

Universität Augsburg  
Fakultät für Angewandte Informatik

# Masterarbeit

im Studiengang Informatik

zur Erlangung des akademischen Grades  
Master of Science

<b>Thema:</b>	<Thema der Arbeit>
<b>Autor:</b>	Gerald Siegert MatNr. 1450117
<b>Version vom:</b>	15. Januar 2018
<b>Betreuer:</b>	M.Sc. Benedikt Eberhardinger
<b>1. Prüfer:</b>	Prof. Dr. X
<b>2. Prüfer:</b>	Prof. Dr. Y

# Zusammenfassung

## Abstract

# Inhaltsverzeichnis

Abbildungsverzeichnis	III
Listingverzeichnis	IV
Abkürzungsverzeichnis	V
<b>1 Einleitung</b>	<b>1</b>
1.1 Testen von Software und Systemen . . . . .	1
1.2 Problemstellung . . . . .	2
1.2.1 Aufbau des Modells . . . . .	3
1.2.2 Mapping zum realen System . . . . .	3
1.2.3 Erstellung der Lastprofile . . . . .	4
1.2.4 Erstellen und Ausführen der Tests . . . . .	4
1.2.5 Evaluierung der Ergebnisse . . . . .	4
<b>2 Relevante Technik</b>	<b>6</b>
2.1 Model Checking . . . . .	6
2.2 S# . . . . .	8
2.3 Apache Hadoop . . . . .	9
<b>3 Aufbau des Modells</b>	<b>12</b>
3.1 Grundlegende Architektur . . . . .	12
3.2 YARN-Modell . . . . .	13
3.3 SSH-Treiber . . . . .	15
3.3.1 Kommandozeilen-Parser . . . . .	15
3.3.2 REST-API-Parser . . . . .	15
3.3.3 Connector . . . . .	15
3.3.4 SSH-Verbindung . . . . .	15
3.4 Reales Hadoop-Cluster . . . . .	15
<b>Literaturverzeichnis</b>	<b>16</b>

# Abbildungsverzeichnis

1.1	Aufbau des V-Modells . . . . .	2
2.1	Schematischer Aufbau beim MC . . . . .	7
2.2	Architektur von YARN . . . . .	10
2.3	Architektur des HDFS . . . . .	11
3.1	Grundlegende Architektur des Gesamtmodells . . . . .	12
3.2	Aufbau des YARN-Modells . . . . .	14

# Listingverzeichnis

2.1	Grundlegender Aufbau einer S#-Komponente . . . . .	8
-----	--	---

# Abkürzungsverzeichnis

<b>MC</b>	Model Checker
<b>MC</b>	Model Checker
<b>DCCA</b>	Deductive Cause-Consequence Analysis
<b>RM</b>	ResourceManager
<b>NM</b>	NodeManager
<b>AM</b>	ApplicationManager
<b>AppMstr</b>	ApplicationMaster

# Kapitel 1

## Einleitung

### 1.1 Testen von Software und Systemen

Softwaretests sind in der heutigen Zeit eine wichtige Grundlage im Bereich der Qualitätssicherung bei Softwareprojekten. Abhängig von Komplexität und Sicherheitsanforderungen werden meist zwischen 30 und 60 Prozent der Kosten einer Software für die Qualitätssicherung und somit das Testen der Software verwendet (13). Ohne Softwaretests hätte heutige Software zahlreiche Fehler. Um die Wichtigkeit von Softwaretests zu unterstreichen sieht z.B. das V-Modell (vgl. Abbildung 1.1) für jeden Entwicklungsschritt am Ende eine entsprechende Testphase vor, woher auch der Name V-Modell stammt.

Nun ist es natürlich sehr aufwändig und daher mit einem vertretbaren Aufwand kaum machbar, jeden Test manuell durchzuführen. Daher wird versucht, möglichst viel zu automatisieren. Vor allem bei Unit-Tests ist dies mithilfe des *xUnit*-Framework (zu dem JUnit<sup>1</sup> für Java und NUnit<sup>2</sup> für .NET zählen) sehr einfach möglich. Dabei werden zunächst einzelne Testfälle erstellt und können im Anschluss aufgrund der aktuellen Codebasis jederzeit ausgeführt werden. Automatisierte Tests können auch dazu genutzt werden, um einen einzelnen Test mit verschiedenen Eingaben durchzuführen. Dadurch können verschiedene Eingabeklassen (wie negative oder positive Ganzzahlen) mit sehr geringem Aufwand in einem Test genutzt werden und somit verschiedene Testfälle direkt ausgeführt werden. Durch diese Testautomatisierung können somit zahlreiche Kosten eingespart werden (13).

Weitere Testframeworks gibt es auch zum Testen von sicherheitskritischen Systemen, welche nicht immer reine Software sein müssen. Dies ist einer der wichtigsten Anwendungsfälle für sogenannte Model Checker (MC). Model Checker (MC) wird dazu verwendet, um ein Modell auf seine Spezifikation zu testen. Dazu werden abhängig von den aktivierten Komponentenfehlern verschiedene Zustände des Modells vollautoma-

---

<sup>1</sup><https://junit.org>

<sup>2</sup><https://nunit.org/>



Abbildung 1.1: Aufbau des V-Modells (14)

tisch ausgeführt und anhand der dadurch im System auftretenden Fehler entschieden, ob die Spezifikation erfüllt wird (vgl. Abschnitt 2.1) (10, 11).

