



Hochschule für Angewandte Wissenschaften Hamburg
Hamburg University of Applied Sciences

Bachelorarbeit

Daniel Kirchner

Skalierbare Datenanalyse mit Apache Spark

*Fakultät für Technik und Informatik
Studiendepartment Informatik*

*Faculty of Engineering and Computer Science
Department of Computer Science*

Daniel Kirchner

Skalierbare Datenanalyse mit Apache Spark

Bachelorarbeit eingereicht im Rahmen der Bachelorprüfung

im Studiengang Bachelor of Science Angewandte Informatik
am Department Informatik
der Fakultät Technik und Informatik
der Hochschule für Angewandte Wissenschaften Hamburg

Betreuender Prüfer: Prof. Dr. Kahlbrandt
Zweitgutachter: Prof. Dr. Zukunft

Eingereicht am: 1. Januar 2345

Daniel Kirchner

Thema der Arbeit

Skalierbare Datenanalyse mit Apache Spark

Stichworte

Schlüsselwort 1, Schlüsselwort 2

Kurzzusammenfassung

Dieses Dokument ...

Daniel Kirchner

Title of the paper

Scalable Data Analysis with Apache Spark

Keywords

keyword 1, keyword 2

Abstract

This document ...

Inhaltsverzeichnis

1	Einführung	1
1.1	Motivation	1
1.2	Kontextabgrenzung	3
2	Vorstellung von Apache Spark	4
2.1	Überblick	5
2.2	Kernkonzepte	5
2.2.1	Nutzung von Arbeitsspeicher	6
2.2.2	Nutzung von persistentem Speicher	6
2.2.3	Nutzung von CPUs	6
2.2.4	Scheduling/Shuffling	6
2.2.5	Kern-API	6
2.3	Standardbibliotheken	6
2.3.1	Dataframes/Spark SQL	6
2.3.2	MLlib	6
2.3.3	Streaming	6
2.3.4	GraphX	6
2.4	Betrieb und Security	6
2.5	Entwicklergemeinschaft	6
2.6	Verwandte Produkte	7
2.6.1	YARN	7
2.6.2	Mesos	7
2.6.3	Flink	7
3	Spark in der Praxis	8
3.1	Echtzeitbewertung von Twitter-Accounts nach ihrer Relevanz für Spark	8
3.1.1	Beschreibung des Problems	8
3.1.2	Hardwarekontext und Performance-Basisdaten	9
3.1.3	Architekturübersicht	11
3.1.4	Detaillierte Lösungsbeschreibung	11
3.1.5	Ergebnisse	13
3.2	Evaluierung einer spark-basierten Implementation von CDOs auf einem HPC Cluster mit nicht-lokalem Storage	14
3.2.1	Beschreibung des Problems	14
3.2.2	Hardwarekontext und Performance-Basisdaten	14
3.2.3	Architekturübersicht	14

3.2.4	Detaillierte Lösungsbeschreibung	14
3.2.5	Ergebnisse	14
4	Schlussbetrachtung	15
4.1	Kritische Würdigung der Ergebnisse	15
4.2	Ausblick und offene Punkte	15
	Glossar	17

Tabellenverzeichnis

3.1 –DUMMY– Netzwerkdurchsatz 9

Abbildungsverzeichnis

1.1	Google Trends	2
2.1	Verteilung der Software-Komponenten im Cluster	5
3.1	svg image	10

Listings

2.1	Word Count in der Spark Konsole	4
3.1	Treiber für Testanwendung (Programmiersprache Scala)	12

1 Einführung

1.1 Motivation

Die Entwicklung TESTETSTETS22222222 **Request Evaluate Reply Loop** und Verbesserung von Frameworks zur Verarbeitung großer Datenmengen ist zur Zeit hochaktuell und sehr im Fokus von Medien und Unternehmen [VERWEIS]. Verschiedene Programme und Paradigmen konkurrieren um die schnellste, bequemste und stabilste Art großen Datenmengen einen geschäftsfördernden Nutzen abzuringen.

