http://dx.doi.org/10.1145/1629175.1629198>. – DOI 10.1145/1629175.1629198. – ISSN 0001–0782
- [Doc] DOCKER INC.: *Docker Documentation*. <https://docs.docker.com/>. – Accessed: 2020-12-06
- [Far17] FARCIC, Viktor: *The DevOps 2.1 Toolkit: Docker Swarm*. Packt Publishing, 2017. – ISBN 9781787289703
- [Far18] FARCIC, Viktor: *The DevOps 2.2 Toolkit*. Packt Publishing, 2018. – ISBN 9781788991278

- [FH10] FARLEY, David ; HUMBLE, Jez: *Continuous Delivery: Reliable Software Releases through Build, Test, and Deployment Automation, Video Enhanced Edition*. Addison-Wesley Professional, 2010. – ISBN 9780321670250
- [GBR11] GOSCINSKI, Andrzej ; BROBERG, James ; RAJKUMAR, Buyya: *Cloud Computing: Principles and Paradigms*. Wiley, 2011. – ISBN 9780470887998
- [Git] GITLAB INC.: *GitLab User Documentation*. <https://docs.gitlab.com/>. – Accessed: 2020-12-19
- [GOKB16] GANELIN, Ilya ; ORHIAN, Ema ; KAI, Sasaki ; BRENNON, York: *Spark*. Wiley, 2016. – ISBN 9781119254010
- [Goo] GOOGLE LLC: *cAdvisor Github Repository*. <https://github.com/google/cadvisor>. – Accessed: 2020-12-31
- [Gre20] GREGG, Brendan: *Systems Performance, 2nd Edition*. Pearson, 2020. – ISBN 9780136821694
- [HBK18] HIDRI, Klodjan K. ; BILAS, Angelos ; KOZANITIS, Christos: HetSpark: A Framework that Provides Heterogeneous Executors to Apache Spark. In: *Procedia Computer Science* 136 (2018), S. 118 – 127. <http://dx.doi.org/https://doi.org/10.1016/j.procs.2018.08.244>. – DOI <https://doi.org/10.1016/j.procs.2018.08.244>. – ISSN 1877-0509. – 7th International Young Scientists Conference on Computational Science, YSC2018, 02-06 July2018, Heraklion, Greece
- [JSAP04] JACOB, Bart ; SUDIPTO, Basu ; AMIT, Tuli ; PATRICIA, Witten: *A First Look at Solution Installation for Autonomic Computing*. IBM Redbooks, 2004
- [KC03] KEPHART, J. O. ; CHESS, D. M.: The vision of autonomic computing. In: *Computer* 36 (2003), Nr. 1, S. 41–50
- [LBMAL13] LORIDO-BOTRÁN, Tania ; MIGUEL-ALONSO, Jose ; LOZANO, Jose: *Comparison of Auto-scaling Techniques for Cloud Environments*, 2013
- [Lig12] LIGUS, Slawek: *Effective Monitoring and Alerting*. O'Reilly Media, Inc., 2012. – ISBN 9781449333522
- [LT11] L., Abbott M. ; T., Fisher M.: *Scalability Rules: 50 Principles for Scaling Web Sites*. Addison-Wesley Professional, 2011. – ISBN 9780132614016
- [LT15] L., Abbott M. ; T., Fisher M.: *The Art of Scalability: Scalable Web Architecture, Processes, and Organizations for the Modern Enterprise, Second Edition*. Addison-Wesley Professional, 2015. – ISBN 9780134031408

