

Design and Implementation of a Tool to Collect Execution- and
Service-Data of Big Data Analytics Applications

Bachelor's Thesis

for obtaining the academic degree
Bachelor of Science (B.Sc.)

at

Beuth Hochschule für Technik Berlin
Department Informatics and Media VI
Degree Program Medieninformatik

1. Examiner and Supervisor: Prof. Dr. Stefan Edlich
2. Examiner: Prof. Dr. Elmar Böhler

Submitted by: Markus Lamm
Matriculation number: s786694
Date of submission: 06.09.2016

Acknowledgements

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Objective	1
1.3	Structure of thesis	2
2	Theoretical Foundations	3
2.1	Big-Data-Analytics-Architectures	3
2.2	Stream-Processing	3
2.2.1	Apache Flink	3
2.2.2	Apache Kafka	3
3	Data Analysis	4
3.1	System data	4
3.2	Java Management Extensions (JMX)	4
3.3	Representational State Transfer (REST)	4
3.4	Data Quality	4
4	Requirements	5
4.1	Collection	5
4.2	Transport	5
4.3	Persistence	5
5	Architecture and Implementation	6
5.1	Collected data as time-series based stream	6
5.2	Microservices and Service-Discovery	6
5.3	System components	6
5.3.1	CollectorClient	6

5.3.2	Service-Discovery	6
5.3.3	CollectorManager	6
5.3.4	Message-Broker	6
5.3.5	Indexer	7
5.3.6	Persistence	7
6	Evaluation	8
7	Conclusion	9
8	Beispiele	10
8.1	Quelltext	10
8.2	Bild	10
8.3	Text Formatierungen und sonstiges	11
8.3.1	Listen	11
8.3.2	Text Hervorhebungen	12
8.4	Tabelle	13
8.5	Long-Table	13
8.6	Literaturverweis	14
8.7	Onlineverweise	14
8.8	Glossar	14
8.9	Abkürzungsverzeichnis	14
	List of Figures	A
	List of Tables	B
	List of Source Codes	C
	Index	D
	Bibliography	E
	Online resources	F
	Image resources	G

Anhang A	H
A.1 Diagrams	H
A.1.1 Use Case diagram	H
A.1.2 Class diagrams	I
A.1.3 Sequence diagrams	N
A.1.4 Component diagram	O
A.1.5 Deployment diagram	P
A.2 Tabelle	P
A.3 Screenshot	P
A.4 Graph	P
Eigenständigkeitserklärung	Q

1 Introduction

1.1 Motivation

According to a survey in Germany, nine out of ten companies (89 percent) analyze large volumes of data for operational decision-making processes using modern Big Data Analytics Architectures, where 48 percent of the respondents see the greatest potential of Big Data [Jör14]. The analysis of continuous data streams is taking up a growing importance for companies and therefore constitutes an important factor for business success.

Collecting, storing and analyzing system and operational data of Big Data Architectures is therefore an essential tool in order to ensure successful operation. Even though logfiles are usefull for tracing problems in software systems, problems can be tracked and potential sources of error can be identified much earlier by collecting and analyzing execution and service data at runtime to describe the state of the system at a given point in time.

Due to the distributed character of Big Data Applications, where a system is composed of several interacting components, the examination of log data is not an adequate choice to gain insight into the entire system [Les14].

1.2 Objective

The main goal of the thesis is the design and implementation of a software system to ingest and store system and operational data of Big Data Analytics Applications on the example of the streaming frameworks Apache Flink and Apache Kafka. It should be examined which data is available and can be collected at all, what data is relevant and how to collect from source systems. Furhermore, the collected data must be stored

in a persistence system to become available for possible consumers like visualization applications, analytical processes or as a data source for applications from the context of Machine Learning for example.

1.3 Structure of thesis

After a short introduction to the topics and the main goals of the present thesis in this chapter, the Chapter 2 covers the theoretical foundations of Big Data Analytics Applications, discusses the concept of "stream-processing" and introduces Apache Flink and Apache Kafka as representatives of widely used stream-processing frameworks.

Chapter 3 investigates which sources for collecting data exist for Apache Flink and Kafka and which data should be collected and stored in a persistence system regarding to its relevance and data quality.

The requirements and the target definition of the software-system will be introduced in Chapter 4, Chapter 5 describes the software solution by giving a detailed conceptional overview of the software components and providing implementation details for selected items.

In chapter 6 we'll see how to setup the technical environment for the usage of the prototype to verify the correct functionality related to the requirements defined in Chapter 4.

The last Chapter 7 covers a conclusion and summary of the present work.

2 Theoretical Foundations

2.1 Big-Data-Analytics-Architectures

2.2 Stream-Processing

Stream processing is the real-time processing of data continuously, concurrently, and in a record-by-record fashion. in which data is treated not as static tables or files, but as a continuous infinite stream of data integrated from both live and historical sources.

Benefits:

- Accessibility: live data can be used while still in motion, before being stored.
- Completeness: historical data can be streamed and integrated with live data for more context.
- High throughput: high-velocity, high-volume data can be processed with minimal latency.

2.2.1 Apache Flink

2.2.2 Apache Kafka

3 Data Analysis

3.1 System data

3.2 Java Management Extensions (JMX)

3.3 Representational State Transfer (REST)

3.4 Data Quality

4 Requirements

4.1 Collection

4.2 Transport

4.3 Persistence

5 Architecture and Implementation

5.1 Collected data as time-series based stream

5.2 Microservices and Service-Discovery

5.3 System components

5.3.1 CollectorClient

The CollectorClient tier is our entry point for bringing data into the system...

