

Vorwissenschaftliche Arbeit

# Programme und Fallbeispiele zur automatisierten Extraktion von Daten aus Webseiten

Verfasser

Rafael Vrečar, 8A

Schule

BRG19 Krottenbachstraße 11-13, 1190 Wien

Betreuer

Dipl.-Ing. Christian Schöbel

Abgabedatum

8. Februar 2016

## Abstract

[Eine deutschsprachige Fassung des Abstracts ist im Anhang (A.I) angeführt.]

The automated extraction of data from web pages has the objective of saving time and avoiding mistakes. This prescientific paper is concerned with an automation process of this type.

The paper starts with the description of the proceedings which take place while surfing the web. Thereby, the Hypertext Transfer Protocol (HTTP) and its functional units client and server are being discussed. Furthermore, an empirical classification of the different kinds of web pages is taken. Subsequently, various tools which are necessary for this automation process are being examined. These are divided into the classes “Tools for retrieval” and “Tools for processing”.

After the theoretical description of the tools the automation process is being analysed based on two case examples. In the first example the list of consultation hours of BRG19 is being extracted, whereas in the second example it is the individual list of modules which is being extracted. Both lists are being extracted from the downloaded files of the Hypertext Markup Language format (HTML) and converted into files of the Comma-separated values format (CSV).

Lastly, evaluations and conclusions concerning the automation process are being made. These considerations analyse which kind of knowledge is required to set up an automation of this kind and whether automation is reasonable or not.

## Vorwort

Die Automatisierung von Arbeitsprozessen und die damit verbundene Effizienzmaximierung genoss bei mir schon lange reges Interesse, und da für mich schon vor dem Prozess der Themenfindung klar war, dass ich eine Arbeit mit informatischem Hintergrund schreiben möchte, war es für mich wenig überraschend, im Frühjahr 2015 folgende Themenstellung einzureichen:

Programme und Fallbeispiele zur automatisierten Extraktion von Daten aus Webseiten.

An dieser Stelle möchte ich mich bei meiner Familie und meinen Freunden bedanken, die mich stets unterstützt haben. Der größte Dank gilt allerdings meinem Betreuer Dipl.-Ing. Christian Schöbel, der sich aller Fragen und Anliegen annahm, welche ich ihm im Schreibprozess dieser Arbeit entgegenbrachte.

Wien, 19. Jänner 2016

Rafael Vrečar

## Wichtige Hinweise und Erklärungen

Im Vorfeld dieser vorwissenschaftlichen Arbeit sind einige Hinweise und Erklärungen anzustellen, um das Auftreten etwaiger Fragen, die innerhalb der eigentlichen Arbeit nicht beantwortet werden, zu vermeiden.

### Hinweise zu Zitaten und Quellen(angaben)

In dieser vorwissenschaftlichen Arbeit befinden sich keine wörtlichen Zitate, da in einer technischen Arbeit, im Gegensatz zu einer geisteswissenschaftlichen Arbeit, in der Regel der Inhalt eines Texts, nicht aber dessen Wortlaut von Bedeutung ist. Sämtliche Quellen wurden grundsätzlich in der Form

vgl. URL (Datum des letzten Zugriffs)

(Anmerkung: weitere Spezifikation durch: „u. a.“, Kapitel- und Seitenangaben et cetera möglich)

in Fußnoten angegeben, da diese Form für eine auf Internetquellen basierende Arbeit als einzig sinnvolle erscheint. Diese Arbeit basiert zur Gänze auf Internetquellen, da die autoritativen Quellen ihre Inhalte in der Regel nicht in Papierform, sondern auf ihren Webseiten veröffentlichen. Besagte Quellen sind in einem Literaturverzeichnis (C Literaturverzeichnis) aufsteigend nach Fußnotennummern sortiert aufgelistet. Sie wurden im Zeitraum vom 19. bis zum 29. Jänner 2016 sorgfältig auf ihre Gültigkeit überprüft.