Mit dem Begriff „große Datenmengen“ oder „Big Data“ werden in dieser Arbeit solche Datenmengen zusammengefasst, die die Kriterien Volume, Velocity, Variety¹ erfüllen oder „Datenmengen, die nicht mehr unter Auflage bestimmter SLAs auf einzelnen Maschinen verarbeitet werden können“ (Vgl. [SW14]).

Als Unternehmen, das früh mit solchen Datenmengen konfrontiert war implementierte Google das Map-Reduce Paradigma² als Framework zur Ausnutzung vieler kostengünstiger Rechner für verschiedene Aufgaben (u.a. Indizierung von Webseiten und PageRank³).

In Folge der Veröffentlichung dieser Idee im Jahr 2004 wurde Map-Reduce in Form der OpenSource Implementation Hadoop (gemeinsam mit einer Implementation des Google File Systems GFS, u.a.)⁴ zum de-facto Standard für Big-Data-Analyseaufgaben.

Reines Map-Reduce (nach Art von Hadoop) als Programmierparadigma zur Verarbeitung großer Datenmengen zeigt jedoch in vielen Anwendungsfällen Schwächen:

¹[Lan01].

²[DG04].

³[Pag01].

⁴[SG03].

- Daten, die in hoher Frequenz entstehen und schnell verarbeitet werden sollen erfordern häufiges Neustarten von Map-Reduce-Jobs. Die Folge ist kostspieliger Overhead durch Verwaltung/Scheduling der Jobs und gegebenenfalls wiederholtem Einlesen von Daten.
- Algorithmen die während ihrer Ausführung iterativ Zwischenergebnisse erzeugen und auf vorherige angewiesen sind (häufig bei Maschinenlernalgorithmen) können nur durch persistentes Speichern der Daten und wiederholtes Einlesen zwischen allen Iterations-schritten implementiert werden.
- Anfragen an ein solches Map-Reduce-System erfolgen imperativ in Form von kleinen Programmen. Dieses Verfahren ist offensichtlich nicht so intuitiv und leicht erlernbar wie deklarative Abfragesprachen klassischer Datenbanken (z.B. SQL).

In der Folge dieser Probleme entstanden viele Ansätze dieses Paradigma zu ersetzen, zu ergänzen oder durch übergeordnete Ebenen und High-Level-APIs zu vereinfachen⁵.

Eine der Alternativen zu der Map-Reduce-Komponente in Hadoop ist die „general engine for large-scale data processing“ Apache Spark.

Ein Indiz für das steigende Interesse an diesem Produkt liefert unter anderem ein Vergleich des Interesses an Hadoop und Spark auf Google:

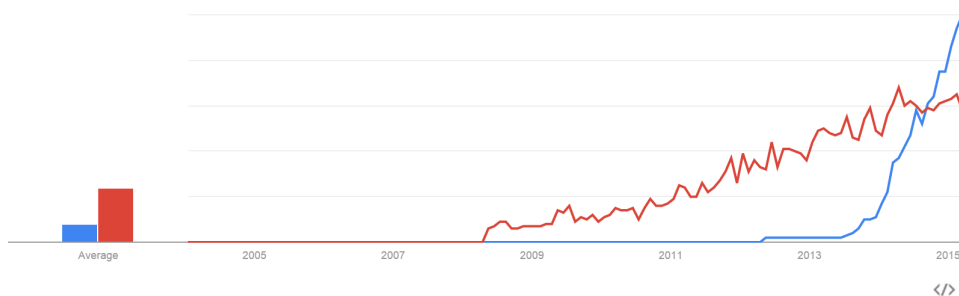


Abbildung 1.1: Vergleich der Suchanfragen zu Spark und Hadoop, Stand 24.03.2014 [Goo]

⁵[SR14].

1.2 Kontextabgrenzung

Das Ziel dieser Arbeit ist es eine Analyse und Bewertung der grundlegenden Konzepte und Anwendungsmöglichkeiten von Apache Spark zu vermitteln.

Für ein tieferes Verständnis werden Installation, Cluster-Betrieb und die Entwicklung von Treiberprogrammen beispielhaft durchgeführt, dokumentiert und bewertet. Hierbei kommt Apache Spark Version 1.3.0 zum Einsatz.

