

# RECHERCHEBERICHT

## APRT-17

Im Rahmen dieses Projektes soll eine Software erstellt werden, die zu einem Themengebiet (bspw. Städte, Künstler, Pflanzenarten) Informationen aus dem Internet bezieht und daraus Präsentationen generiert. So könnte bereits aus dem Wikipedia Artikel zu Leipzig eine Menge an Informationen abgeleitet werden. Diese strukturierten Informationen, wie die historische Entwicklung (incl. Bildern), berühmte Persönlichkeiten und Zahlen / Fakten zur Stadt, sollen anschließend in einem für das Projekt sinnvollen Datenform in Präsentationen überführt.

Teilnehmer: Lucas Lange, Andrea Niebuhr, Johannes Bräuer, Kilian Schwabe, Hunter Foo und Francesco Mandry

# Inhaltsverzeichnis

<b>Begriffe (Glossar)</b> .....	3
<b>Konzepte</b> .....	5
Teamaufteilung .....	5
Java .....	5
Anwendungstyp .....	5
Präsentationsformat .pptx .....	5
DBpedia .....	6
Auslesen von Informationen in Wikipedia (1) .....	6
Recherche zum Aussehen einer guten Präsentation.....	6
Präsentationsthemen und Quellen .....	7
<b>Aspekte</b> .....	7
GUI .....	7
Python .....	7
NLP .....	7
Auslesen von Informationen in Wikipedia (2) .....	8
<b>Quellen</b> .....	8

# Begriffe (Glossar)

**Python<sup>1</sup>**

Eine universelle, üblicherweise interpretierte höhere Programmiersprache.

**Jython<sup>2</sup>**

eine reine Java<sup>30</sup>-Implementierung der Programmiersprache Python<sup>1</sup> und ermöglicht somit die Ausführung von Python<sup>1</sup>-Programmen auf jeder Java<sup>30</sup>-Plattform<sup>6</sup>.

**GUI<sup>3</sup>**

eine grafische Benutzeroberfläche.

**Webanwendung<sup>4</sup>**

Die Datenverarbeitung und -auswertung findet [...] hauptsächlich auf einem entfernten Webserver statt.

**Desktopanwendung<sup>5</sup>**

[...] [werden] lokal auf dem Rechner des Benutzers installiert und dort ausgeführt.

**Plattform<sup>6</sup>**

bezeichnet in der Informatik eine einheitliche Grundlage, auf der [...] [Programme] ausgeführt. und entwickelt werden können.

**PowerPoint<sup>7</sup>**

Präsentationssoftware<sup>8</sup> von Microsoft.

**Präsentationssoftware<sup>8</sup>**

Computerprogramm, das die Erarbeitung und Präsentation eines Vortrages oder Referates unterstützt.

**API<sup>9</sup>**

(engl. „application programming interface“) - ein Programmteil, mit dessen Hilfe sich andere Programme an das Softwaresystem anbinden können.

**HTML<sup>10</sup>**

(engl. Hypertext Markup Language) - textbasierte Auszeichnungssprache zur Strukturierung digitaler Inhalte wie Texte, Hyperlinks<sup>29</sup>, Bilder und anderen.

**Office Open XML<sup>11</sup>**

von Microsoft entwickelte Dateiformate zur Speicherung von Bürodokumenten auf XML<sup>15</sup>-Basis.

**Library<sup>12</sup>**

(auch (Programm-)bibliothek) - Sammlung von Unterprogrammen/-Routinen; nicht eigenständig lauffähig.

**parsen/Parser<sup>13</sup>**

Programm zur Zerlegung und Umwandlung einer Eingabe in ein bestimmtes Format.

**String<sup>14</sup>**

eine Folge von Zeichen aus einem definierten Zeichensatz; in Java<sup>30</sup> ein Datentyp.

**XML<sup>15</sup>**

(Extensible Markup Language) eine Auszeichnungssprache zur Darstellung hierarchisch strukturierter Daten in Form von Textdateien.

**automatic summarization<sup>16</sup>**

automatisierte Zusammenfassung von Texten.

**extraction<sup>17</sup>**

wichtige Thesen aus einem Text werden extrahiert, ohne dabei die Sätze oder Wörter zu verändern.

**abstraction<sup>18</sup>**

Umschreibung von wichtigen Kernpunkten („paraphrasing“).