### Hinweis zu Quellenangaben Glossarbegriffe betreffend

Im Verlauf dieser vorwissenschaftlichen Arbeit werden einige Begriffe auftreten, die in einem separaten Glossar erklärt sind. Diese sind entsprechend gekennzeichnet, eine genaue Erklärung diesbezüglich erfolgt noch in diesem Kapitel. Sofern Erklärungen schon im Text angestellt werden, da sie für das Verständnis der Arbeit erforderlich sind, werden die Quellen für die Erklärung allgemeiner Begriffe nicht nochmals angegeben. Besagte Quellen sind im Glossar angeführt. Der Text referenziert also auf das Glossar, das Glossar auf die entsprechende Quelle. Webseiten zu diversen spezifischen Begriffen, wie zum Beispiel Hypertext Transfer Protocol (HTTP), PHP Hypertext Processor (PHP) et cetera, werden hingegen bereits im Text als Fußnote angegeben.

### Hinweis zum Abrufzeitpunkt der URLs von [www.brg19.at](http://www.brg19.at)

Diverse in den Fallbeispielen vorkommende URLs von [www.brg19.at](http://www.brg19.at) wurden allesamt, sofern nicht anders angegeben, am 20. Jänner 2016 auf ihre Gültigkeit überprüft. Änderun-

gen sind möglich und liegen nicht im Ermessen des Autors dieser Arbeit. Die auf dem beiliegenden Speichermedium vorhandenen Dateien wurden allesamt unter Absprache mit Dipl.-Ing. Christian Schöbel am 21. Jänner 2016 heruntergeladen.

## Erklärungen zur Darstellung einzelner Textstücke

Im Folgenden werden Erklärungen zur Darstellung diverser Textstücke abgegeben. Dies betrifft Schriftart, Schriftgröße und Einrückung.

### Erklärung zur Kennzeichnung von im Glossar enthaltenen Begriffen

Begriffe, die im Glossar enthalten sind, werden ab jetzt bei jeder Verwendung kursiv in der Schriftart Cambria dargestellt. Dabei sind auch Abwandlungen der Begriffe vollständig markiert. Zwischen Einzahl und Mehrzahl wird nicht unterschieden. Begriffe in Tabellen werden ebenfalls in dieser Form markiert.

Von dieser Markierung ausgeschlossen sind hingegen:

- Überschriften und im Text gesetzte Referenzen auf diese
- Bild-, Tabellen- und *Quelltext*unterschriften
- Begriffe in Grafiken
- Namen von Programmen, die im Rahmen dieser Arbeit als *Befehle* im *Terminal* ausgeführt wurden, sowie deren *Parameter*
- diverse Verzeichnisse

Abkürzungen sind ebenfalls im Glossar enthalten und werden in der gleichen Form dargestellt. Sie sind dort durch die Voranstellung von „Abk.“ entsprechend gekennzeichnet.

### Erklärung zur Darstellung von URLs und Dateinamen

Vollständige *URLs*, also jene, die auch die Angabe des *Protokolls* enthalten, sowie Dateinamen, werden in dieser vorwissenschaftlichen Arbeit in der Schriftart **Consolas** dargestellt.

### Erklärung zur Darstellung von Programmnamen, Befehlsketten und Quelltext

Die Namen von Programmen, die im Rahmen dieser Arbeit als *Befehle* im *Terminal* ausgeführt wurden, sowie deren *Parameter* sind ausnahmslos in der Schriftart **Consolas** dargestellt. Alle anderen Programme werden normal in der Schriftart des Texts dargestellt.

*Quelltext*beispiele sind ebenfalls ausnahmslos in der Schriftart **Consolas** dargestellt.

### Erklärung zur Darstellung spezieller Zeichen

Zeichen sind dann in der Schriftart *Consolas* dargestellt, wenn sie sich auf einen *Befehl* beziehen.

