

MASTER THESIS  
COMPUTER SCIENCE

# Scaling UIMA

*Simon Gehring*

Am Jesuitenhof 3  
53117 Bonn  
simon.gehring@fkie.fraunhofer.de  
Matriculation Number 2553262

At the  
RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT BONN  
in cooperation with the  
FRAUNHOFER-INSTITUT FÜR KOMMUNIKATION,  
INFORMATIONSVERRARBEITUNG UND ERGONOMIE

supervised by  
Prof. Dr. Heiko RÖGLIN and 

February 28, 2018

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Motivation . . . . .	4
1.2	Basics . . . . .	4
1.2.1	UIMA . . . . .	4
1.2.2	Docker . . . . .	4
1.2.3	Hadoop . . . . .	5
1.2.4	Spark . . . . .	5
1.2.5	Kafka . . . . .	5
1.3	Problem . . . . .	5
1.3.1	Scaling UIMA . . . . .	5
1.3.2	Implementation Requirements . . . . .	5
1.4	Related Work . . . . .	5
1.4.1	Watson . . . . .	5
1.4.2	Something else that warrants another subsection . . . . .	5
1.4.3	GATE? . . . . .	5
1.5	Outline . . . . .	5
<b>2</b>	<b>Implementation</b>	<b>6</b>
2.1	Concrete Application . . . . .	6
2.2	Documents -> Kafka . . . . .	6
2.3	Kafka -> Spark . . . . .	6
2.4	Spark -> UIMA . . . . .	6
2.5	UIMA -> Java . . . . .	6
2.6	Kafka -> Output . . . . .	6
2.7	Bottlenecks . . . . .	6
<b>3</b>	<b>Evaluation</b>	<b>7</b>
3.1	Computation Speed . . . . .	7
3.2	Memory Usage . . . . .	7
3.3	Extensibility . . . . .	7
3.4	Maintainability . . . . .	7

<b>4</b>	<b>Summary</b>	<b>8</b>
4.1	The Judgement . . . . .	8
<b>5</b>	<b>Future Work</b>	<b>9</b>

# Chapter 1

## Introduction

Natural language is part of everyone’s everyday life and is most commonly used to transmit information human-to-human. While most of this interaction takes place orally or written on paper, the digital revolution and the rise of social media increased the amount of digitally stored natural language tremendously. Gantz and Reinsel predicted 2012 that the amount of digital data stored globally will double about every two years until at least the year 2020 [GR12].

Many opportunities arise from this amount of digital data, specifically in the field of machine learning. In 2011, IBM’s QA (Question Answering) system “Watson” famously outmatched professional players in the quiz show “Jeopardy!” [Fer12, ESI<sup>+</sup>12]. Kudesia et al. proposed 2012 an algorithm to detect so called CAUTIs<sup>1</sup>, common hospital-acquired infections, by utilizing a NLP (Natural Language Processing) analysis on the medical records of patients [KSDG12].

Apache UIMA (Unstructured Information Management Architecture) is one of few general approaches to implement NLP solutions. With a very modular architecture, UIMA is a popular tool that can easily be applied to a majority of NLP problems. A large part of the popularity of UIMA stems from the large DKPro Core (Darmstadt Knowledge Processing Software Repository) collection of components, containing hundreds of analysis modules and precomputed language models [EdCG14], which are easily imported into existing Java projects with the build automation tool Apache Maven [DKP].

A common problem with UIMA in non-academic environments is scaling [DCR<sup>+</sup>15, ESI<sup>+</sup>12, RBB<sup>+</sup>10]. UIMA itself provides two distinct interfaces to analyse larger collections of unstructured data, with one being UIMA-AS (UIMA Asynchronous Scaleout) and the other being the more dated and less flexible CPE (Collection Processing Engine) [FLVN09].

In this thesis, we will evaluate different means of scaling UIMA, using modern technologies like Docker, a container virtualization solution, Apache Spark, a cluster computing framework, and Apache Kafka, an information stream processing software. We will compare said implementations with the native UIMA-AS and CPE approach in terms of

---

<sup>1</sup>Catheter-associated Urinary Tract Infections

processor and memory efficiency, ease of implementation and maintainability.

## 1.1 Motivation

## 1.2 Basics

In this chapter, we will cover the basics for the necessary technologies used throughout the evaluation. All of these are concrete implementations of more general concepts and may be exchanged for similar products. However, the following products were chosen, mainly because they are Open Source<sup>1,2,3,4,5</sup> but also because of their popularity and relevance in the industry.

### 1.2.1 UIMA

UIMA is a data mining and NLP framework, created in 2005 by IBM [FL04] and maintained by Apache since 2006 [Apac]. It is available in Java and C++ and contains various scale-out options. Strictly speaking, one has to differentiate between the UIMA specification and Apache UIMA, an open source implementation of said specification [FLVN09]. Since both terms are often used interchangeably, this thesis will also not differentiate between Apache UIMA and UIMA. Unless specified else, we will always reference the Apache UIMA implementation.