5.3.2 Service-Discovery

Registration for CollectorClients

5.3.3 CollectorManager

Gives overview, uses Consul as service-discovery

5.3.4 Message-Broker

Transport, "Event-Log"

5.3.5 Indexer

Receive messages from Kafka, route data, create ES index

5.3.6 Persistence

ES as search index for time-series based data, easy visualization with Kibana

6 Evaluation

7 Conclusion

8 Beispiele

Im Kapitel Beispiele (siehe chapter 8) werden die möglichen Funktionen und Möglichkeiten des LaTeX-Dokuments demonstriert.

8.1 Quelltext

Nachfolgend der Codeauszug 8.1.

```
1  /**
2  * The HelloWorldApp class implements an application that
3  * simply prints "Hello World!" to standard output.
4  */
5  class HelloWorldApp {
6      public static void main(String[] args) {
7          System.out.println("Hello World!"); // Display the string.
8      }
9  }
```

Codeauszug 8.1: Hello World

8.2 Bild

Die rechts zu sehende Grafik demonstriert die Möglichkeiten des Paketes „wrapfig“. Grafiken innerhalb einer „wrapfigure“ können entweder links oder rechts von Text umlaufen werden.

Die nachfolgende Figure 8.2 demonstriert die Darstellung eines „*.jpg“ Bildes innerhalb des Textes (beim Einfügen kann auf die Endung verzichtet werden, solange der Name einzigartig ist). Zusätzlich enthält dieses einen Untertitel der über das bereits verwendete Label verlinkt werden kann. Der Untertitel erscheint im Abbildungsverzeichnis (Abbvz.).



Figure 8.1: Beispielbild [PEX]

8.3 Text Formatierungen und sonstiges

Dieser Text enthält eine Fußnote¹.

8.3.1 Listen

Listen könne sowohl mit Bullet points als auch mit Zahlen erstellt werden

- Eine Liste mit Bullet points
 - Ein weiteres Element
1. Eine Liste mit Zahlen
 2. Ein weiteres Element

¹Fußnoten sind Anmerkungen, die im Druck-Layout aus dem Fließtext ausgelagert werden, um den Text flüssig lesbar zu gestalten.

8.3.2 Text Hervorhebungen

The problem with internet quotes is that you can't always depend on their accuracy

— Abraham Lincoln, 1864

"Inspirierende Zitate können mit epigraph eingefügt werden

The problem with internet quotes is
that you can't always depend on their
accuracy

Abraham Lincoln, 1864

Seitenumbrüche können nur direkt nach Text geschrieben werden, sonst lässt sich das Latex nicht mehr compilieren.



Figure 8.2: Beispielbild [PEX]

8.4 Tabelle

Nachfolgend Table 8.1.

Inhaber: Alice
Peer (Ersteller): Bob
Öffentlicher Schlüssel des Inhabers: F2 D2 0E ED FA 4E 9E 0A F2 DD 23 8A 32 44 F3 E9
Gültigkeit: 2015-07-01 – 2016-06-30

Table 8.1: Digitales Zertifikat

8.5 Long-Table

Die „Long-Table“ kann über definierte Header und Footer über Seitenumbrüche hinweg angezeigt werden.

Version	Codename	API	Verteilung
2.2	Froyo	8	0.1%
2.3.3 - 2.3.7	Gingerbread	10	2.7%
4.0.3 - 4.0.4	Ice Cream Sandwich	15	2.5%
4.1.x	Jelly Bean	16	8.8%
4.2.x		17	11.7%
4.3		18	3.4%
4.4	KitKat	19	35.5%

Fortsetzung auf nachfolgender Seite

Fortsetzung - Verteilung der Androidversionen (Stand 01.02.2016)

Version	Codename	API	Verteilung
5.0	Lollipop	21	17.0%
5.1		22	17.1%
6.0	Marshmallow	23	1.2%

Table 8.2: Verteilung der Androidversionen (Stand: 01.02.2016)

8.6 Literaturverweis

Weil für die alte und die neue Rechtschreibung verschiedene Trennregeln gelten, sind Deutsch mit alter Rechtschreibung und Deutsch mit neuer Rechtschreibung zwei verschiedene Sprachen ([Kna09], S. 192).

8.7 Onlineverweise

Siehe Google.de [Goo].

8.8 Glossar

Der Glossar enthält die Beschreibung verwendeter Begriffe für das bessere Verständnis gegenüber dem Leser. Beispiele sind: Berlin, Outsourcing, Application Service Providing, Policy und PCI Express.

8.9 Abkürzungsverzeichnis

Das Abkürzungsverzeichnis listet alle verwendeten Abkürzungen auf. Einige Beispiele sind Serial Attached SCSI (SAS), Compact Disk (CD), Local Area Network (LAN) und

Internationale Organisation für Normung (ISO). Die erneute Verwendung zeigt nur noch die Abkürzung: SAS, CD, LAN und ISO.

List of Figures

8.1	Beispielbild [PEX]	11
8.2	Beispielbild [PEX]	12
A.1	Use Case Diagramm	H
A.2	Class diagram 'JvmCollector'	I
A.3	Class diagram 'DStatCollector'	J
A.4	Class diagram 'FlinkRestCollector'	K
A.5	Class diagram 'FlinkJmxCollector'	L
A.6	Class diagram 'KafkaBrokerJmxCollector'	L
A.7	Class diagram 'CollectorClient'	M
A.8	Class diagram 'CollectorManager'	N
A.9	Sequence diagram 'Client discovery'	N
A.10	Sequence diagram 'Client scheduling'	O
A.11	Component diagram	O
A.12	Deployment diagram	P

List of Tables

8.1	Digitales Zertifikat	13
8.2	Verteilung der Androidversionen (Stand: 01.02.2016)	14

List of Source Codes

8.1	Hello World	10
-----	-----------------------	----

Index

A

alte 14

D

Darstellung 11

T

Trennregeln 14

U

und 10, 15

Untertitel 11

Bibliography

- [Kna09] Joerg Knappen. *Schnell ans Ziel mit LATEX 2e* -. ueberarbeitete und erweiterte Auflage. Muenchen: Oldenbourg Verlag, 2009. ISBN: 978-3-486-59015-9.
- [Les14] Tammo van Lessen. “Wissen, was läuft - Mit Laufzeitmetriken den Überblick behalten”. In: *Javamagazin* 10000.11 (2014), pp. 48–52.