## 1.2 Problemstellung

Nun gibt es zahlreiche MC, wie z. B. *LTSmin*<sup>3</sup>. Auch am Institute for Software & Systems Engineering der Universität Augsburg wurde in den letzten Jahren mit S#<sup>4</sup> (sprich *Safety Sharp*) ein entsprechendes Framework zum testen von sicherheitskritischen und selbst-adaptiven Systemen basierend auf dem MC-Ansatz entwickelt (11, 12). Nun ist das aktuelle Vorhaben, mithilfe von S# ein reales Serversystem zu testen, welches bereits als theoretisches Modell basierend auf der ZNN.vom-Fallstudie, bekannt aus der Dissertation von Shang-Wen Cheng (7), implementiert wurde. Als reales System soll nun *Apache Hadoop*<sup>5</sup> getestet werden, welches in der Industrie im Bereich Datenverarbeitung eingesetzt wird. Mit Hadoop ist es möglich, ein Servercluster zu erstellen, auf dem anschließend dafür entwickelte Anwendungen ausgeführt werden und somit große Datenmengen zu verarbeiten. Es soll daher nun getestet und analysiert werden, wie sich Hadoop unter verschiedenen Lastprofilen verhält und dabei bestimmte Fehler auftreten, wenn z. B. einer der Hadoop-Nodes ausfällt.

Bei der ZNN.com-Fallstudie als reines Modell gab es bereits eine ähnliche Aufgabenstellung, welche im Positionspapier (8) genauer erläutert wurde. Das Hauptziel dieses Projektes ist daher nun, anstatt eines reinen Modells ein reales System zu testen. Dafür wird ein modellbasierter Ansatz als Testkonzept genutzt und Hadoop zunächst als Modell nachgebildet. Dieses Modell wird dann dazu genutzt, um ein reales Hadoop-Cluster entsprechend zu steuern und mithilfe des MC-Ansatzes zu testen, wie sich das

<sup>3</sup><http://fmt.cs.utwente.nl/tools/ltsmin/>

<sup>4</sup><https://safetysharp.isse.de>

<sup>5</sup><https://hadoop.apache.org/>



reale System unter bestimmten Bedingungen verhält und dabei zu ermitteln, wann es nicht mehr funktionsfähig ist.

### 1.2.1 Aufbau des Modells

Zunächst muss natürlich erst einmal der Versuchsaufbau selbst in S# modelliert werden. Ein Modell beinhaltet in S# zunächst einmal die Komponenten des Systems und deren Zusammenhänge, also wie die Komponenten miteinander agieren. Wichtig sind in einem S#-Modell aber auch mögliche Komponentenfehler, welche bekannt sind und jederzeit auftreten können. Komponentenfehler werden bereits in der Designphase eines Modells eingearbeitet und können bei der späteren Ausführung flexibel aktiviert und deaktiviert werden, um die Probleme des zu testenden Systems zu ermitteln (vgl. Abschnitt 2.2).

Um nun Hadoop in S# zu modellieren wird zunächst ein Konzept erstellt, in dem ausgearbeitet wird, welche Komponenten und Komponentenfehler relevant sind. Anschließend müssen deren Zusammenhänge und wesentlichen Eigenschaften ausgearbeitet werden und in das Konzept übernommen werden. Sobald das Konzept fertig ausgearbeitet ist, kann das Modell in S# implementiert werden.

### 1.2.2 Mapping zum realen System

Nachdem die Basis des Modells steht, kann die Funktionalität entwickelt werden. Dazu werden in S# nur Basisfunktionen eingebaut, um mit dem realen System kommunizieren zu können. Um für das reale System möglichst viele Ressourcen zur Verfügung zu stellen, wird das reale System auf einem eigenen PC installiert. Die Verbindung zwischen Modell und realem System wird mit einer Art Treiber realisiert, welcher mithilfe von mehreren SSH-Verbindungen mit dem realen Hadoop-Cluster kommuniziert und so das Mapping zwischen Modell und realem System übernimmt. Jede der SSH-Verbindungen ist dabei nur für einen Einsatzzweck gedacht, sodass es Verbindungen für u. A. folgende Einsatzzwecke gibt:

- Starten von Benchmark-Anwendungen
- Monitoring des realen Cluster
- Injektion von Komponentenfehler

Der Vorteil von mehreren Verbindungen liegt darin, dass jede Verbindung unabhängig ist und nicht auf die Antwort des zuvor ausgeführten Befehls einer anderen Verbindung warten muss. So ist es möglich, mithilfe mehrerer Verbindungen mehrere Anwendungen parallel zu starten und jede Rückgabe auszuwerten und währenddessen verschiedene Komponentenfehler zu aktivieren.