Nur am Rande wird betrachtet:

- Der Vergleich mit ähnlichen Produkten
- Die Empirische Messung des Skalierungsverhaltens
- Details zu Installation und Betrieb

2 Vorstellung von Apache Spark

Aus Sicht eines Nutzers ist Apache Spark eine API zum Zugriff auf Daten und deren Verarbeitung.

Diese API (wahlweise für die Sprachen Scala, Java und Python verfügbar), kann im einfachsten Fall über eine eigene Spark Konsole mit REPL-Pattern¹ verwendet werden.

Die Zählung von Wortvorkommen in einem Text - das „Hello World“ der Big Data Szene - lässt sich dort mit zwei Befehlen realisieren:

```
1 $ ./spark-shell
2 [...]
3   /  __/___  ____  ____/  /___
4   _\  \/_  _\  _\  ' /  ___/  ' /
5   /___/  .___/\_,_/_/_/  /_/\_ \   version 1.3.0
6   /_/_/
7 Using Scala version 2.10.4 (OpenJDK 64-Bit Server VM, Java 1.7.0_75)
8 Type in expressions to have them evaluated.
9 [...]
10 scala> val text = sc.textFile("../Heinrich Heine - Der Ex-Lebendige")
11 [...]
12 scala> :paste
13 text.flatMap(line => line.split(" "))
14 .map(word => (word, 1))
15 .reduceByKey(_ + _)
16 .collect()
17 [...]
18 res0: Array[(String, Int)] = Array((Tyrann,,1), (im,2), (Doch,1) ...)
```

Listing 2.1: Word Count in der Spark Konsole

Aus Sicht eines Administrators ist Apache Spark eine Applikation auf einem Rechencluster, die sich in der Anwendungsschicht befindet und charakteristische Anforderungen insbesondere

¹[MH99].

an Lokalität des Storages und die Netzwerkperformance stellt.

Was das konkret bedeutet, welche Mechanismen und Konzepte dahinterstehen und in welchem Ökosystem von Anwendungen sich Apache Spark bewegt wird im nächsten Abschnitt beleuchtet.

2.1 Überblick

— Was ist es, wofür ist es, wie funktioniert es, (JVM, Scala) —

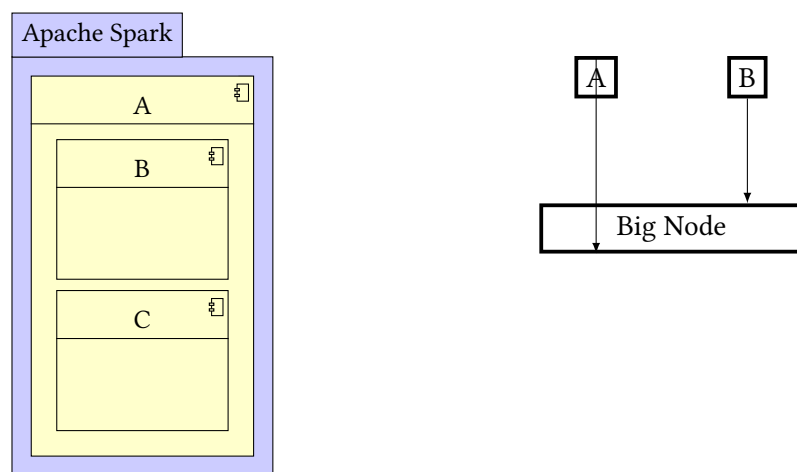


Abbildung 2.1: Verteilung der Software-Komponenten im Cluster

2.2 Kernkonzepte

— Warum ist Spark schnell, ausfallsicher, secure (und wo vielleicht nicht)? — — Annahme: Durchsatzfaktoren RAM, Netzwerk, Festplatte —

2.2.1 Nutzung von Arbeitsspeicher

2.2.2 Nutzung von persistentem Speicher

2.2.3 Nutzung von CPUs

2.2.4 Scheduling/Shuffling

2.2.5 Kern-API

2.3 Standardbibliotheken

— Warum ist Spark so einfach (und wo vielleicht nicht)? —

Die vier Standardbibliotheken erweitern die Kern-API für bestimmte, häufig genutzte Aufgaben aus Bereichen der Datenanalyse.