Online resources

- [Goo] Google. *Google*. URL: <http://www.google.de> (visited on 10/06/2015).
- [Jör14] u.a Jörg Bartel Axel Mester. *Big Data Technologien – Wissen für Entscheider*. 2014. URL: <https://www.bitkom.org/Publikationen/2014/Leitfaden/Big-Data-Technologien-Wissen-fuer-Entscheider/140228-Big-Data-Technologien-Wissen-fuer-Entscheider.pdf> (visited on 08/06/2016).

Image resources

[PEX] PEXELS. *Black and white branches tree*. URL: <https://www.pexels.com/photo/black-and-white-branches-tree-high-279/>.

A

A.1 Diagrams

A.1.1 Use Case diagram

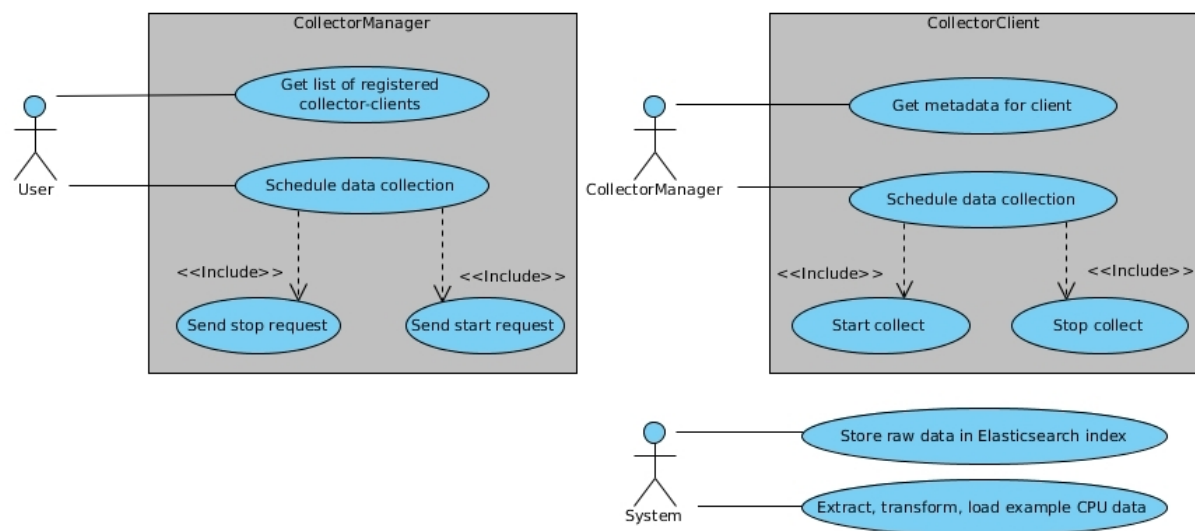


Figure A.1: Use Case Diagramm

A.1.2 Class diagrams

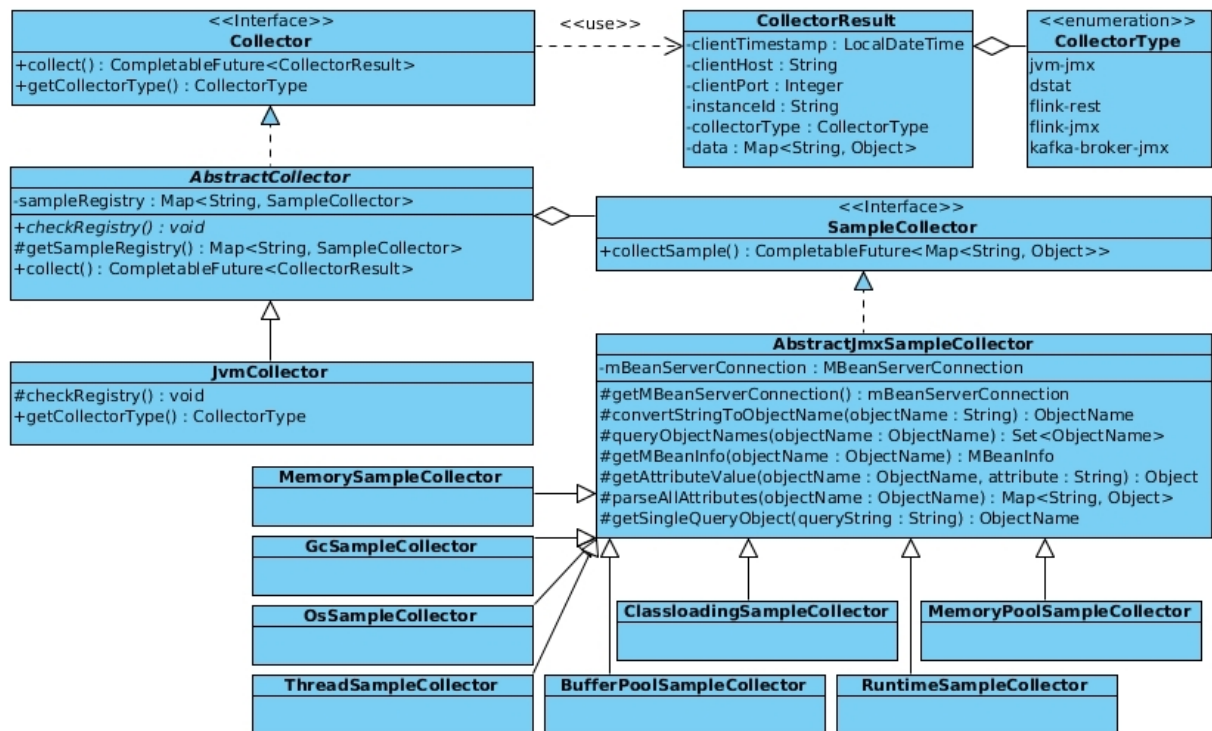


Figure A.2: Class diagram 'JvmCollector'