### 1.2.3 Erstellung der Lastprofile

Sobald das Grundmodell steht, können die Testfälle selbst entwickelt werden. Als Testfälle dienen dazu unterschiedliche Lastprofile, um verschiedene Auslastungsgrade und Nutzungsszenarien zu simulieren. Dazu sollen die Lastprofile verschiedene Benchmarks beinhalten, deren einzelne Anwendungen kombiniert oder alleine auf dem realen System ausgeführt werden:

- Hadoop Mapreduce Examples
- Intel HiBench
- SWIM (Statistical Workload Injector for Mapreduce)

Eine Besonderheit bildet der SWIM-Benchmark, welcher sehr Ressourcenintensiv ist und daher auf einem *Single Node Cluster*, also einem kompletten Hadoop-Cluster auf nur einem Computer, sehr zeitintensiv sein kann. Der Intel HiBench basiert auf einzelnen Bestandteilen der Mapreduce Examples. Die Examples wiederum sind zahlreiche voneinander unabhängige Beispielanwendungen für Hadoop. Dadurch besteht die Möglichkeit, abhängig davon, welche Anwendungen einzeln oder parallel gestartet werden, unterschiedliche Profile zu simulieren. Daher muss zunächst auch geprüft werden, welcher Benchmark welche Möglichkeiten bietet, um die benötigten Testfälle bzw. Lastprofile zu erstellen und so den dynamischen Teil des zu testenden Modells zu erstellen.

### 1.2.4 Erstellen und Ausführen der Tests

Sobald Modell und Testfälle stehen, kann mit der Erstellung der Tests fortgefahren werden. Die Tests müssen nun so erstellt werden, dass sie sich einerseits auf veränderte Bedingungen des realen Clusters anpassen, aber auch automatisiert die einzelnen Anwendungen der Lastprofile aktivieren und ausführen. Dies schließt auch unterschiedliche Profile für die Aktivierung der Komponentenfehler ein. Zum einen kann nur eine Simulation ohne Fehler gestartet werden, zum anderen aber auch unterschiedliche Komponentenfehler aktiviert werden. Der MC von S# besitzt dazu auch Möglichkeiten, um Komponentenfehler kombiniert auszuführen. Dazu werden basierend auf zuvor definierte *Constraints* Komponentenfehler aktiviert, um so typische Probleme des realen Systems zu simulieren. Basierend darauf wird dann ermittelt, welche Fehler nun im realen System auftauchen.

### 1.2.5 Evaluierung der Ergebnisse

Je nachdem welche *Constraints* bei der Ausführung genutzt werden, sind nun unterschiedliche Fehler und Daten im realen System ermittelt worden, welche zum Abschluss evaluiert werden müssen. Einige Erwartungen sind da natürlich bereits im Voraus klar:

Sollte es zu einem Netzwerk- oder Serverausfall eines Hadoop-Nodes kommen, muss das System dies selbstständig erkennen und die Anwendung an einen anderen Node abgeben. Dabei sollte das System auch erkennen, welche anderen Nodes bereits beschäftigt sind und entsprechend auf dem von Hadoop genutzten *Load Balancer* einen Node auswählen. Neben einer Fehleranalyse können aber auch die Laufzeiten unter bestimmten Bedingungen analysiert werden.

# Kapitel 2

## Relevante Technik

### 2.1 Model Checking

MC ist eine Möglichkeit, um Systeme zu testen und zu verifizieren. Dazu werden vom MC alle möglichen Systemzustände in einem *brute-force*-ähnlichem Vorgehen getestet und somit alle möglichen Szenarien getestet. Die Anzahl der Zustände kann sehr schnell  $10^{120}$  oder mehr betragen (10, 6).

Ein MC nutzt, wie der Name schon sagt, ein Modell des Systems, um das System zu testen. Wie bei jeder anderen modellbasierten Technik ist daher die Qualität des MC nur so gut wie das darauf zugrunde liegende Modell. Ein Modell kann auch als endlicher Automat angesehen werden, da ein Modell ebenfalls eine endliche Anzahl an möglichen Zuständen und dazugehörige Übergänge besitzt. Für jede Eigenschaft eines Zustandes muss zudem mithilfe einer sog. *temporalen Logik*, also mathematisch bzw. formal, festgelegt werden, was gültige Werte dieser Eigenschaft sind. Die dazu benötigten Informationen werden aus den Anforderungen des Systems ermittelt und dem MC übergeben. So können später verschiedene Eigenschaften des gesamten Systems (z. B. die formale Korrektheit, die Ausführbarkeit ohne Deadlocks oder die Einhaltung von Sicherheitsvorgaben) geprüft werden.

Zur Ausführung wird das gesamte Modell zunächst initialisiert und dann automatisch und systematisch vom MC auf Fehler und ungültige Zustände geprüft. In der Regel ist aber auch eine Ausführung als reine Simulation des Systems möglich, ohne explizit nach Fehlern zu suchen.

Wenn alle Zustände und deren Eigenschaften die Anforderungen erfüllen, erfüllt auch das Modell die Spezifikation. Wenn ein Zustand bzw. Eigenschaft die Anforderungen nicht erfüllt, prüft der MC anhand eines Gegenbeispiels den Ausführungspfad zum Fehler. Dadurch kann ermittelt werden, wo die Fehlerursache liegt. Einige der wesentlichen Fehlertypen und Ursachen sind:

**Modelling Error** Der Fehler liegt im Modell, welches korrigiert werden muss.

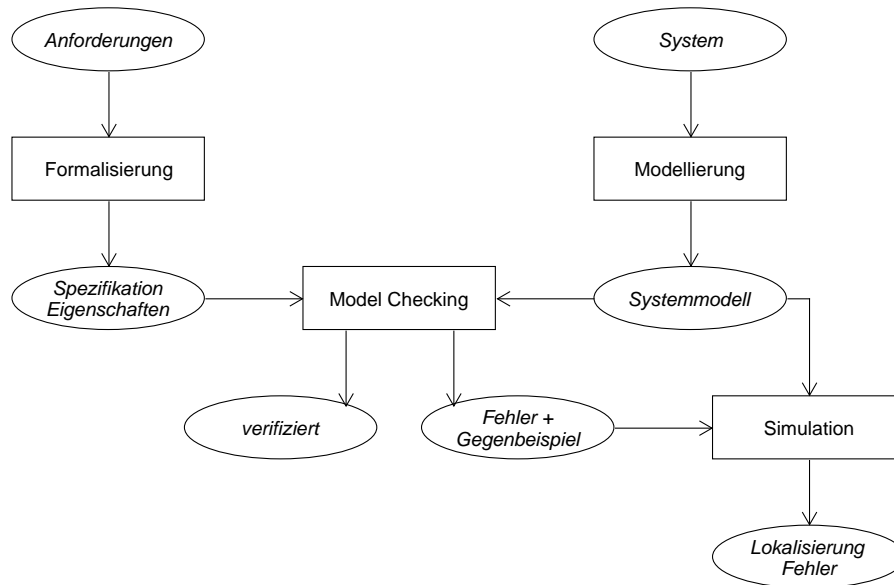


Abbildung 2.1: Schematischer Aufbau beim MC, nach (6)

**Design Error** Der Fehler liegt in den formellen oder informellen Anforderungen, dadurch muss das Modell und/oder die temporale Logik korrigiert werden.

**Property Error** Der Fehler ist wirklich ein Fehler im System, welcher gefunden werden soll.

Möglich ist aber auch, dass die Ressourcen nicht ausreichen, um alle Zustände zu prüfen. In so einem Fall gibt es mehrere Möglichkeiten, damit umzugehen, z. B. können Heuristiken oder Abstraktionen vom Modell genutzt werden (6, 9).

MC besitzt durch seine Charakteristik einige Vorteile, u. A. (6):

- MC ist universell nutzbar, z. B. für Software, Hardware oder eingebettete Systeme
- Partielle Verifikation ist möglich ohne das gesamte System testen zu müssen
- Vollständig automatisierbar und benötigt kaum Benutzerinteraktion oder hohe Expertise

Natürlich gibt es aber auch einige Nachteile, u. A. (6):

- Mit MC wird nur ein Systemmodell und nicht das eigentliche System getestet, was weitere Fehler nicht ausschließt
- Hauptsächlich für steuerungsbasierte Anwendungen und nicht für datenbasierte Anwendungen geeignet
- Anzahl der möglichen Zustände kann zu hoch sein, um alle zu testen

Es gibt zahlreiche MC-Frameworks, die bereits erwähnten *LTSmin* und *S#* sind nur zwei davon.

```
1 public class YarnNode : Component
2 {
3     // fault definition, also possible: new PermanentFault()
4     public readonly Fault NodeConnectionError = new TransientFault();
5
6     // interaction logic (Fields, Properties, Methods...)
7
8     // fault effect
9     [FaultEffect(Fault = nameof(NodeConnectionError))]
10    internal class NodeConnectionErrorEffect : YarnNode
11    {
12        // fault effect logic
13    }
14 }
```

Listing 2.1: Grundlegender Aufbau einer S#-Komponente

## 2.2 S#

Wie bereits erwähnt, ist S# das am Institute for Software & Systems Engineering der Universität Augsburg entwickelte MC-Framework. Da es in C# entwickelt wurde und C# auch zum Entwickeln von Modellen und dazugehörigen Testszenarien genutzt wird, können zahlreiche Features des .NET-Frameworks bzw. der Sprache C# im Speziellen genutzt werden. S# vereint dabei die Simulation, die Visualisierung, modellbasierte Tests sowie das MC der Modelle (11, 12). Dadurch können alle Schritte einer vollständigen Analyse inkl. Modellierung direkt im Visual Studio ausgeführt werden und somit auch alle Features der IDE und von .NET, wie z. B. die Debugging-Werkzeuge, genutzt werden. Um das MC durchzuführen, hat S# jedoch einige Einschränkungen, u. A. sind Schleifen und Rekursionen nur eingeschränkt bzw. nicht möglich. Die größte Einschränkung ist allerdings, dass während der Laufzeit keine neuen Objektinstanzen erzeugt werden können, sodass alle benötigten Instanzen bereits während der Initialisierung des Modells erzeugt werden müssen (11).

Um nun ein System testen zu können, muss dieses zunächst mithilfe von C#-Klassen und -Instanzen modelliert werden. Die dafür verwendeten Modelle sind meist stark vereinfacht und bilden nur die wesentlichen Aspekte der realen Systeme ab. Für einen korrekten Test ist es jedoch wichtig, dass das Modell des Systems vergleichbar mit dem echten System ist.

Listing 2.1 zeigt den typischen, grundlegenden Aufbau einer S#-Komponente. Jede Komponente des Modells muss von **Component** erben, um als S#-Komponente definiert zu sein. Jede Komponente kann nun temporäre (**TransientFault**) oder dauerhafte (**PermanentFault**) Komponentenfehler enthalten, welche zunächst innerhalb der Komponente als Member definiert werden. Der Effekt eines Komponentenfehlers wird anschließend in der entsprechenden Unterklasse definiert, welche von der Hauptklasse (hier **YarnNode**) erbt und mithilfe des Attributs **FaultEffect** dem dazugehörigen Komponentenfehler zugeordnet wird (12).