- [Luu18] LUU, Hien: *Beginning Apache Spark 2: With Resilient Distributed Datasets, Spark SQL, Structured Streaming and Spark Machine Learning library*. 1st. Apress, 2018. – ISBN 9781484235799
- [McD20] McDONALD, Carol: *ACCELERATING APACHE SPARK 3.X*. 2020
- [MPK⁺17] MAVRIDIS, Stelios ; PAVLIDAKIS, Manos ; KOZANITIS, Christos ; CHYSOS, Nikos ; STAMOULIAS, Ioannis ; KACHRIS, Christoforos ; SOUDRIS, Dimitrios ; BILAS, Angelos: VineTalk: Simplifying Software Access and Sharing of FPGAs in Datacenters. In: *2017 27th International Conference on Field Programmable Logic and Applications (FPL 2017)* (2017). – ISSN 1946–147X
- [Mur04] MURCH, Richard: *Autonomic Computing*. IBM Press, 2004. – ISBN 9780131440258
- [NK19] NICKOLOFF, Jeffrey ; KUENZLI, Stephen: *Docker in Action*. Manning Publications, 2019. – ISBN 9781617294761
- [NVI] NVIDIA CORPORATION: *Spark-Rapids Online Documentation*. <https://nvidia.github.io/spark-rapids/>. – Accessed: 2020-12-04
- [PNKM20] POTDAR, Amit ; NARAYAN, DG ; KENGOND, Shivaraj ; MULLA, Mohammed: Performance Evaluation of Docker Container and Virtual Machine. In: *Procedia Computer Science* 171 (2020), 01, S. 1419–1428. <http://dx.doi.org/10.1016/j.procs.2020.04.152>. – DOI 10.1016/j.procs.2020.04.152
- [PYY15] PEILONG LI ; YAN LUO ; NING ZHANG ; YU CAO: HeteroSpark: A heterogeneous CPU/GPU Spark platform for machine learning algorithms. In: *2015 IEEE International Conference on Networking, Architecture and Storage (NAS)*, 2015, S. 347–348
- [Ros17] ROSSEL, Sander: *Continuous Integration, Delivery, and Deployment*. Packt Publishing, 2017. – ISBN 9781787286610
- [SAP20] SRIRAMA, Satish N. ; ADHIKARI, Mainak ; PAUL, Souvik: Application deployment using containers with auto-scaling for microservices in cloud environment. In: *Journal of Network and Computer Applications* 160 (2020), 102629. <http://dx.doi.org/https://doi.org/10.1016/j.jnca.2020.102629>. – DOI <https://doi.org/10.1016/j.jnca.2020.102629>. – ISSN 1084–8045
- [Sin06] SINREICH, D.: *An architectural blueprint for autonomic computing*, 2006

- [SP20] SABHARWAL, Navin ; PANDEY, Piyush: *Monitoring Microservices and Containerized Applications: Deployment, Configuration, and Best Practices for Prometheus and Alert Manager*. Apress, 2020. – ISBN 9781484262160
- [Thea] THE APACHE FOUNDATION: *Arrow. A cross-language development platform for in-memory data*. <https://arrow.apache.org/>. – Accessed: 2020-12-03
- [Theb] THE APACHE SOFTWARE FOUNDATION: *Spark 3.0.1 Documentation*. <https://spark.apache.org/docs/3.0.1/>. – Accessed: 2020-09-12
- [Thec] THE LINUX FOUNDATION: *Prometheus Online Documentation*. <https://prometheus.io/docs/>. – Accessed: 2020-12-03
- [Vad18] VADAPALLI, Sricharan: *DevOps: Continuous Delivery, Integration, and Deployment with DevOps*. Packt Publishing, 2018. – ISBN 9781789132991
- [VMD⁺13] VAVILAPALLI, Vinod K. ; MURTHY, Arun C. ; DOUGLAS, Chris ; AGARWAL, Sharad ; KONAR, Mahadev ; EVANS, Robert ; GRAVES, Thomas ; LOWE, Jason ; SHAH, Hitesh ; SETH, Siddharth ; SAHA, Bikas ; CURINO, Carlo ; O’MALLEY, Owen ; RADIA, Sanjay ; REED, Benjamin ; BALDESCHWIELER, Eric: *Apache Hadoop YARN: Yet Another Resource Negotiator*. New York, NY, USA : Association for Computing Machinery, 2013 (SOCC ’13). – ISBN 9781450324281
- [Wil12] WILDER, Bill: *Cloud Architecture Patterns*. O’Reilly Media, Inc., 2012. – ISBN 9781449319779
- [YSH⁺16] YUAN, Y. ; SALMI, M. F. ; HUAI, Y. ; WANG, K. ; LEE, R. ; ZHANG, X.: *Spark-GPU: An accelerated in-memory data processing engine on clusters*. In: *2016 IEEE International Conference on Big Data (Big Data)*, 2016, S. 273–283
- [ZCD⁺12] ZAHARIA, Matei ; CHOWDHURY, Mosharaf ; DAS, Tathagata ; DAVE, Ankur ; MA, Justin ; MCCAULY, Murphy ; FRANKLIN, Michael J. ; SHENKER, Scott ; STOICA, Ion: *Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing*. In: *9th USENIX Symposium on Networked Systems Design and Implementation (NSDI 12)*. San Jose, CA : USENIX Association, April 2012. – ISBN 978-931971-92-8, 15–28
- [ZCF⁺10] ZAHARIA, Matei ; CHOWDHURY, Mosharaf ; FRANKLIN, Michael J. ; SHENKER, Scott ; STOICA, Ion: *Spark: Cluster Computing with Working Sets*. In: *Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing*. USA : USENIX Association, 2010 (HotCloud’10), S. 10