Sofern Zeichen, wie zum Beispiel Bindestriche, *Slashes* und ähnliche, im Alltagsgebrauch auch in Texten vorkommen können und im Zusammenhang der Arbeit fälschlicherweise mit besagtem Gebrauch assoziiert werden könnten, sind sie unter Anführungszeichen gesetzt.

### Erklärung zur Darstellung in verschiedenen Schriftgrößen

Ist es im Kontext dieser Arbeit erforderlich, dass *Befehlsketten* oder *Zeichenketten* in einer Zeile abgebildet werden, so ist die Schriftart, sofern notwendig, entsprechend verkleinert.

### Erklärung zur Angabe von allgemeiner Syntax

Die allgemeine *Syntax* der einzelnen *Tools* hält sich an folgende Regeln:

- Die *Zeichenfolge* „...“ bedeutet, dass an dieser Stelle mehrere *Optionen* nebeneinander stehen können.
- *Optionen*, die zwischen den Zeichen „[“ und „]“ und stehen, sind nicht zwingend erforderlich.

### Erklärung zum Terminal

In das *Terminal* eingegebene *Befehle* werden immer auf jene Dateien angewendet, welche im *Arbeitsverzeichnis* abgelegt sind. In anderen *Verzeichnissen* abgelegte Dateien mit dem gleichen Dateinamen sind davon nicht betroffen. Es wurde gänzlich mit *Bash*<sup>1</sup> gearbeitet.

### Erklärung zur Entwicklungsumgebung des Quelltexts

Mit Ausnahme des *Python*<sup>2</sup>-*Skripts* im Anhang (A.III), welches unter *Microsoft Windows XP*<sup>3</sup> erstellt wurde, wurden sämtliche Codezeilen dieser Arbeit unter *Ubuntu Linux 14.04.3 LTS*<sup>4</sup> erstellt. Die Installation einzelner Programme ist eventuell notwendig, wird im Rahmen dieser vorwissenschaftlichen Arbeit jedoch nicht erläutert.

---

<sup>1</sup> vgl. <https://www.gnu.org/software/bash/manual/bashref.html> (19.01.2016)

<sup>2</sup> vgl. <https://www.python.org/> (21.01.2016)

<sup>3</sup> vgl. <http://windows.microsoft.com/de-de/windows/windows-help#windows=windows-xp> (21.01.2016)

<sup>4</sup> vgl. <http://releases.ubuntu.com/14.04/> (21.01.2016)

# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung .....</b>	<b>10</b>
<b>2</b>	<b>Einführung in die Thematik.....</b>	<b>11</b>
2.1	Internetsurfen .....	11
2.1.1	Internetsurfen als Kommunikation.....	12
2.2	Transport von Webseiten gemäß HTTP .....	13
2.2.1	Funktionseinheiten des HTTP .....	13
2.2.2	Kommunikationsablauf aus Client-Sicht.....	14
2.2.3	HTTP-Request.....	14
2.2.4	HTTP-Response .....	15
2.2.5	HTTPS .....	16
2.2.6	Session-Management.....	16
2.2.7	Content-Type / MIME-Type .....	16
2.3	Woraus Webseiten bestehen .....	17
2.3.1	Hypertext Markup Language (HTML).....	17
2.3.2	Cascading Style Sheets (CSS) .....	19
2.3.3	JavaScript .....	19
2.4	Manueller Abruf und manuelle Datenextraktion .....	20
2.5	Automatisierter Abruf und automatisierte Datenextraktion .....	20
2.6	Reflexion über das Thema Webseiten und daraus zu ziehende Schlüsse.....	21
<b>3</b>	<b>Klassifizierung von Webseiten .....</b>	<b>22</b>
3.1	Statische Webseiten .....	22
3.1.1	Manuell aktualisierte statische Webseiten .....	23
3.1.2	Automatisch aktualisierte statische Webseiten.....	23
3.2	Dynamische Webseiten.....	23
3.2.1	Serverseitig generierte dynamische Webseiten .....	24
3.2.2	Clientseitig generierte dynamische Webseiten.....	25