Um die Modelle zu testen, kommt in S# die Deductive Cause-Consequence Analysis (DCCA) zum Einsatz. Die DCCA ermöglicht eine vollautomatisch und MC-basierte Sicherheitsanalyse, wodurch selbstständig die Menge der aktivierten Komponentenfehler ermittelt wird, mit denen sich das Gesamtsystem nicht mehr rekonfigurieren kann und somit ausfällt. Je nach Konfiguration können dazu auch Heuristiken genutzt werden, welche die Analyse beschleunigen und genauer machen können (9). Dabei werden die verschiedenen aktivierten Komponentenfehler während der Analyse in tolerierbare und nicht-tolerierbare Fehler unterschieden. Tolerierbare Komponentenfehler werden dazu genutzt, die Grenzen der Selbstkonfiguration des Systems zu ermitteln. Dabei wird für jeden Systemzustand nach einer Rekonfiguration durch die DCCA eine neue Fehlermenge ermittelt, mit der das System gerade noch so lauffähig ist. Das Auftreten eines tolerierbaren Komponentenfehler ist also gleichbedeutend mit einem einfachen Fehler im System, welcher die gesamte Funktionsweise des Systems nicht massiv einschränkt und es sich noch selbst rekonfigurieren kann. Sobald jedoch ein Fehler auftritt, durch den es dem System nicht mehr möglich ist, sich selbst zu rekonfigurieren, wurde ein nicht-tolerierbarer Fehler gefunden, durch den das System nicht mehr funktionsfähig ist (11).

## 2.3 Apache Hadoop

Apache Hadoop ist ein Open-Source-Software-Projekt, mit dessen Hilfe ermöglicht wird, Programme zur Datenverarbeitung mit großen Ressourcenbedarf auf verteilten System auszuführen. Hadoop wird von der *Apache Foundation* entwickelt und bietet verschiedene Komponenten an, welche vollständig skalierbar sind, von einer einfachen Installation auf einem PC bis hin zu einer Installation über mehrere Server in einem Serverzentrum. Hadoop besteht hauptsächlich aus folgenden Kernmodulen (5):

**Hadoop Common** Gemeinsam genutzte Kernkomponenten

**Hadoop YARN** Framework zur Verteilung und Ausführung von Tasks und das dazugehörige Ressourcen-Management

**Hadoop Distributed File System** Kurz HDFS, Verteiltes Dateisystem

**Hadoop MapReduce** YARN-Basiertes System zum Verarbeiten von großen Datenmengen

Hadoop ermöglicht es dadurch, sehr einfach mit Tasks umzugehen, welche große Datenmengen verarbeiten. Da es für Hadoop nicht relevant ist, auf wie vielen Servern es läuft, kann es beliebig skaliert werden, wodurch entsprechend viele Ressourcen zur Bearbeitung und Speicherung von großen Datenmengen zur Verfügung stehen können.



Abbildung 2.2: Architektur von YARN (1)

Die Kernidee der Architektur von **YARN** ist die Trennung vom Ressourcenmanagement und Scheduling. Dazu besitzt der Master bzw. **Controller** den **ResourceManager** (RM), welcher für das gesamte System zuständig ist und die Anwendungen im System überwacht. Er besteht aus zwei Kernkomponenten, dem **ApplicationManager** (AM) und dem *Scheduler*. Der AM ist für die Annahme und Ausführung von einzelnen Anwendungen zuständig, denen der Scheduler die dafür notwendigen Ressourcen zuteilt und überwacht. Für jeden *Slave-Node* im Hadoop-System gibt es dazu einen **NodeManager** (NM), welcher die lokalen Ressourcen des Nodes überwacht und dem RM mitteilt. Jede Anwendung besitzt jeweils einen eigenen **ApplicationMaster** (AppMstr), welcher für das Monitoring und die Kommunikation mit dem RM und NM zuständig ist und die dazu notwendigen Informationen bereit stellt. Jede YARN-Anwendung bzw. Job besteht aus einem oder mehreren *Containern*, in denen die einzelnen Tasks der Anwendung auf einem beliebigen Node ausgeführt werden (1).

Das **HDFS** basiert auf der gleichen Architektur wie YARN und besitzt ebenfalls einen Master und mehrere Slaves, welches in der Regel die gleichen Nodes sind wie bei YARN sind. Der *NameNode* ist als Master für die Verwaltung des Dateisystems zuständig und reguliert den Zugriff auf die darauf gespeicherten Daten. Die Daten selbst werden in mehrere Blöcke aufgeteilt auf den *DataNodes* gespeichert. Um den Zugriff auf die Daten im Falle eines Node-Ausfalls zu gewährleisten, wird jeder Block auf anderen Nodes repliziert. Dateioperationen (wie Öffnen oder Schließen) werden direkt auf den *DataNodes* ausgeführt, sie sind darüber hinaus auch dafür verantwortlich, dass Clients die Daten lesen oder beschreiben können (2).