Die bedienten Bereiche sind

- Deklaratives Abfragen und On-the-Fly-Transformationen (*Spark SQL*)
- Maschinenlernverfahren (*MLlib*)
- Echtzeitbehandlung von eingehenden Daten (*Streaming*)
- Operationen auf Graph-artigen Strukturen (*GraphX*)

2.3.1 Dataframes/Spark SQL

2.3.2 MLlib

2.3.3 Streaming

2.3.4 GraphX

2.4 Betrieb und Security

2.5 Entwicklergemeinschaft

— Herkunft, Apache Foundation, Entwicklungsphilosophien, Anzahl Entwickler, ... —

2.6 Verwandte Produkte

— Ergänzende oder konkurrierende Produkte —

2.6.1 YARN

2.6.2 Mesos

2.6.3 Flink

3 Spark in der Praxis

Im Folgenden wird Apache Spark im Rahmen zweier grundsätzlich verschiedener Anwendungsfälle betrachtet.

Beispiel 1: Eine typische Anwendung mit verteiltem lokalem Storage (HDFS) und Spark als „Client“ eines bestehenden Yarn Clustermanagers. — Commodity Hardware (Raspberry Pi Cluster). —

Beispiel 2: Eine untypische Anwendung mit verteiltem entfernten Storage und dem Spark Standalone Clustermanager. — HPC Hardware („Thunder“ des Hamburger KlimaCampus). —

3.1 Echtzeitbewertung von Twitter-Accounts nach ihrer Relevanz für Spark

— Fusion von Tweets und Mailinglisten <https://spark.apache.org/docs/1.3.0/mllib-feature-extraction.html>
Implementation auf einem Raspberry Pi Cluster mit HDFS und Yarn Clustermanager —

3.1.1 Beschreibung des Problems

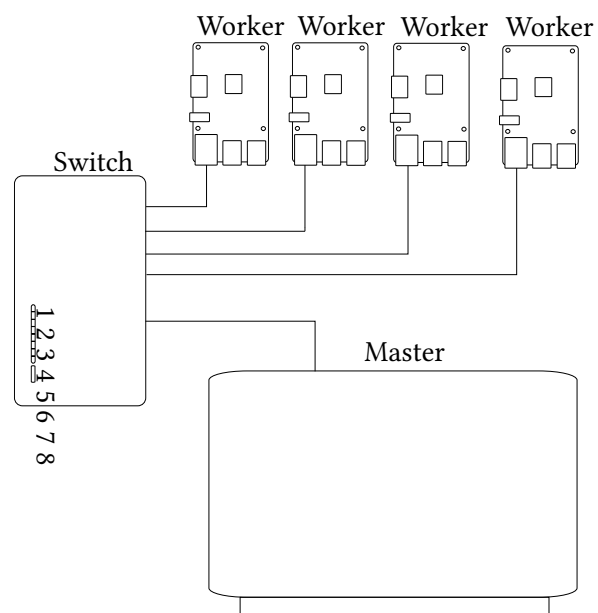
- Es sollen die beiden Spark Mailingslisten (Developer, User) zur Identifikation relevanter und aktueller Themen genutzt werden. Mit den so bewerteten Begriffen können wiederum Tweets bewertet werden. Mit den Tweets können dann ganze Accounts nach ihrer Relevanz beurteilt werden. —
- Zwei Datenquellen: Tweets (Nahe-Echtzeit), Entwickler-E-mails (Sporadisch) —
- Stichworte: HDFS, Yarn, Rasperry Pi Cluster, Machine Learning, Feature Extraction, Big Data Life Cycle —

