

Der TextImager als Front- und Backend für das verteilte NLP von Big Digital Humanities Data

Hemati, Wahed

hemati@em.uni-frankfurt.de
Goethe-Universität Frankfurt, Deutschland

Mehler, Alexander

mehler@em.uni-frankfurt.de
Goethe-Universität Frankfurt, Deutschland

Uslu, Tolga

uslu@em.uni-frankfurt.de
Goethe-Universität Frankfurt, Deutschland

Abrami, Giuseppe

abrami@em.uni-frankfurt.de
Goethe-Universität Frankfurt, Deutschland

Immer mehr Disziplinen benötigen Natural Language Processing (NLP) Werkzeuge, um automatische Textanalysen auf verschiedenen Ebenen der Sprache durchzuführen. Die Anzahl der NLP-Werkzeuge wächst rasant¹. Auch die Anzahl der frei oder anderweitig zugänglichen Ressourcen wächst. Angesichts dieser wachsenden Zahl an Werkzeugen und Ressourcen ist es schwierig, den Überblick zu behalten; gleichzeitig ist ein Computational-Linguistic-Framework, das große Datenmengen aus verschiedenen Quellen verarbeiten kann, noch nicht etabliert. Ein solches Framework sollte in der Lage sein, Daten verteilt zu verarbeiten und gleichzeitig eine standardisierte Programmier- und Modellschnittstelle bereitzustellen. Darüber hinaus sollte es modular und leicht erweiterbar sein, um die ständig wachsende Palette neuer Ressourcen und Tools zu integrieren. Das Framework muss offen genug für Erweiterungen Dritter sein, wobei jede Erweiterung für die gesamte Community zugänglich bleibt. Das Framework sollte es zudem Dritten ermöglichen, den Zugang zu ihren Erweiterungen zu beschränken, wenn dies beispielsweise durch Urheberrecht, geistiges Eigentum oder Datenschutz erforderlich ist. Um diesen Anforderungen gerecht zu werden, haben wir den TextImager (Hemati 2016, Mehler et al. 2018) um ein verteiltes Serversystem mit Cluster-Computing-Funktionen auf der Basis von UIMA (Ferrucci and Lally 2004) weiterentwickelt.

UIMA ist ein Framework zur Verwaltung von Datenflüssen zwischen Komponenten. Es bietet standardisierte Interfaces zur Erstellung von Komponenten

an. Dabei können die Komponenten einzeln oder im Verbund in einer Pipeline-Struktur ausgeführt werden. UIMA bietet weitgehende Möglichkeiten der sequenziellen Ordnung von NLP-Werkzeugen und verspricht, auch in Zukunft von der Community weiterentwickelt zu werden: Prozess-Management auf der Basis von UIMA erscheint nach derzeitigem Stand daher als erste Wahl im Bereich von NLP und DH.

TextImager bietet eine Vielzahl von UIMA-basierten NLP-Komponenten an, darunter unter anderen einen Tokenisierer, einen Lemmatisierer, einen Part-Of-Speech-Tagger, einen Named-Entity-Parser und einen Dependency Parser, und zwar für eine Vielzahl von Sprachen, darunter Deutsch, Englisch, Französisch und Spanisch. Dieses Spektrum an Werkzeugen besteht allerdings nicht ausschließlich aus Eigenentwicklungen, sondern wird maßgeblich um Entwicklungen Dritter erweitert, wozu unter anderem die Tool-Palette von Stanford CoreNLP (Manning 2014), OpenNLP (OpenNLP 2010) und DKpro (Eckart de Castilho 2014) zählen.

In Zeiten von Big Data wird es immer relevanter, Daten schnell zu verarbeiten. Aus diesem Grund ist TextImager als Multi-Server- und zugleich als Multi-Instanz-Clusteraufgebaut, um das verteilte Verarbeiten von Daten zu ermöglichen. Dafür setzt TextImager auf UIMAs Cluster-Management-Dienste UIMA-AS² und UIMA-DUCC³ auf.