---

3.2.3	Hybride dynamische Webseiten .....	26
3.2.4	Serverseitige Generierung versus clientseitige Generierung.....	27
3.2.5	Login-Bereiche.....	27
3.3	Relevanz der Art einer Webseite im Kontext dieser Arbeit .....	27
<b>4</b>	<b>Tools zum Abruf.....</b>	<b>28</b>
4.1	<b>wget</b> .....	28
4.1.1	Optionen von <b>wget</b> .....	29
4.1.2	Beispieleingaben zu <b>wget</b> .....	31
4.1.3	Anwendung von <b>wget</b> auf die verschiedenen Arten von Webseiten.....	31
4.2	<b>curl</b> .....	31
4.2.1	Optionen von <b>curl</b> .....	32
4.2.2	Beispieleingaben zu <b>curl</b> .....	32
4.2.3	Anwendung von <b>curl</b> auf die verschiedenen Arten von Webseiten.....	32
4.3	Selenium .....	33
4.3.1	Selenium IDE.....	33
4.3.2	Befehle von Selenium .....	34
4.3.3	Anwendung von Selenium auf die verschiedenen Arten von Webseiten .....	34
<b>5</b>	<b>Tools zur Verarbeitung.....</b>	<b>35</b>
5.1	GNU Coreutils und ähnliche .....	35
5.1.1	<b>cat</b> .....	35
5.1.2	<b>cut</b> .....	36
5.1.3	<b>grep</b> .....	37
5.1.4	<b>head</b> und <b>tail</b> .....	38
5.1.5	<b>sed</b> .....	39
5.1.6	<b>sort</b> .....	40
5.1.7	<b>tr</b> .....	42
5.1.8	Abschlussbemerkungen zu den GNU Coreutils.....	42



---

5.2	Parser .....	43
<b>6</b>	<b>Fallbeispiele .....</b>	<b>46</b>
6.1	Extraktion der Sprechstundenliste des BRG19 .....	46
6.2	Extraktion der eigenen Modulliste von elischa.brg19.at .....	50
<b>7</b>	<b>Bewertungen und Schlussfolgerungen.....</b>	<b>56</b>
7.1	Kenntnissen über das Internetsurfen und den Aufbau von Webseiten .....	56
7.2	Die Art einer Webseite.....	56
7.3	Sinnvolle Verwendung geeigneter Tools .....	57
7.4	Abschlussbewertung.....	57
<b>A</b>	<b>Anhang.....</b>	<b>58</b>
<b>B</b>	<b>Glossar .....</b>	<b>61</b>
<b>C</b>	<b>Literaturverzeichnis .....</b>	<b>70</b>
<b>D</b>	<b>Abbildungsverzeichnis .....</b>	<b>78</b>
<b>E</b>	<b>Tabellenverzeichnis .....</b>	<b>79</b>
<b>F</b>	<b>Quelltextverzeichnis .....</b>	<b>80</b>
<b>G</b>	<b>Beiliegendes Speichermedium .....</b>	<b>81</b>

# 1 Einleitung

Die Automatisierung bestimmter Arbeitsprozesse ist heutzutage essenziell, um ökonomisch mit Ressourcen haushalten zu können. Arbeitsprozesse, die automatisiert ablaufen, bedeuten eine große Zeit- und Geldersparnis. Die dadurch gewonnene Zeit kann bei Bedarf anderweitig verwendet werden. Das dadurch ersparte Geld eröffnet die Möglichkeit, in Güter oder Dienstleistungen zu investieren.

Diese vorwissenschaftliche Arbeit beschäftigt sich mit einem Teilbereich der Automatisierung, der automatisierten Extraktion von Daten aus *Webseiten*.

Diese Arbeit ist als hybride Variante aus Literaturarbeit und empirischer Arbeit anzusehen. Der Literaturteil befasst sich mit dem Analysieren und Vergleichen der einzelnen *Tools*. Der empirische Teil beschäftigt sich mit der Anwendung der einzelnen *Tools* im geeigneten Kontext. Dabei werden Fallbeispiele angeführt und praxisnahe in einem jeweils geeigneten Rahmen abgehandelt und analysiert.