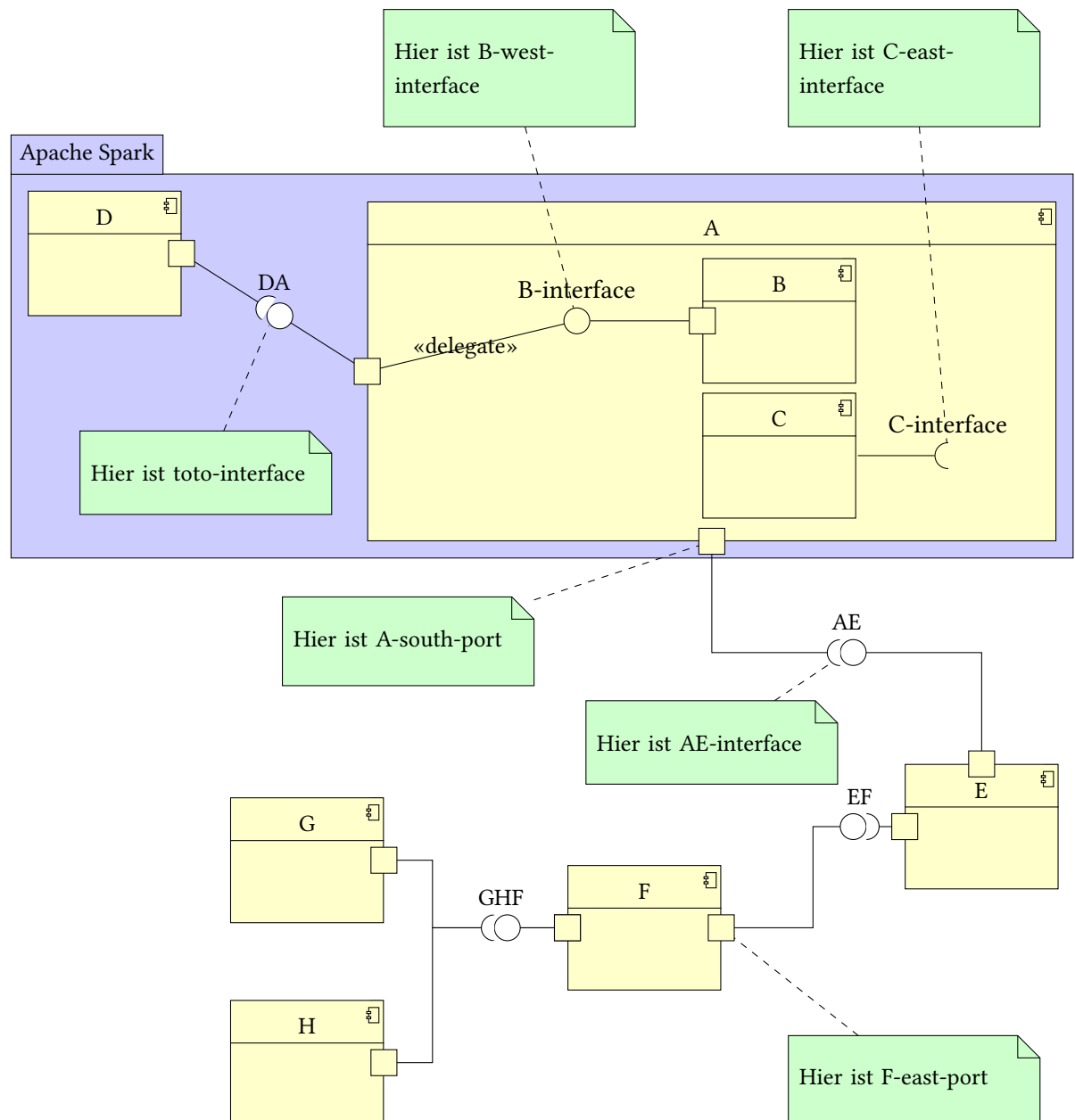
3.1.2 Hardwarekontext und Performance-Basisdaten

— hier kommen die eingesetzten systeme, und relevante laufzeitmessungen (netzwerk, storage, cpu) hin —

Tabelle 3.1: –DUMMY– Netzwerkdurchsatz

Nachrichtengröße	Worker → Worker	Master → Worker	Worker → Master
1kB	50ms	837ms	970ms
64kB	47ms	877ms	230ms
1MB	31ms	25ms	415ms
64MB	35ms	144ms	2356ms