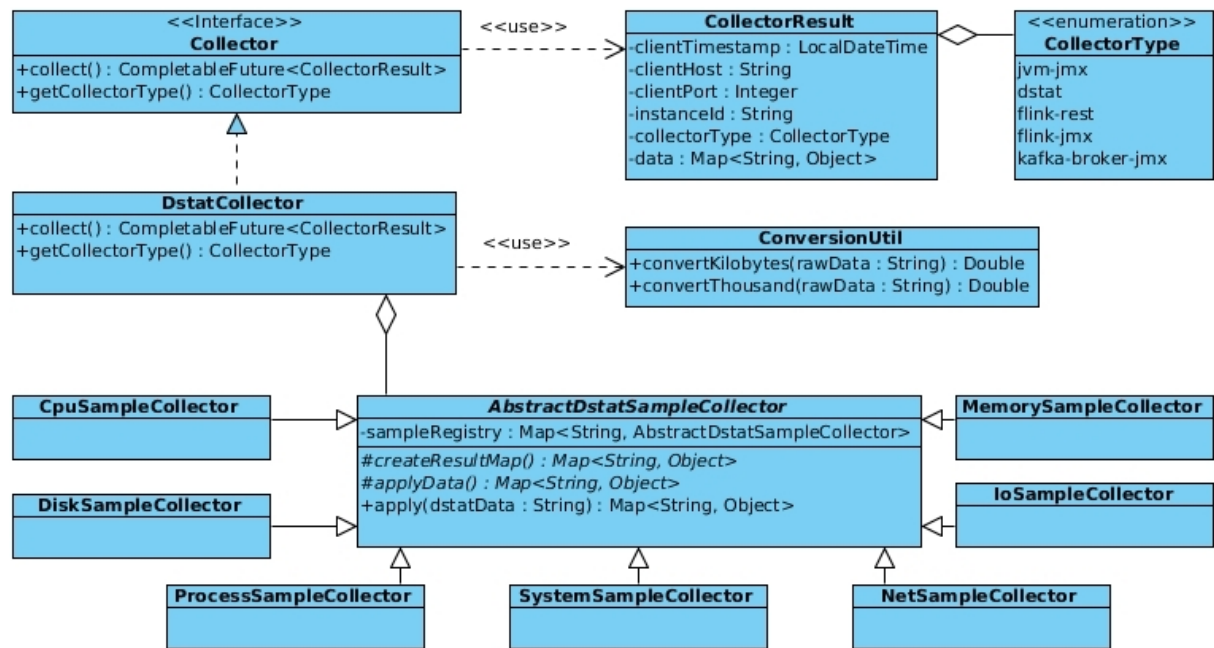


Figure A.3: Class diagram 'DStatCollector'

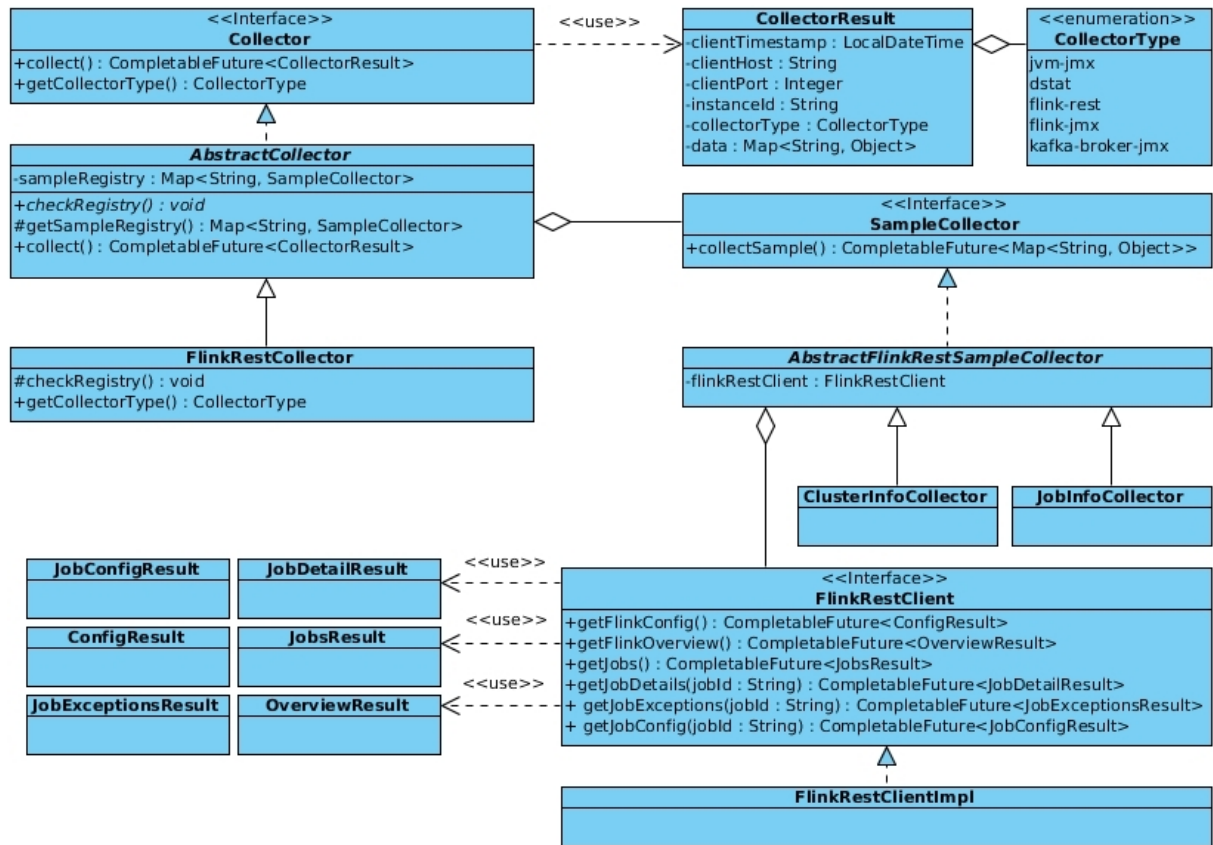


Figure A.4: Class diagram 'FlinkRestCollector'