Ferner werden das *Internetsurfen* im Allgemeinen und die verschiedenen Arten von *Webseiten* besprochen, wobei hier empirisch eine eigenständige Klassifizierung vorgenommen wird.

Ziel dieser Arbeit ist es, Schlussfolgerungen bezüglich der automatisierten Extraktion von Daten aus *Webseiten* zu ziehen. Dabei soll der gesamte Automatisierungsprozess beleuchtet werden. Die Weiterverarbeitung der gewonnenen Daten nach Abschluss des Extraktionsprozesses wird aufgrund des ansonsten überdimensionierten thematischen Rahmens nicht besprochen. Aus ebendiesem Grund ist es nicht Ziel dieser Arbeit, einen Überblick über die Entwicklungsgeschichte des *HTTP* oder jene der einzelnen *Tools* zu geben. Besagte Inhalte werden hinsichtlich jener Aspekte analysiert, die für die Arbeit relevant sind.

Als Quellen für diese Arbeit dienen hauptsächlich die offiziellen Dokumentationen und Spezifikationen der einzelnen *Tools* sowie die offiziellen *Webseiten* des *World Wide Web Consortiums (W3C)*<sup>5</sup> und der *Internet Engineering Task Force (IETF)*<sup>6</sup>.

---

<sup>5</sup> vgl. <http://www.w3.org/> (19.01.2016)

<sup>6</sup> vgl. <https://www.ietf.org/> (19.01.2016)

## 2 Einführung in die Thematik

Im Vorfeld dieser vorwissenschaftlichen Arbeit ist es notwendig, einige Begriffe und Prozesse hinsichtlich ihrer Bedeutung für diese Thematik zu erklären. Das folgende Kapitel zielt darauf ab, den Gesamtprozess und die dabei vonstattengehenden Vorgänge vom manuellen Abrufen von *Webseiten* beim *Internetsurfen* bis hin zur automatisierten Extraktion von Daten aus diesen zu veranschaulichen.

### 2.1 Internetsurfen

*Internetsurfen* bezeichnet die Tätigkeit der wahllosen oder gezielten Suche von Informationen im Internet. Voraussetzung für das *Internetsurfen* ist, neben einer Internetverbindung, ein Programm, welches die Benutzerin beziehungsweise den Benutzer befähigt, dieser Tätigkeit nachzugehen. Dieses Programm ist der *Webbrowser*, kurz *Browser*. Beispiele für *Browser* sind *Mozilla Firefox*<sup>7</sup>, *Google Chrome*<sup>8</sup>, *Windows Internet Explorer*<sup>9</sup>, *Microsoft Edge*<sup>10</sup>, *Safari*<sup>11</sup> und *Opera*<sup>12</sup>. Der *Browser* muss gestartet werden, um mit dem *Internetsurfen* beginnen zu können.

Beim *Internetsurfen* werden *Websites* abgerufen. *Websites* bestehen aus einer oder mehreren *Webseiten*.

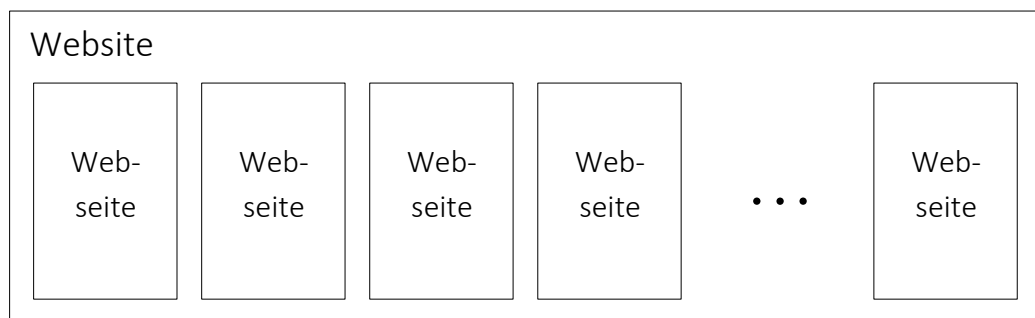


Abbildung 1: Eine Website besteht in der Regel aus mehreren Webseiten

---

<sup>7</sup> vgl. <https://www.mozilla.org/de/firefox/new/> (19.01.2016)