The UIMA specification describes NLP applications as a collection of (mostly) independent components. Such a component is called an AE (Analysis Engine) and enriches a given document by inferred information. To modularize said components, UIMA provides the notion of a CAS (Common Analysis System). A CAS is an object containing the Sofa (Subject of Analysis), the analysis results and the used type system.

Alles ein wenig  
wirr. Braucht  
rewrite. Weiß  
aber nicht genau  
wie tief in die  
Materie ich hier  
gehen muss.

### 1.2.2 Docker

Docker is a software to create and run applications in virtual environments, without the need to setup a complete virtual operating system for each application. Based on Cgroups and namespaces, Docker runs applications in a containerized environment. To create a containerized application, a special markup file, called a Dockerfile has to be written by the programmer. The Dockerfile defines exactly how the applications environment needs to look like.

Source?

Listing 1.1: hiho

```
1 FROM stain/jena-fuseki:latest
2
```

TODO: Besseres  
Dockerfile finden.

<sup>1</sup><https://svn.apache.org/viewvc/uima/>, last accessed on 2018-02-27.

<sup>2</sup><https://github.com/docker>, last accessed on 2018-02-27.

<sup>3</sup><https://github.com/apache/hadoop>, last accessed on 2018-02-27.

<sup>4</sup><https://github.com/apache/spark>, last accessed on 2018-02-27.

<sup>5</sup><https://github.com/apache/kafka>, last accessed on 2018-02-27.

```
3 COPY ./fuseki /fuseki
4
5 COPY ./docker-entrypoint.sh /docker-entrypoint2.sh
6 COPY ./ReAn.ttl /staging/ReAn.ttl
7 COPY ./inferred.ttl /staging/inferred.ttl
8
9 ENTRYPOINT ["/docker-entrypoint2.sh"]
```

### 1.2.3 Hadoop

### 1.2.4 Spark

### 1.2.5 Kafka

## 1.3 Problem

### 1.3.1 Scaling UIMA

UIMA-AS

UIMA-CPM

### 1.3.2 Implementation Requirements

## 1.4 Related Work

### 1.4.1 Watson

### 1.4.2 Something else that warrants another subsection

### 1.4.3 GATE?

## 1.5 Outline

## Chapter 2

# Implementation

**2.1 Concrete Application**

**2.2 Documents -> Kafka**

**2.3 Kafka -> Spark**

**2.4 Spark -> UIMA**

**2.5 UIMA -> Java**

**2.6 Kafka -> Output**

**2.7 Bottlenecks**

## Chapter 3

# Evaluation

3.1 Computation Speed

3.2 Memory Usage

3.3 Extensibility

3.4 Maintainability



## Chapter 4

# Summary

### 4.1 The Judgement

## Chapter 5

# Future Work

# Glossary

## **Analysis Engine**

An Analysis Engine is a component of an UIMA NLP pipeline. As such it analyses given documents and enriches it with inferred information. An Analysis Engine may contain other Analysis Engines. 4

## **Apache Kafka**

Apache Kafka is an open-source stream processing software platform developed by the Apache Software Foundation written in Scala and Java. The project aims to provide a unified, high-throughput, low-latency platform for handling real-time data feeds. 3

## **Apache Spark**

Apache Spark is an open-source cluster-computing framework. Originally developed at the University of California, Berkeley's AMPLab, the Spark codebase was later donated to the Apache Software Foundation, which has maintained it since. Spark provides an interface for programming entire clusters with implicit data parallelism and fault tolerance. 3

## **Collection Processing Engine**

Collection Processing Engines (CPE) are the first generation of UIMA native scaling solutions. A CPE contains a collection reader, which knows how to read the underlying collection, and CAS Consumers for the final analysis result extraction [Apac]. 3

## **Common Analysis System**

The Common Analysis System is a type of object in the UIMA framework. It contains the subject of analysis, the analysis result and a corresponding type system. 4

## **Darmstadt Knowledge Processing Software Repository**

A collection of UIMA components for natural language processing. This includes analysis engines, language models and custom type systems [DKP, EdCG14]. 3

**Docker**

Docker is a virtualization solution based on containers. By using containers instead of fully fledged virtual machines Docker tries to reduce the system overhead per running application [doc15]. 3, 4

**Natural Language Processing**

Natural-language processing (NLP) is the discipline of collecting and analysing natural language. This includes for example speech recognition, natural language understanding and generation [Lid01]. 3

**Question Answering**

Being a subfield of NLP, Question Answering (QA) is about extracting and understanding questions from natural language and answering them accordingly [JM14]. 3

**Subject of Analysis**

The Subject of Analysis is the document that gets analyzed by a given UIMA application. It is contained in its corresponding CAS. 4

**UIMA Asynchronous Scaleout**

UIMA-AS is the second generation of UIMA native scaling solutions. It is based on a shared queue based service architecture [Apab] 3