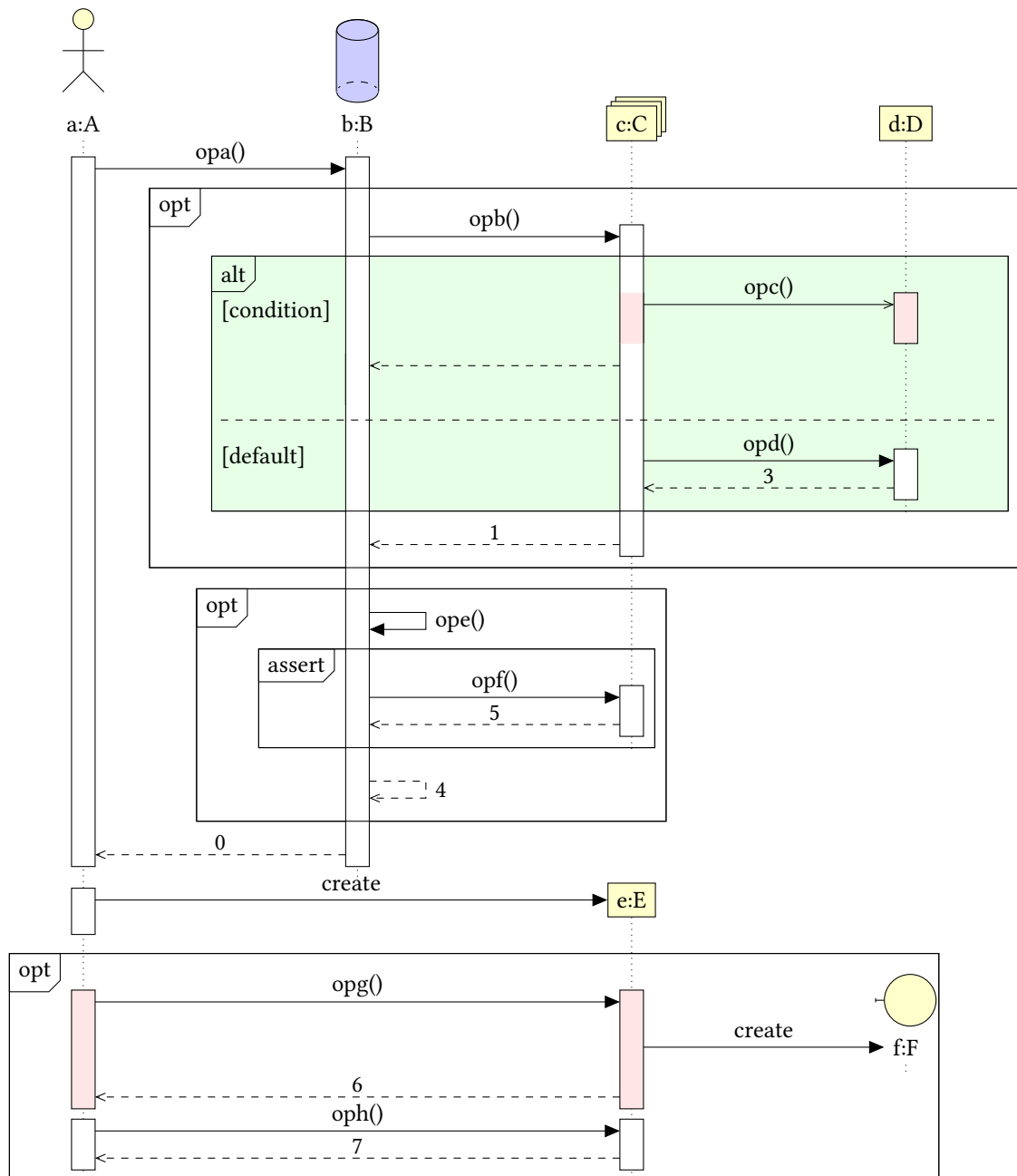




– hier kommen diagramme und codeschnipsel hin –

```
1 import org.apache.spark.SparkContext
2 import org.apache.spark.SparkContext._
3 import org.apache.spark.SparkConf
4
5 object ScalaApp {
6   val my_spark_home = "/home/daniel/projects/spark-1.1.0"
7
8   def main(args: Array[String]): Unit = {
9     val logFile = my_spark_home + "/README.md"
10
11     val conf = new SparkConf().setAppName("ScalaApp")
12     val sc = new SparkContext(conf)
13
14     val parList1 = sc.parallelize(List(1,2,3,4,5,6))
15     val parList2 = sc.parallelize(List(5,6,7,8,9,10))
16     val str1 =
17       "RDD1:_\%s".format(parList1.collect().deep.mkString("_"))
18     val str2 =
19       "RDD2:_\%s".format(parList2.collect().deep.mkString("_"))
20     val str3 =
21       "#_of_RDD1:_\%s".format(parList1.count())
22     val str4 =
23       "Intersect:_\%s".format(parList1.intersection(parList2).collect())
24     val str5 =
25       "Intersect:_\%s".format(parList1.cartesian(parList2).collect())
26   }
27 }
```

Listing 3.1: Treiber für Testanwendung (Programmiersprache Scala)



3.1.5 Ergebnisse

— Tabellen und Diagramme Ergebnissen, evt. Skalierungsverhalten — — Bewertung —

3.2 Evaluierung einer spark-basierten Implementation von CDOs auf einem HPC Cluster mit nicht-lokalem Storage

— Implementation ausgewählter CDOs (sehr wenige, möglicherweise nur 1-2) mit der Core-API von Spark. Testlauf auf einem HPC Cluster mit nicht-lokalem, allerdings per Infiniband angeschlossenen Storage. Insbesondere Betrachtung des Skalierungsverhaltens und der „Sinnhaftigkeit“. —

3.2.1 Beschreibung des Problems

— Erläuterung von CDOs (Climate Data Operators). —