**NLP<sup>19</sup>**

natural language processing, automatisierte Textverarbeitung und/oder –generierung.

**JSON<sup>20</sup>**

JavaScript Object Notation, kompaktes Datenformat in einer einfach lesbaren Textform zum Zweck des Datenaustauschs zwischen Anwendungen.

**URL<sup>21</sup>**

Uniform Resource Locator, lokalisiert eine Ressource, z.B.: eine Webseite.

**String Operatoren<sup>22</sup>**

Funktionen zur Manipulation von Zeichenketten.

**Framework<sup>23</sup>**

ist ein Programmiergerüst, das in der Softwaretechnik [...] verwendet wird.

**Folie<sup>24</sup>**

eine einzelne Seite einer am PowerPoint<sup>7</sup>-Präsentation(Präsentations-) Thema: eine Kategorie, zu der das zu erstellende Programm Präsentationen anbietet, z.B. „Städte“.

**(Präsentations-) Instanz<sup>25</sup>**

Subjekt einer spezifischen Präsentation, z.B. „Leipzig“.

**Datenbank<sup>26</sup>**

eine strukturierte Sammlung von Information; Daten sind systematisch katalogisiert und können gezielt abgerufen werden.

**SPARQL<sup>27</sup>**

eine graphenbasierte Abfragesprache für RDF<sup>28</sup>.

**RDF<sup>28</sup>**

(Resource Description Framework) bezeichnet eine technische Herangehensweise im Internet zur Formulierung logischer Aussagen über beliebige Dinge (Ressourcen).

**Hyperlink<sup>29</sup>**

ist ein Querverweis [...], der funktional einen Sprung zu einem anderen elektronischen Dokument oder an eine andere Stelle innerhalb eines Dokuments ermöglicht.

**Java<sup>30</sup>**

eine objektorientierte Programmiersprache [...].

**Apache Jena<sup>31</sup>**

Java<sup>30</sup> Bibliothek zur Arbeit mit unter Anderem RDF<sup>28</sup> und SPARQL<sup>27</sup>.

# Konzepte

## Teamaufteilung

Wir teilen das Team in zwei Gruppen, wobei sich die Eine für die Extraktion der nötigen Informationen aus dem Netz und die Andere für die Ausgabe dieser Daten als angemessene Präsentation zuständig sieht, so dass diese Bereiche besser untereinander besprochen werden können. Die Bindeglieder dieser Teilaufgaben sind gute Kommunikationswege und das wöchentliche Treffen, welches Raum für regen Austausch untereinander bietet.

## Java<sup>30</sup>

Als Hauptsächliche Programmiersprache für das Projekt wählen wir „Java<sup>30</sup>“ auf Grund der guten Vertrautheit aller Gruppenmitglieder mit deren Syntax, sowie der weitreichenden Möglichkeiten, welche auch auf den vielen verfügbaren Bibliotheken<sup>12</sup> beruht.

## Anwendungstyp

Wichtig ist die Frage, ob eine Web- oder Desktopanwendung als Ergebnis des Projekts das Ziel sein sollte. Desktopanwendungen<sup>5</sup> bieten meist eine schnellere Laufzeit, interagieren besser mit der PC-Hardware und eine GUI<sup>3</sup> ist nicht speziell nötig, aber auch leicht zu bauen. Webanwendungen<sup>4</sup> sind hingegen sehr portabel und plattformunabhängig einsetzbar. Des Weiteren fällt eine Wartung dieser meist leichter. In Hinsicht auf unser Projekt sollten wir den Kunden entscheiden lassen, doch gib es keinen Grund hier eine Webanwendung zu entwickeln, da dies für unser Projekt eher von geringem Mehrwert wäre.

## Präsentationsformat .pptx

Auf jede Präsentationssoftware<sup>8</sup> kommt mindestens ein eigenes Präsentationsformat, jedoch sollte für unser Projekt nur ein Ausgabeformat eine Rolle spielen. Aufgrund der weiten Verbreitung, besonders auch in Bildungseinrichtungen, stehen hier die PowerPoint<sup>7</sup>-Präsentation und ihre Formate an erster Stelle. Im Besonderen das „.pptx“-Format.