Abbildung 2.3: Architektur des HDFS (2)

**MapReduce** bietet analog zu YARN die Möglichkeit, Anwendungen mit einem großen Ressourcenbedarf, welche große Datenmengen verarbeiten, auf einem gesamten Cluster auszuführen. Dazu werden bei einem MapReduce-Job die Eingabedaten aufgeteilt, anschließend von den sog. *Map Tasks* verarbeitet und deren Ausgaben von den sog. *Reduce Tasks* geordnet. Für die Ein- und Ausgabe der Daten wird in der Regel das HDFS, für die Ausführung der einzelnen Tasks YARN genutzt (4). MapReduce kann man auch als Vorgänger von YARN ansehen, da YARN auch als *MapReduce Next Gen* oder *MRv2* bezeichnet wird und erst in einer späteren Version von Hadoop implementiert wurde (3).

# Kapitel 3

## Aufbau des Modells

### 3.1 Grundlegende Architektur

Die grundlegende Architektur des gesamten Aufbaus besteht aus den drei rechts abgebildeten Schichten. Die oberste Schicht bildet das S#-Modell von Hadoop YARN, welches die relevanten YARN-Komponenten und Komponentenfehler abbildet. Das reale Pendant dazu bildet das reale Hadoop-Cluster auf einem eigenen PC als unterste Schicht. Die Verbindung zwischen Modell und realem Cluster bildet der Treiber als eigenständige Schicht. Der Treiber besteht aus folgenden Komponenten:

**Parser** Verarbeitet die Monitoring-Ausgaben vom realen Cluster und konvertiert diese für die Nutzung im Modell.

**Connector** Abstrahiert die SSH-Verbindung

**SSH-Verbindung** Eigentliche Verbindung zum PC mit dem realen Cluster

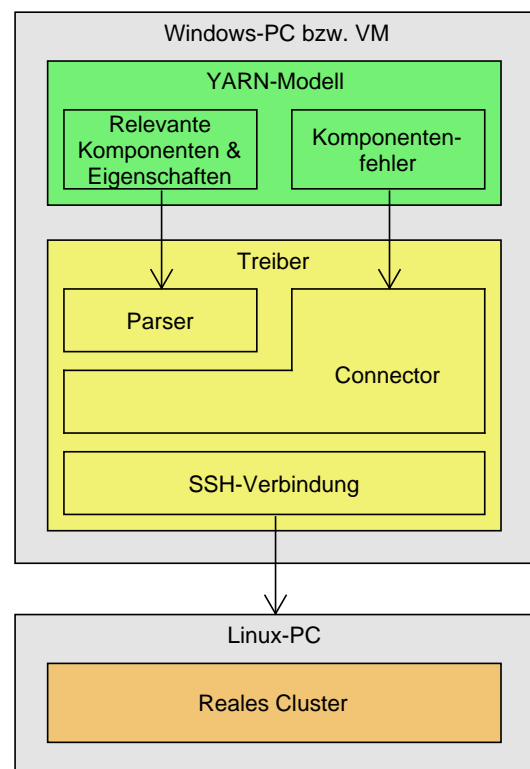


Abbildung 3.1: Grundlegende Architektur des Gesamtmodells

Auf den Treiber bzw. das reale System wird meist über den Parser zugegriffen. Lediglich zum Starten von Anwendungen, zum Aktivieren bzw. Deaktivieren von Komponentenfehlern u. Ä. auf dem realen Cluster wird direkt auf den Connector zugegriffen.

## 3.2 YARN-Modell

Abbildung 3.2 beschreibt im Grunde bereits das gesamte von S# verwendete YARN-Modell. Enthalten sind alle hier relevanten Komponenten sowie deren Eigenschaften. Als Eigenschaften wurden jeweils alle Daten aufgenommen, welche mithilfe von Shell-Kommandos bzw. mithilfe der REST-API von YARN ermittelt werden kann.

Die abstrakte Basisklasse **YarnHost** stellt die Basis für alle Hosts des Clusters dar, also dem **YarnController** mit dem RM, und dem **YarnNode**, was einen Node darstellt, auf dem die Anwendungen bzw. deren Container ausgeführt werden. Die abstrakte Eigenschaft **YarnHost.HttpPort** dient als Hilfs-Eigenschaft, da Controller und Nodes unterschiedliche Ports für die Weboberfläche nutzen, deren URL mit Port in der Eigenschaft **YarnHost.HttpUrl** abrufbar ist. Sie wird daher vom Controller bzw. Node mit dem entsprechenden Port versehen. Die Felder **YarnNode.NodeConnectionError** und **YarnNode.NodeDead** bilden die Komponentenfehler, wenn ein Node seine Netzwerkverbindung verliert bzw. beendet wird. Die Effekte der Komponentenfehler werden über entsprechende innere Klassen realisiert.

Die Anwendungen, welche mit **YarnApp** dargestellt werden, werden mithilfe des **Client** gestartet, was den zu startenden Client darstellt. Die Anwendungen selbst enthalten neben grundlegenden Daten wie z. B. den Namen auch einige Daten zum Ressourcenbedarf (Speicher und CPU). Zwar gibt Hadoop nicht direkt die zu der Anwendung gehörigen Job-Ausführungen an, allerdings können diese mithilfe der **YarnApp.AppId** sehr einfach ermittelt werden und dann in der Liste **YarnApp.Attempts** gespeichert werden. Das Feld **YarnApp.IsKillable** gibt an, ob die Ausführung der Anwendung mit den aktuellen Daten im Modell durch den Komponentenfehler **YarnApp.KillApp** abgebrochen werden kann. Abhängig ist das durch **YarnApp.FinalStatus**, was angibt, ob eine Anwendung erfolgreich ausgeführt wurde oder die Ausführung noch nicht abgeschlossen ist (durch **EFinalStatus.UNDEFINED**).