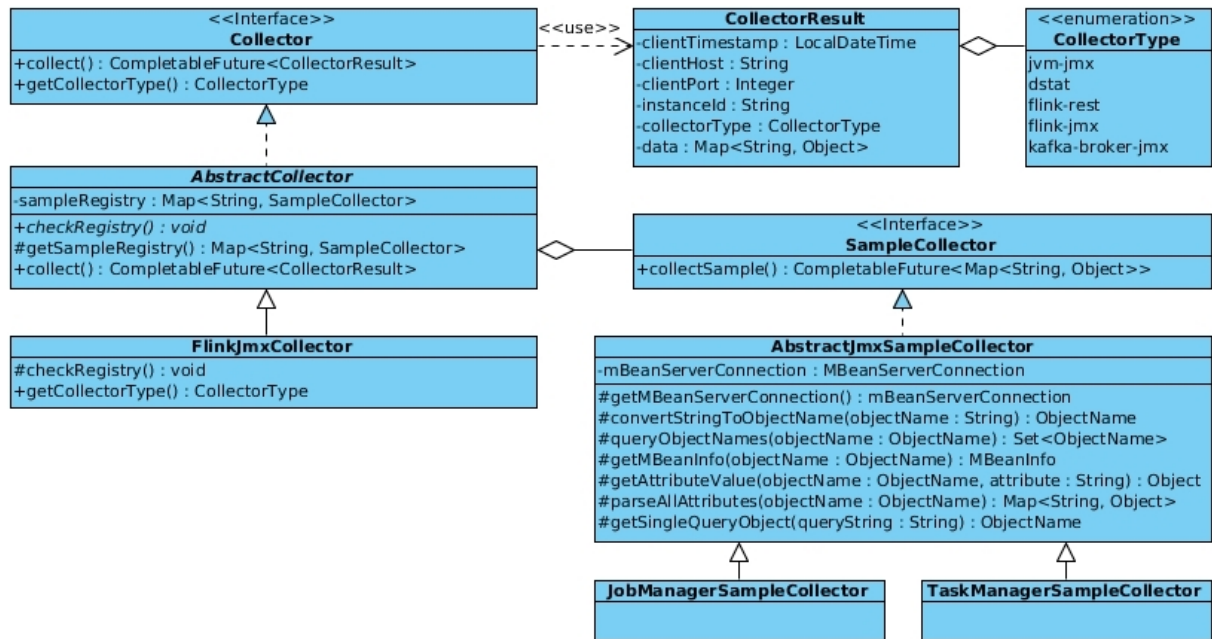


Figure A.5: Class diagram 'FlinkJmxCollector'

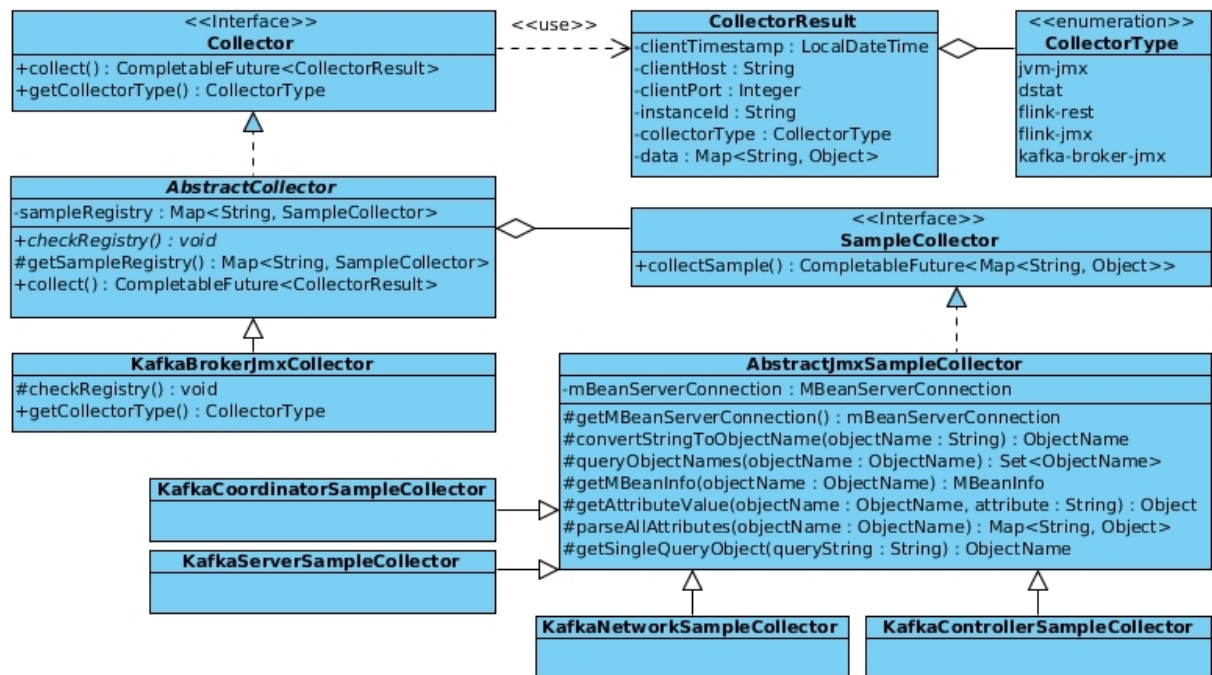


Figure A.6: Class diagram 'KafkaBrokerJmxCollector'