Beim Präsentationsformat .pptx handelt es sich um eines der bekanntesten Formate für PowerPoint<sup>7</sup> Präsentationen, welches im Office Open XML<sup>11</sup>-Format gespeichert wird. Die API<sup>9</sup> „python<sup>1</sup>-pptx“, eine sehr gut dokumentierte Python<sup>1</sup> library<sup>12</sup> mit vielen Features, um Präsentationen zu erstellen und zu bearbeiten. (<https://python-pptx.readthedocs.io/en/latest>) Da unser Programm jedoch in Java<sup>30</sup> wird, sei „python<sup>1</sup>-pptx“ hier nur wegen der guten Kompatibilität von Java<sup>30</sup> und Python<sup>1</sup> aufgeführt.

Eine ebenfalls ausgezeichnet dokumentierte API<sup>9</sup> für Java<sup>30</sup> ist die library<sup>12</sup> docx4j, mit der man neben .docx auch .pptx, sowie .xlsx-Dateien erstellen und Bearbeiten kann.

(<http://www.docx4java.org/trac/docx4j>)

Mit dieser library<sup>12</sup> erstellte Präsentationen werden vor dem parsen<sup>13</sup> zu .pptx in einem String<sup>14</sup> erstellt und gespeichert, welcher ein ähnliches Layout wie ein HTML<sup>10</sup>-Dokument aufweist.

## DBpedia

Im Projekt DBpedia wird versucht Informationen aus Wikipedia strukturiert zu extrahieren und diese online zugänglich zu machen.

Suchvorgänge auf DBpedia werden mit der Abfragesprache „SPARQL<sup>27</sup>“ durchgeführt. Diese ermöglicht es bestimmte Daten aus der Wissensdatenbank zu extrahieren. So lassen sich Beispielsweise alle Informationen zur Stadt „Leipzig“ mit nur einem Zugriff extrahieren. Problematisch ist nur, dass man bis auf einen allgemein gehaltenen Textabschnitt, nur sehr faktenhaft notierte Informationen finden und extrahieren kann. Folgt man den Hyperlinks<sup>29</sup>, welche für viele Fakten existieren, kommt man zwar zu tiefreichenden Informationen, doch wird mit steigender Suchtiefe auch der Zusammenhang zum Eigentlichen Thema in der Regel immer geringer. Dies führt dazu, dass man sich themenabhängig nicht alleinig auf DBpedia als Informationsquelle für eine ansprechende Präsentation stützen kann. Der letztlich erwartete Informationsgehalt und dessen geplante Darstellung spielen hier eine wichtige Rolle. Speziellere Themen könnten ein Problem darstellen, da die Datenbank hier möglicherweise nur wenige interessante Fakten bieten kann.

Der sinnvolle Einsatz von DBpedia hängt also von der effektiven Nutzung der Sprache „SPARQL<sup>27</sup>“, den Präsentationsthemen und der Art der zu gewinnenden Informationen ab. Noch zu erwähnen ist, dass der Zugriff mit „SPARQL<sup>27</sup>“ in Java<sup>30</sup> beispielsweise über die Bibliothek „Apache Jena“ einfach möglich gemacht wird, was die Suche auf DBpedia für unser Projekt zugänglich macht.

## Auslesen von Informationen in Wikipedia (1)

Im Web findet sich eine Vielzahl verschiedener Daten die unterschiedlich gut strukturiert sind. Plattformen<sup>6</sup> wie Wikipedia geben verschiedensten Themenbereichen bzw. Artikeln eine Struktur aus Inhaltsverzeichnis, Infobox, Überschriften, Text und Bilder. Mittels zentraler Datenspeicher wie Wikidata und DBpedia, die auf Wikipedia basieren, werden Artikel stichpunktartig zusammengefasst. Auf diese Informationen kann man mittels Datenbanksprachen (SPARQL<sup>27</sup>) oder anderen APIs<sup>9</sup> zugreifen. Für spezifische Texte in Wikipedia reicht es auch Artikel in JSON<sup>20</sup>- oder XML<sup>15</sup>-Formate zu parsen<sup>13</sup>.

## Recherche zum Aussehen einer guten Präsentation

Wenn man das Internet nach Tipps für gut aussehende PowerPoint<sup>7</sup>-Präsentationen durchsucht, stellt man schnell fest, dass die meisten dieser Hilfen ein spezielles Thema mit einbeziehen. Sie sollen den Leser also unterstützen, seine eigene Präsentation anschaulicher zu gestalten. Nun mögen einige dieser Tipps auch in Bezug auf dieses Projekt durchaus Anwendung finden, andere sind nur schwer umsetzbar.