**Unstructured Information Management Architecture**

UIMA is a general purpose framework to extract information from unstructured data [Apaa, FLVN09]. Although any data format is supported, natural language texts are the most common one. 3

# Bibliography

- [Apaa] The Apache Software Foundation. Apache uima – apache uima. <https://uima.apache.org/>. Accessed: 2018-02-26.
- [Apab] The Apache Software Foundation. Getting started: Apache uima asynchronous scaleout. <https://uima.apache.org/doc-uimaas-what.html>. Accessed: 2018-02-26.
- [Apac] The Apache Software Foundation. Uima tutorial and developers’ guides. [https://uima.apache.org/d/uimaj-2.4.0/tutorials\\_and\\_users\\_guides.html](https://uima.apache.org/d/uimaj-2.4.0/tutorials_and_users_guides.html). Accessed: 2018-02-26.
- [DCR<sup>+</sup>15] Guy Divita, M Carter, A Redd, Q Zeng, K Gupta, B Trautner, M Samore, and A Gundlapalli. Scaling-up nlp pipelines to process large corpora of clinical notes. *Methods of information in medicine*, 54(06):548–552, 2015.
- [DKP] The DKPro Core Team. Dkpro core<sup>TM</sup> user guide. [https://zoidberg.ukp.informatik.tu-darmstadt.de/jenkins/job/DKProCoreDocumentation\(GitHub\)/de.tudarmstadt.ukp.dkpro.core\\$de.tudarmstadt.ukp.dkpro.core.doc-asl/doclinks/6/user-guide.html](https://zoidberg.ukp.informatik.tu-darmstadt.de/jenkins/job/DKProCoreDocumentation(GitHub)/de.tudarmstadt.ukp.dkpro.core$de.tudarmstadt.ukp.dkpro.core.doc-asl/doclinks/6/user-guide.html). Accessed: 2018-02-26.
- [doc15] What is docker? <https://www.docker.com/what-docker>, 2015. Accessed: 2018-02-26.
- [EdCG14] Richard Eckart de Castilho and Iryna Gurevych. A broad-coverage collection of portable nlp components for building shareable analysis pipelines. In *Proceedings of the Workshop on Open Infrastructures and Analysis Frameworks for HLT*, pages 1–11, Dublin, Ireland, August 2014. Association for Computational Linguistics and Dublin City University.
- [ESI<sup>+</sup>12] Edward A Epstein, Marshall I Schor, BS Iyer, Adam Lally, Eric W Brown, and Jaroslaw Cwiklik. Making watson fast. *IBM Journal of Research and Development*, 56(3.4):15–1, 2012.
- [Fer12] David A Ferrucci. Introduction to “this is watson”. *IBM Journal of Research and Development*, 56(3.4):1–1, 2012.

- [FL04] David Ferrucci and Adam Lally. Uima: an architectural approach to unstructured information processing in the corporate research environment. *Natural Language Engineering*, 10(3-4):327–348, 2004.
- [FLVN09] David Ferrucci, Adam Lally, Karin Verspoor, and Eric Nyberg. Unstructured information management architecture (UIMA) version 1.0. OASIS Standard, mar 2009.
- [GR12] John Gantz and David Reinsel. The digital universe in 2020: Big data, bigger digital shadows, and biggest growth in the far east. *IDC iView: IDC Analyze the future*, 2007(2012):1–16, 2012.
- [JM14] Dan Jurafsky and James H Martin. *Speech and language processing*, volume 3. Pearson London;, 2014.
- [KSDG12] Valmeek Kudesia, Judith Strymish, Leonard D’Avolio, and Kalpana Gupta. Natural language processing to identify foley catheter–days. *Infection control and hospital epidemiology*, 33(12):1270–1272, 2012.
- [Lid01] Elizabeth D Liddy. Natural language processing. 2001.
- [RBJB<sup>+</sup>10] Cartic Ramakrishnan, William A Baumgartner Jr, Judith A Blake, Gully APC Burns, K Bretonnel Cohen, Harold Drabkin, Janan Eppig, Eduard Hovy, Chun-Nan Hsu, Lawrence E Hunter, et al. Building the scientific knowledge mine (sciknowmine): a community-driven framework for text mining tools in direct service to biocuration. *Language Resources and Evaluation*, page 33, 2010.

## Eidesstattliche Erklärung

Hiermit versichere ich, Simon Gehring, dass ich die vorliegende Masterarbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Die Stellen meiner Arbeit, die dem Wortlaut oder dem Sinne nach anderen Werken und Quellen, einschließlich Quellen aus dem Internet, entnommen sind, habe ich in jedem Fall unter Angabe der Quelle deutlich als Entlehnung kenntlich gemacht. Dasselbe gilt sinngemäß für Tabellen, Karten und Abbildungen.

Unterschrift: \_\_\_\_\_  
Simon Gehring, Student  
Universität Bonn