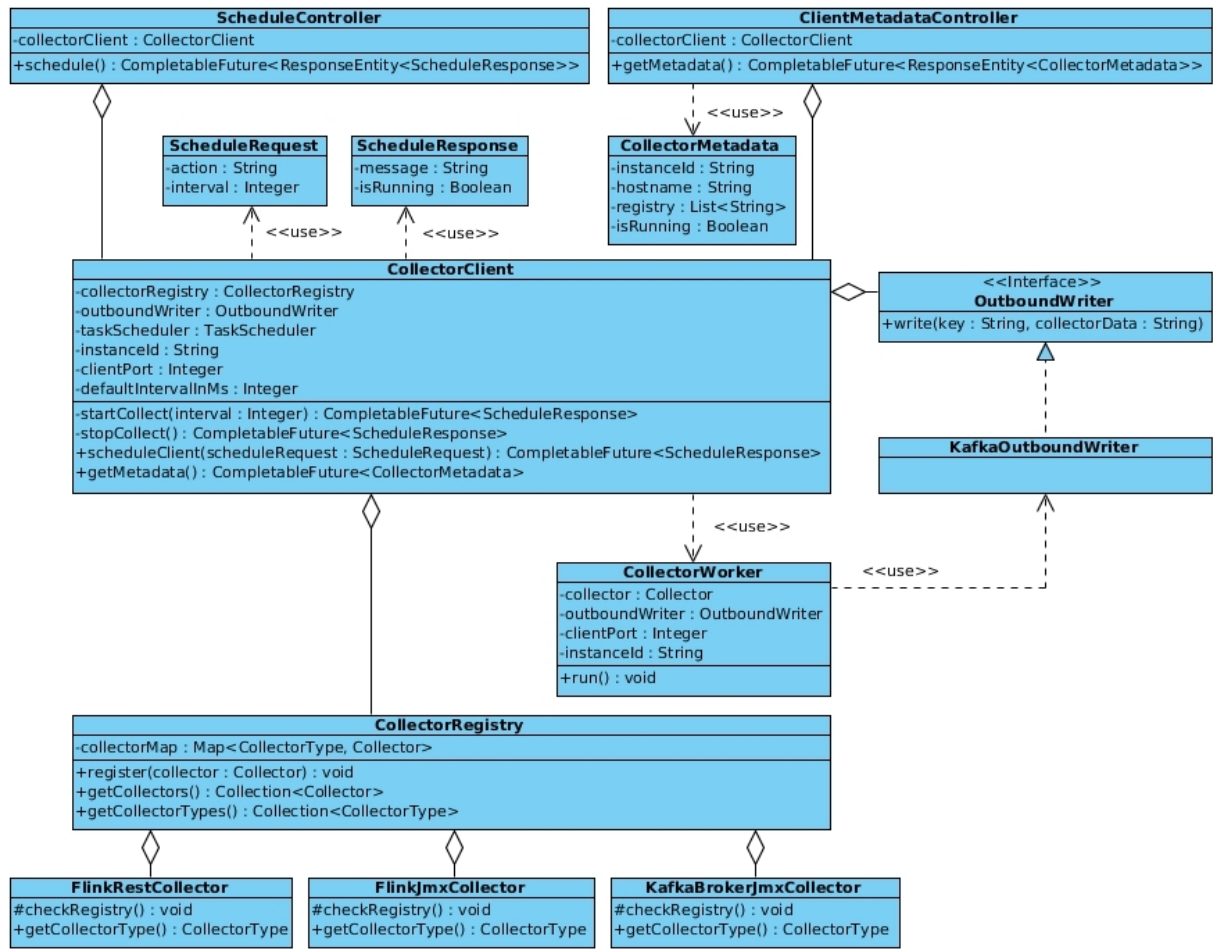


Figure A.7: Class diagram 'CollectorClient'

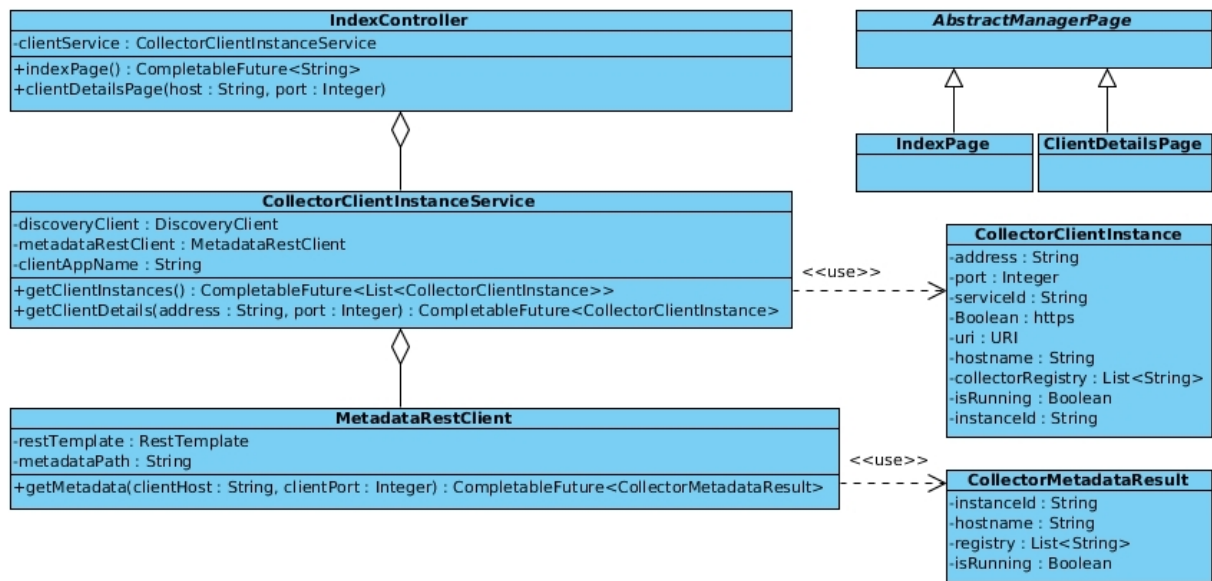


Figure A.8: Class diagram 'CollectorManager'