So liest man immer wieder, wie wichtig große Schrift ist und geringe Textmenge je Folie<sup>24</sup> einzuhalten ist. Dagegen das Design der Folien<sup>24</sup>-Hintergründe oder die Schriftart auf das Thema anzupassen, ist aufgrund der Themenvielfalt, die im Projekt gesetzt werden nicht umsetzbar. Letztendlich hängen Design, so wie einige andere Eckpunkte der Präsentation (z.B. die Einbindung von Bildern) noch von Vorgaben der Betreuer und Umfang des Präsentation-APIs<sup>9</sup> ab.

## Präsentationsthemen und Quellen

Für die Präsentationen sind Themen (z.B. „Städte“) zu finden, zu deren einzelnen Instanzen<sup>25</sup> (z.B. „Leipzig“) Präsentationen erstellt werden können. Diese Themen sollen einige Eigenschaften erfüllen. Sie sollen weit gefasst sein, d.h. zu jedem Thema sollen möglichst viele Instanzen<sup>25</sup> existieren, um viele Präsentationen daraus erstellen zu können. Die Themen sollen hauptsächlich durch strukturierte Parameter definiert und die einzelnen Instanzen<sup>25</sup> durch ihre verschiedenen Werte abgegrenzt werden, um das automatisierte Einlesen der Information und die sinnvolle Aufarbeitung in der Präsentation zu erleichtern. Information zu den Themen soll aus möglichst wenigen (ideal: einer) vertrauenswürdigen Quellen verfügbar sein, in der Information zu möglichst vielen Instanzen<sup>25</sup> möglichst strukturiert gesammelt ist. Da die Präsentationen Grundstock und Testobjekte für eine wissenschaftlich orientierte Plattform<sup>6</sup> sind, sollen die Themen eine entsprechende wissenschaftliche Relevanz haben. Diese Anforderungen legen nahe, gepflegte Online-Datenbanken<sup>26</sup> als Quellen zu verwenden, wie z.B. DBpedia oder Wikidata. Auch gut strukturierte (Elemente von) Webseiten, wie z.B. Infoboxen auf Wikipedia, sind als Quellen interessant. Für Informationen, die nicht kurz und tabellarisch auffindbar sind, muss auf Natural Language Processing (NLP<sup>19</sup>) zurückgegriffen werden, um Information aus Fließtexten, z.B. Wikipediaartikeln, zu extrahieren. Mögliche Themen, die die genannten Anforderungen erfüllen, sind Städte und Minerale.

# Aspekte

## GUI<sup>3</sup>

Eine schöne, intuitive GUI<sup>3</sup> sehen wir eher als optionales Ziel des Projekts, da sie für die Arbeit des Programms nicht existenzielle Bedeutung hat. Dabei gehen wir nach dem Verwendungszweck, der von spezieller Natur ist und somit keine massentaugliche Bedienbarkeit vorgibt.

## Python

Beim Blick auf weitere Programmiersprachen lässt sich besonders „Python<sup>14</sup>“ nicht gänzlich ausschließen, da im Hinblick auf dieses Projekt dort ebenfalls viele nützliche Bibliotheken<sup>12</sup> zur Verfügung stehen. Vorteilhaft ist das leichte Zusammenspiel von „Python<sup>14</sup>“-Programmen in „Java<sup>30</sup>“ durch beispielsweise „Jython<sup>24</sup>“.

## NLP<sup>19</sup>

„Automatic summarization<sup>16</sup>“, ein Teilgebiet des NLP<sup>19</sup> (natural language processing), dient dazu, Fließtexte möglichst kurz zusammenzufassen, ohne wichtige Schlüsselaussagen des Textes zu vernachlässigen. Es gibt zwei verschiedene Ansätze: die „extraction<sup>17</sup>“ kopiert einzelne Sätze aus einem Fließtext heraus, die „abstraction<sup>18</sup>“ umschreibt Sätze und verknüpft komplexe Sachverhalte logisch miteinander. Im Rahmen dieses Projekts ist letztere zu bevorzugen, da sie möglichst kurze Fakten aus Informationstexten liefert, ohne

dabei den Sinn des Textes zu verändern. Die Web-Applikation „smmry“ (smmry.com) bietet diese Dienste an. Hier lässt der User einen Fließtext einlesen und bestimmt, in wie vielen Sätzen er zusammengefasst werden soll. Die Anzahl der Sätze passend zum Informationsgehalt automatisiert festzulegen ist problematisch, da sich die Menge von wichtigen Fakten nicht anhand von z.B. der Textlänge ermitteln lässt. Es liegt nahe, stattdessen Werte festzulegen, um so eine bestimmte Anzahl von Folien mit einer bestimmten Menge von Fakten pro Folie zu erstellen.