Jede Ausführung **YarnAppAttempt** hat eine eigene ID und kann einer Anwendung zugeordnet werden. Genau wie bei den Anwendungen selber wird hier direkt der Node gespeichert, auf welchem der AppMstr ausgeführt wird und einen eigenen Container bildet, dessen ID direkt gespeichert wird. Container (dargestellt durch **YarnAppContainer**) existieren in Hadoop nur während der Laufzeit eines Programmes und enthalten nur wenige Daten, darunter ihr ausführender Node. Jede Anwendung, deren Ausführungen und deren Container enthalten zudem den derzeitigen Status, ob die Komponente noch eingereicht wird, bereits ausgeführt wird oder beendet ist. **EAppState.NotStartedYet** dient als Status, den es nur im Modell gibt und angibt, dass die Anwendung im späteren Verlauf der Testausführung gestartet wird.

Alle vier YARN-Kernkomponenten implementieren das Interface **IYarnReadable**, was angibt, dass die Komponente ihren Status aus Hadoop ermitteln kann. Entsprechend wird in allen Komponenten die Methode **GetStatus()** implementiert, in welchem

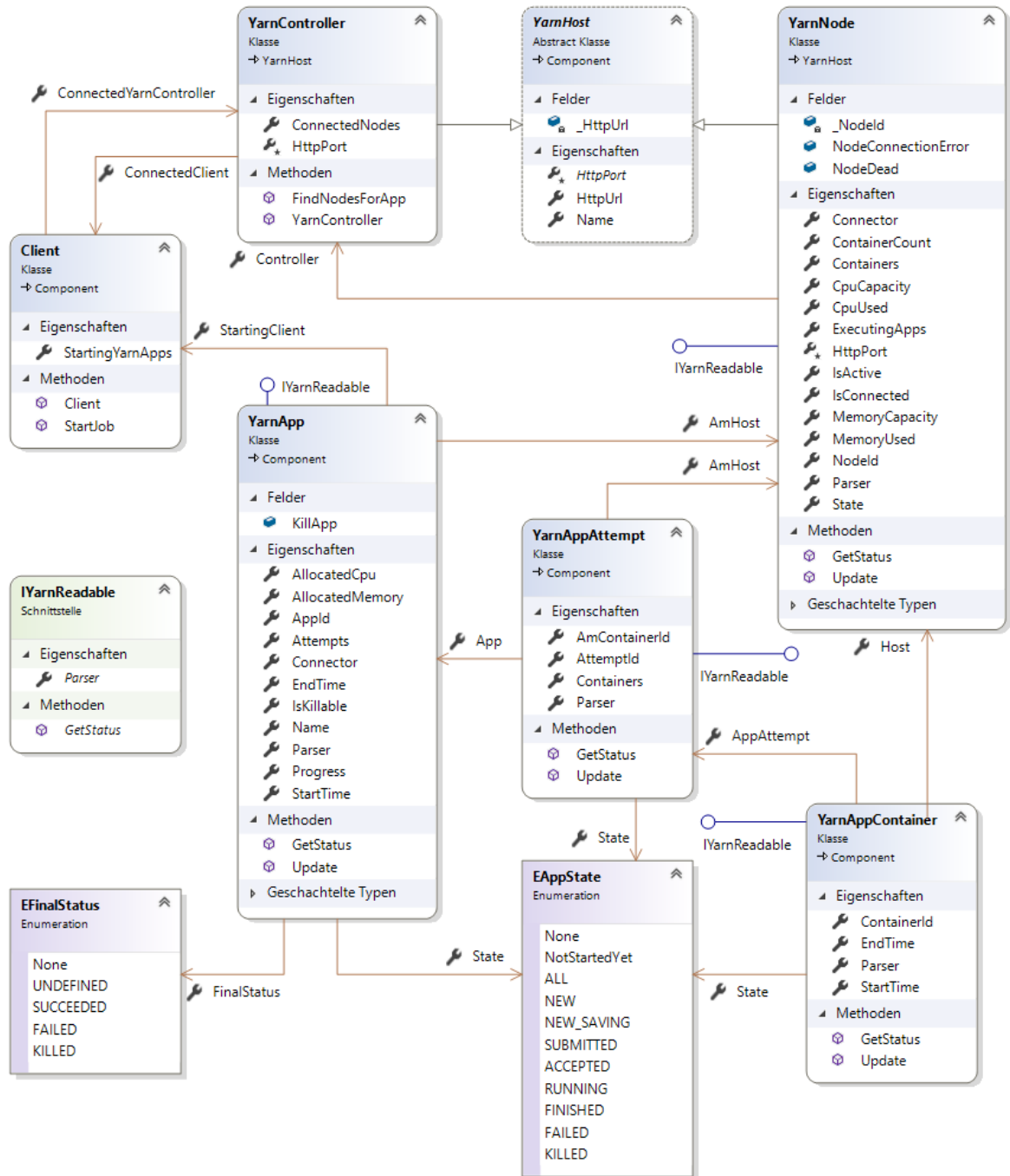


Abbildung 3.2: Aufbau des YARN-Modells. Das Modell wurde mithilfe des Klassendiagramm-Designers in Visual Studio 2017 erstellt. Daher werden Assoziationen aus Listen (wie `YarnApp.Attempts`) im Diagramm nicht angezeigt.

mithilfe des angegebenen Parsers auf den SSH-Treiber zugegriffen werden kann und die Komponenten im Modell so ihre Daten aus dem realen Cluster ermitteln können.