A.1.3 Sequence diagrams

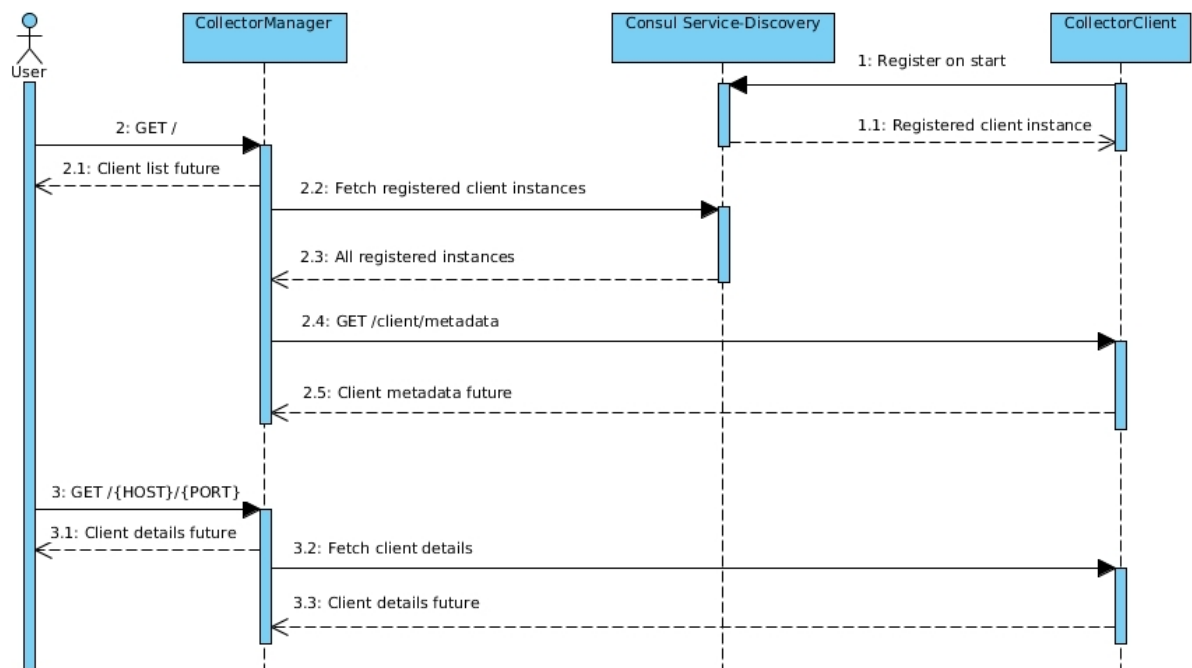


Figure A.9: Sequence diagram 'Client discovery'

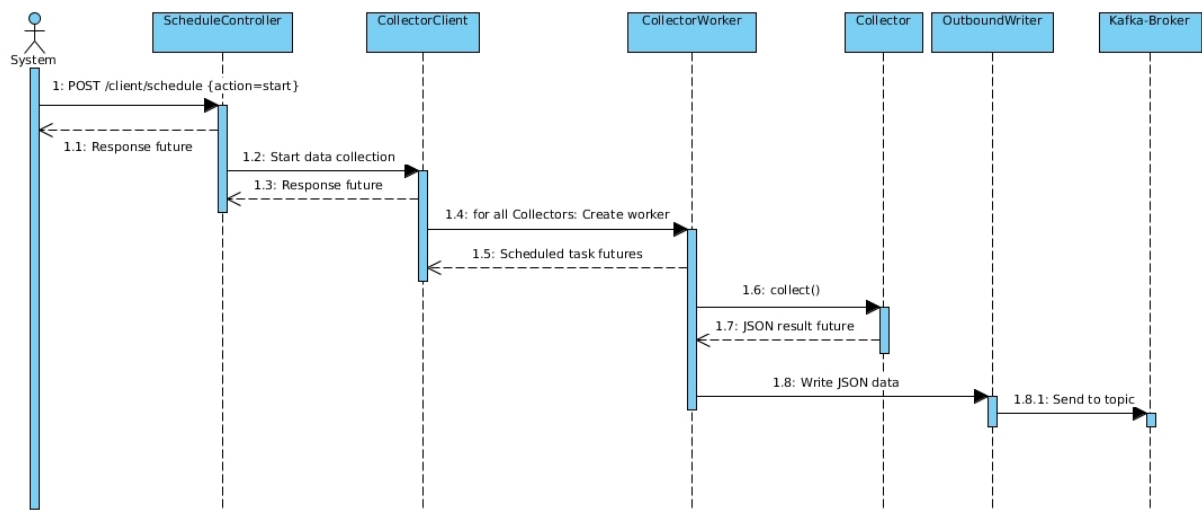


Figure A.10: Sequence diagram 'Client scheduling'