<sup>8</sup> vgl. <https://www.google.de/chrome/browser/desktop/index.html> (19.01.2016)

<sup>9</sup> vgl. <http://windows.microsoft.com/de-at/internet-explorer/download-ie> (19.01.2016)

<sup>10</sup> vgl. <https://www.microsoft.com/de-at/windows/microsoft-edge> (19.01.2016)

<sup>11</sup> vgl. <http://www.apple.com/at/safari/> (19.01.2016)

<sup>12</sup> vgl. <http://www.opera.com/de> (19.01.2016)

Da bei der automatisierten Extraktion von Daten aus *Webseiten* diese Unterscheidung allerdings irrelevant ist, da die Daten immer aus einzelnen *Webseiten* einer *Website* extrahiert werden, wird in dieser Arbeit lediglich der Begriff *Webseite* verwendet.

Um eine *Webseite* abzurufen, ist deren *Adresse* in die *Adresszeile* des *Browsers* einzugeben. Man kann aber auch nach einer *Google-Suche*<sup>13</sup> auf eines der Ergebnisse klicken, um zu einer *Webseite* zu gelangen. Die entsprechende *Webseite* wird in beiden Fällen auf dem Bildschirm ausgegeben.

Die obig angesprochene *Adresse* hält sich an ein Schema, welches *Uniform Resource Locator (URL)*<sup>14</sup> genannt wird. Eine *URL* gibt an, welche Daten<sup>15</sup> von welchem Computer unter Verwendung des angegebenen *Protokolls*<sup>16</sup> übertragen werden sollen, wobei hierbei unter *Protokoll* eine Kommunikationsvorschrift verstanden wird, die vorschreibt, wie eine Kommunikation abzulaufen hat.

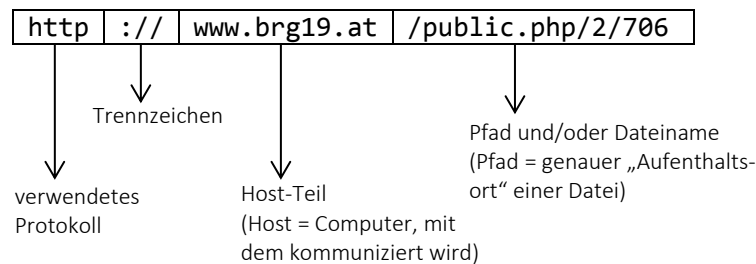


Abbildung 2: Ein konkretes Beispiel für eine URL

Gibt man keinen *Pfad* oder Dateinamen an, gelangt man je nach *Servereinstellung* zur Startseite.

### 2.1.1 Internetsurfen als Kommunikation

Abschließend ist zu sagen, dass das *Internetsurfen* als eine Kommunikation zwischen zwei Instanzen anzusehen ist: der Person, die eine *Webseite* abrufen will, und jener, die diese anbietet. Diese Darstellung ist jedoch vereinfacht und wird im Folgenden konkretisiert.

<sup>13</sup> vgl. <https://www.google.com> (19.01.2016)

<sup>14</sup> vgl. <https://tools.ietf.org/html/rfc1738> (19.01.2016)

<sup>15</sup> vgl. <https://tools.ietf.org/html/rfc1738#section-2> (19.01.2016)

<sup>16</sup> vgl. <https://tools.ietf.org/html/rfc1738#section-2.1> (19.01.2016)

## 2.2 Transport von Webseiten gemäß HTTP

Die obig angesprochene Kommunikation läuft, wenn in der *URL* `http://` vorangestellt ist, gemäß dem *Hypertext Transfer Protocol (HTTP)* ab.