NLP<sup>19</sup> ist ein aktives Forschungsgebiet und die Ergebnisse beider Arten von automatic summarization<sup>16</sup> schwanken in ihrer Qualität. Dennoch ist es eine sinnvolle Methode, um für die Präsentationen relevante Information aus Fließtexten zu extrahieren.

## Auslesen von Informationen in Wikipedia (2)

Mittels Query Abfragen (Abfragesprachen) in der URL<sup>21</sup> kann man konkrete Inhalte aus Wikipedia Artikeln filtern. Informationen der Befehle dazu finden sich in der Dokumentation für die MediaWiki API<sup>9</sup>. <https://en.wikipedia.org/w/api.php>. Das gleiche gibt es auch für die TMDb mit kurzer Dokumentation. <https://www.themoviedb.org/documentation/api>

Durch Abfrage des geparsten Ergebnisses in JAVA<sup>30</sup>, kann der Inhalt in einem String<sup>14</sup> gespeichert werden und mittels String<sup>14</sup>-Operatoren von unnötigen Zeichen und Text befreit werden. Dies kann bei langen Texten umständlich werden, eignet sich jedoch für einzelne Themen eines Artikels, da dieser wenige Text umfassen.

Eine andere Möglichkeit bietet SPARQL<sup>27</sup>, das nicht nur DBPedia, sondern auch Wikidata unterstützt. Apache Jena<sup>31</sup> als Framework<sup>23</sup> kann dafür genutzt werden, mit JAVA<sup>30</sup> und SPARQL<sup>27</sup> die Datenbanken<sup>26</sup> auszulesen.

## Quellen

[https://de.wikipedia.org/wiki/JavaScript\\_Object\\_Notation](https://de.wikipedia.org/wiki/JavaScript_Object_Notation)  
[https://de.wikipedia.org/wiki/Uniform\\_Resource\\_Locator](https://de.wikipedia.org/wiki/Uniform_Resource_Locator)  
<https://de.wikipedia.org/wiki/Programmierschnittstelle>  
[https://de.wikipedia.org/wiki/Hypertext\\_Markup\\_Language](https://de.wikipedia.org/wiki/Hypertext_Markup_Language)  
[https://de.wikipedia.org/wiki/Extensible\\_Markup\\_Language](https://de.wikipedia.org/wiki/Extensible_Markup_Language)  
[https://de.wikipedia.org/wiki/Office\\_Open\\_XML](https://de.wikipedia.org/wiki/Office_Open_XML)  
<https://de.wikipedia.org/wiki/Parser>  
<https://de.wikipedia.org/wiki/Zeichenkette>  
<https://de.wikipedia.org/wiki/Programmierschnittstelle>  
<https://de.wikipedia.org/wiki/Webanwendung>  
[https://de.wikipedia.org/wiki/Plattform\\_\(Computer\)](https://de.wikipedia.org/wiki/Plattform_(Computer))  
<https://de.wikipedia.org/wiki/Pr%C3%A4sentationsprogramm>  
[https://de.wikipedia.org/wiki/Extensible\\_Markup\\_Language](https://de.wikipedia.org/wiki/Extensible_Markup_Language)  
<https://de.wikipedia.org/wiki/Programmbibliothek>  
<https://de.wikipedia.org/wiki/Abfragesprache>  
[https://de.wikipedia.org/wiki/Graph\\_\(Graphentheorie\)](https://de.wikipedia.org/wiki/Graph_(Graphentheorie))  
[https://de.wikipedia.org/wiki/Resource\\_Description\\_Framework](https://de.wikipedia.org/wiki/Resource_Description_Framework)  
<https://de.wikipedia.org/wiki/Hyperlink>  
[https://de.wikipedia.org/wiki/Java\\_\(Programmiersprache\)](https://de.wikipedia.org/wiki/Java_(Programmiersprache))  
[https://jena.apache.org/about\\_jena/about.html](https://jena.apache.org/about_jena/about.html)  
<https://de.wikipedia.org/wiki/SPARQL>