## **3.3 SSH-Treiber**

In Abschnitt 3.1 wurde bereits auf den grundlegenden Aufbau des Treibers eingegangen. Der SSH-Treiber besteht aus drei einzelnen Komponenten, welche mithilfe von Interfaces im YARN-Modell eingebunden sind. Dadurch ist es möglich, unterschiedliche Parser bzw. auch Verbindungen für unterschiedliche Komponenten zu nutzen. Da es zwei Möglichkeiten gibt, die Daten des realen Clusters ausgeben zu lassen, wurde dies auch genutzt und so ein Parser für das Monitoring mittels Kommandozeilen-Befehle, und ein Parser für die Nutzung der REST-API von Hadoop erstellt.

### **3.3.1 Kommandozeilen-Parser**

### **3.3.2 REST-API-Parser**

### **3.3.3 Connector**

### **3.3.4 SSH-Verbindung**

## **3.4 Reales Hadoop-Cluster**

# Literaturverzeichnis

- [1] APACHE SOFTWARE FOUNDATION: *Apache Hadoop NextGen MapReduce (YARN)*. <https://hadoop.apache.org/docs/r2.7.1/hadoop-yarn/hadoop-yarn-site/YARN.html>
- [2] APACHE SOFTWARE FOUNDATION: *HDFS Architecture*. <http://hadoop.apache.org/docs/r2.7.1/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html>
- [3] APACHE SOFTWARE FOUNDATION: *MapReduce NextGen aka YARN aka MRv2*. <http://hadoop.apache.org/docs/r2.7.1/hadoop-yarn/hadoop-yarn-site/index.html>
- [4] APACHE SOFTWARE FOUNDATION: *MapReduce Tutorial*. <http://hadoop.apache.org/docs/r2.7.1/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>
- [5] APACHE SOFTWARE FOUNDATION: *Welcome to Apache<sup>TM</sup>Hadoop®!* <https://hadoop.apache.org/>
- [6] BAIER, J.-P. Christel; K. Christel; Katoen: *Principles of model checking*. MIT Press, 2008. – ISBN 978–0–262–02649
- [7] CHENG, S.-W. : *Rainbow: Cost-effective Software Architecture-based Self-adaptation*. Pittsburgh, PA, USA, Carnegie Mellon University, Diss., 2008. – AAI3305807
- [8] EBERHARDINGER, B. ; HABERMAIER, A. ; REIF, W. : Toward Adaptive, Self-Aware Test Automation. In: *2017 IEEE/ACM 12th International Workshop on Automation of Software Testing (AST)*, 2017, S. 34–37
- [9] EBERHARDINGER, B. ; HABERMAIER, A. ; SEEBACH, H. ; REIF, W. : Back-to-Back Testing of Self-organization Mechanisms. In: WOTAWA, F. (Hrsg.) ; NICA, M. (Hrsg.) ; KUSHIK, N. (Hrsg.): *Testing Software and Systems: 28th IFIP WG 6.1 International Conference, ICTSS 2016, Graz, Austria, October 17-19, 2016, Proceedings*. Springer International Publishing. – ISBN 978–3–319–47443–4, 18–35
- [10] GRUMBERG, O. ; CLARKE, E. ; PELED, D. : *Model checking*. (1999)
- [11] HABERMAIER, A. ; EBERHARDINGER, B. ; SEEBACH, H. ; LEUPOLZ, J. ; REIF, W. : Runtime Model-Based Safety Analysis of Self-Organizing Systems with S#. In: *2015 IEEE International Conference on Self-Adaptive and Self-Organizing Systems Workshops*, 2015, S. 128–133

- [12] HABERMAIER, A. ; LEUPOLZ, J. ; REIF, W. : Unified Simulation, Visualization, and Formal Analysis of Safety-Critical Systems with S#. In: BEEK, M. H. (Hrsg.) ; GNESI, S. (Hrsg.) ; KNAPP, A. (Hrsg.): *Critical Systems: Formal Methods and Automated Verification: Joint 21st International Workshop on Formal Methods for Industrial Critical Systems and 16th International Workshop on Automated Verification of Critical Systems, FMICS-AVoCS 2016, Pisa, Italy, September 26-28, 2016, Proceedings*. Springer International Publishing. – ISBN 978-3-319-45943-1, 150–167
- [13] POLO, M. ; REALES, P. ; PIATTINI, M. ; EBERT, C. : Test Automation. In: *IEEE Software* 30 (2013), Jan, Nr. 1, S. 84–89. <http://dx.doi.org/10.1109/MS.2013.15>. – DOI 10.1109/MS.2013.15. – ISSN 0740-7459
- [14] PÄTZOLD, M. ; SEYFERT, S. : *V-Modell*. <https://commons.wikimedia.org/wiki/File:V-Modell.svg>. Version: Januar 2010. – Lizenz: CC-BY-SA 3.0