3.2.2 Hardwarekontext und Performance-Basisdaten

— hier kommen die eingesetzten systeme, und relevante laufzeitmessungen (netzwerk, storage, cpu) hin —

3.2.3 Architekturübersicht

— hier kommen Verteilungs- und Komponentendiagramm hin —

3.2.4 Detaillierte Lösungsbeschreibung

— hier kommen laufzeitdiagramme und codeschnipsel hin —

3.2.5 Ergebnisse

— Tabellen und Diagramme Ergebnissen, evt. Skalierungsverhalten — — Bewertung —

4 Schlussbetrachtung

4.1 Kritische Würdigung der Ergebnisse

4.2 Ausblick und offene Punkte

Literatur

- [DG04] Jeffrey Dean und Sanjay Ghemawat. “MapReduce: Simplified Data Processing on Large Clusters”. In: *OSDI* (2004).
- [Goo] Google. *Google Trends*. URL: <https://www.google.com/trends>.
- [Lan01] Doug Laney. “3D Data Management: Controlling Data Volume, Velocity and Variety”. In: *Application Delivery Strategies* (2001).
- [MH99] Karl Knight Max Hailperin Barbara Kaiser. *Concrete Abstractions: An Introduction to Computer Science Using Scheme*. 1999, 278ff.
- [Pag01] Lawrence Page. *Method for node ranking in a linked database*. US Patent 6,285,999. 2001. URL: <http://www.google.com/patents/US6285999>.
- [SG03] Shun-Tak Leung Sanjay Ghemawat Howard Gobioff. *The Google File System*. Techn. Ber. Google, 2003.
- [SR14] Dilpreet Singh und Chandan K Reddy. “A survey on platforms for big data analytics”. In: (2014).
- [SW14] Jasson Venner Sameer Wadkar Madhu Siddalingaiah. *Pro Apache Hadoop*. 2014, S. 1.

Glossar

Request Evaluate Reply Loop sadadf dsfsdgm. 1

Hiermit versichere ich, dass ich die vorliegende Arbeit ohne fremde Hilfe selbständig verfasst und nur die angegebenen Hilfsmittel benutzt habe.

Hamburg, 1. Januar 2345 Daniel Kirchner