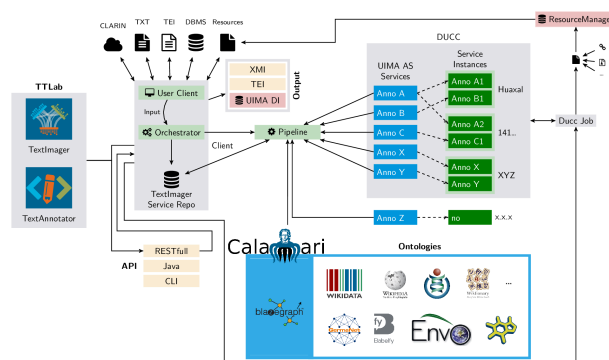


Abbildung 1

Abbildung 1 zeigt eine schematische Darstellung von TextImager. Jede NLP-Komponente läuft als UIMA-AS Webservice auf dem Computing-Cluster des TextImager. Dabei können mehrere Instanzen einer Komponente instanziiert (s. Abbildung 1, Service Instances) werden und dennoch über eine Webservice-Schnittstelle (s. Abbildung 1, UIMA AS Services) angesprochen werden. Dazu wird das Java Messaging Service (JMS) verwendet, das die Kommunikation zwischen verschiedenen Komponenten einer verteilten Anwendung ermöglicht. JMS implementiert ein Point-to-Point-Kommunikationssystem. Dieser Kommunikationstyp basiert auf dem Konzept der message queues (Warteschlangen), senders (Sender) und receivers

(Empfänger). Jedem Dienst ist eine Eingabewarteschlange und eine Ausgabewarteschlange zugeordnet. Um mehrere Instanzen einer Komponente zu verteilen, verbinden sich die Instanzen mit der gleichen Service-Eingangswarteschlange. Die Instanzen erhalten aus dieser Warteschlange Arbeitseinheiten. Nach der Verarbeitung wird das Ergebnis an eine Ausgabewarteschlange zurückgegeben. Die Ausgabewarteschlange eines Dienstes kann an eine Eingabewarteschlange eines anderen Dienstes angeschlossen werden, um eine Pipeline zu erstellen. Aufgrund dieser Ein- und Ausgabewarteschlangen-Systematik kann jeder Service Arbeitseinheiten asynchron bearbeiten. Durch diese Architektur ist TextImager eine Multi-Server-, Multi-Service- und Multi-Service-Instanz-Architektur.

Darüber hinaus bietet TextImager ein Toolkit, das es jedem Entwickler ermöglicht, einen eigenen TextImager-Cluster aufzusetzen und Services im TextImager-System hinzuzufügen. Entwickler können den Zugriff auf die Dienste einschränken, wenn dies wie oben beschrieben erforderlich ist, was mittels der Integration des ResourceManagers (Gleim 2012) und des AuthorityManagers (Gleim 2012) realisiert wird.

Durch Freigabe des Quellcodes des TextImager und die Bereitstellung von Leitlinien für dessen Erweiterung wollen wir es Dritten ermöglichen, ihre NLP-Software über die Webservices von TextImager zu vertreiben, so dass die gesamte wissenschaftliche Gemeinschaft davon profitiert.

Installationsanweisungen und Beispiele für die Einrichtung eines TextImager-Servers finden Nutzer in folgendem GitHub-Repository: <https://github.com/texttechnologylab/textimager-server>.

Der Beitrag erörtert die Möglichkeiten und Grenzen des NLP von Big Data, stellt den TextImager als Werkzeug für diesen Bereich zur Diskussion und zeigt anhand von drei Nutzungsszenarien Einsatzmöglichkeiten in den DH auf.

desktop, in: Proceedings of the Workshop on Service-oriented Architectures (SOAs) for the Humanities: Solutions and Impacts.

Mehler, Alexander / Hemati, Wahed / Gleim, Rüdiger / Baumartz, Daniel (2018): *VienNA: Auf dem Weg zu einer Infrastruktur für die verteilte interaktive evolutionäre Verarbeitung natürlicher Sprache*, in: Forschungsinfrastrukturen und digitale Informationssysteme in der germanistischen Sprachwissenschaft, H. Lobin, R. Schneider, and A. Witt, Eds., Berlin: De Gruyter, 2018, vol. 6.

Hemati, Wahed / Uslu, Tolga / Mehler, Alexander (2016): *Textimager: a distributed uima-based system for nlp*, in: Proceedings of the COLING 2016 System Demonstrations. Federated Conference on Computer Science and Information Systems.

Manning, Christopher / Surden, Mihai / Bauer, John / Finkel, Jenny / Bethard, Steven / McClosky, David (2014): The Stanford CoreNLP natural language processing toolkit, in: Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations 55–60.

OpenNLP (2010): *Apache OpenNLP*, in: <http://opennlp.apache.org> [letzter Zugriff 05. Oktober 2018]

Fußnoten

1. <https://github.com/topics/nlp>
2. <https://uima.apache.org/doc-uimaas-what.html>
3. <https://uima.apache.org/doc-uimaducc-whatit.html>

Bibliographie

de Castilho, Richard Eckart / Gurevych, Iryna (2014): *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 1–11.

Ferrucci, David / Lally, Adam (2004): *UIMA: an architectural approach to unstructured information processing in the corporate research environment.*, in: Natural Language Engineering 10(3–4) 327–348.

Gleim, Rüdiger / Mehler, Alexander / Ernst, Alexandra (2012): *SOA implementation of the ehumanities*