A.1.4 Component diagram

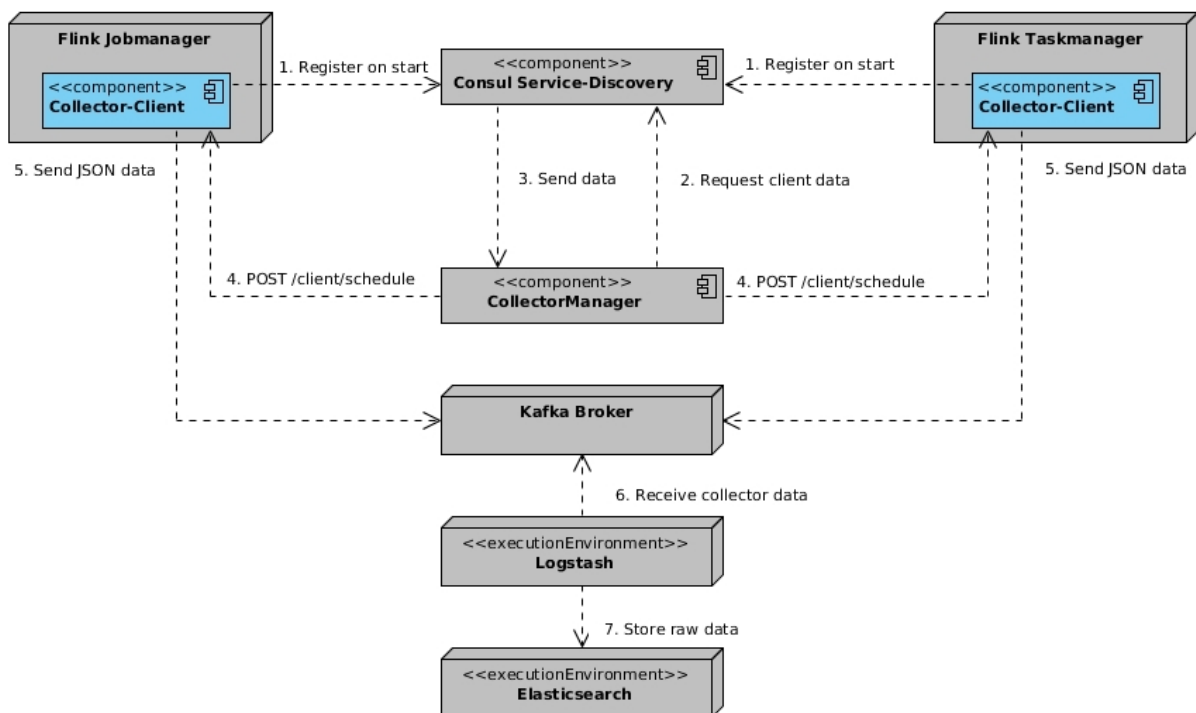


Figure A.11: Component diagram

A.1.5 Deployment diagram

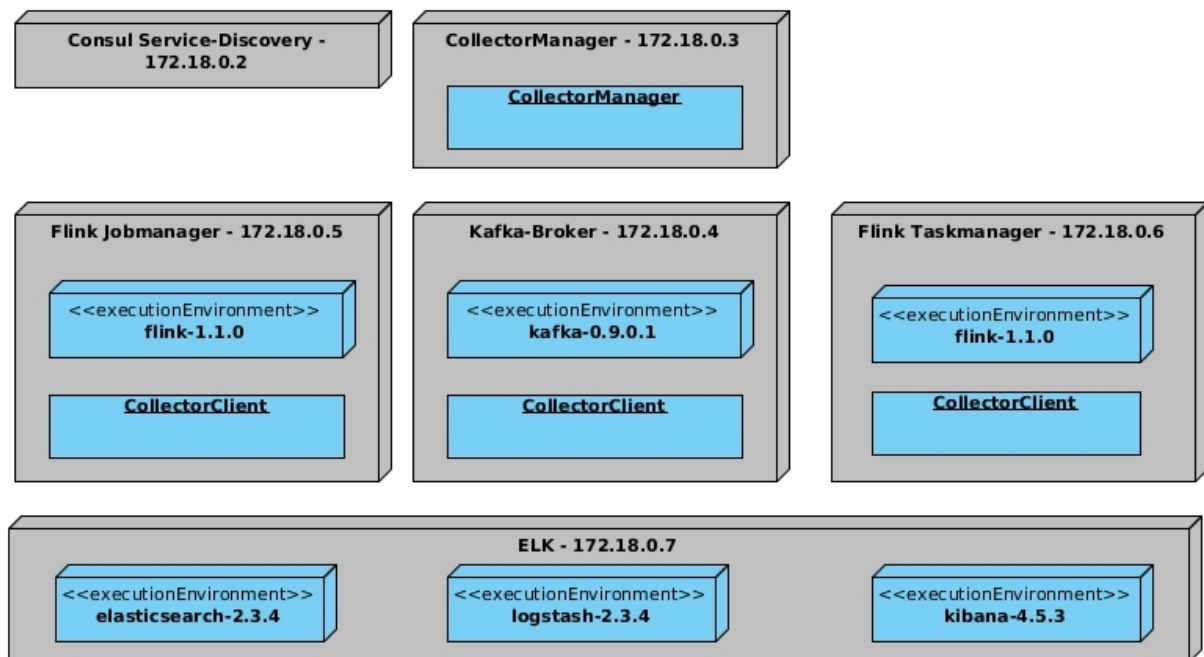


Figure A.12: Deployment diagram

A.2 Tabelle

A.3 Screenshot

A.4 Graph

Eigenständigkeitserklärung

Hiermit versichere ich, dass ich die vorliegende Masterarbeit selbstständig und nur unter Verwendung der angegebenen Quellen und Hilfsmittel verfasst habe. Die Arbeit wurde bisher in gleicher oder ähnlicher Form keiner anderen Prüfungsbehörde vorgelegt.

Stadt, den xx.xx.xxxx

Max Mustermann