Das *HTTP* ist keineswegs das einzige *Protokoll*, mit dem Inhalte über das Internet übertragen werden können. Weitere Beispiele für *Protokolle* sind das *File Transfer Protocol (FTP)*<sup>17</sup>, das *Simple Mail Transfer Protocol (SMTP)*<sup>18</sup>, das *Post Office Protocol (POP)*<sup>19</sup>, das *Internet Message Access Protocol (IMAP)*<sup>20</sup>, das *Lightweight Directory Access Protocol (LDAP)*<sup>21</sup> und das *Network Time Protocol (NTP)*<sup>22</sup>, um nur einige zu nennen.

Aufgrund seiner weiten Verbreitung beim *Internetsurfen* ist das *HTTP* allerdings das für diese Arbeit relevante *Protokoll*.

Das *HTTP* ist ein *zustandsloses Protokoll* zur Übertragung von Daten über ein *Rechnernetz*. *Zustandslos* bedeutet in diesem Fall, dass alle *Anfragen* voneinander unabhängig behandelt werden.<sup>23</sup>

### 2.2.1 Funktionseinheiten des HTTP

Das *HTTP* hat verschiedene Funktionseinheiten, welche im *Request for Comments (RFC) 2616* in Abschnitt 1.3<sup>24</sup> spezifiziert sind.

Eine Kommunikation gemäß dem *HTTP* läuft immer zumindest zwischen den zwei Funktionseinheiten *Client* und *Webserver*, oft einfach *Server*, ab:

- Der *Client* initiiert Verbindungen zu *Servern* und ist in der Regel ein *Browser*.
- Der *Server* wartet auf Verbindungsanfragen von *Clients* und antwortet auf diese durch das Liefern der gewünschten Ressource oder Statusinformation.

---

<sup>17</sup> vgl. <https://tools.ietf.org/html/rfc354> (19.01.2016)

<sup>18</sup> vgl. <https://tools.ietf.org/html/rfc5321> (19.01.2016)

<sup>19</sup> vgl. <https://tools.ietf.org/html/rfc1939> (19.01.2016)

<sup>20</sup> vgl. <https://tools.ietf.org/html/rfc3501> (19.01.2016)

<sup>21</sup> vgl. <https://tools.ietf.org/html/rfc4511> (19.01.2016)

<sup>22</sup> vgl. <https://tools.ietf.org/html/rfc5905> (19.01.2016)

<sup>23</sup> vgl. <https://tools.ietf.org/html/rfc2616>, Abstract (19.01.2016)

<sup>24</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-1.3> (19.01.2016)

Folgende Grafik veranschaulicht den Ablauf einer Kommunikation gemäß dem *HTTP* auf Basis der genannten Funktionseinheiten:

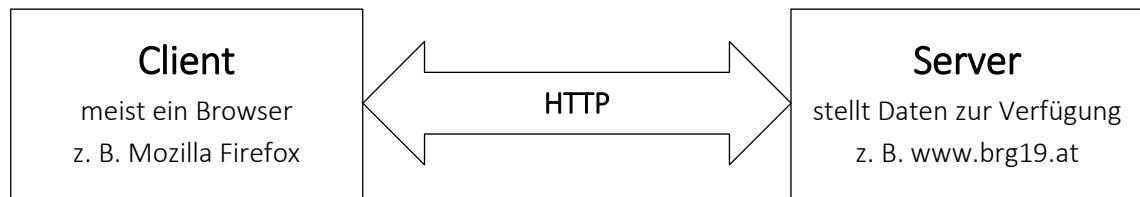


Abbildung 3: Kommunikation zwischen Client und Server gemäß HTTP (1)

### 2.2.2 Kommunikationsablauf aus Client-Sicht

Die Kommunikation läuft bei einer *HTTP*-Verbindung immer nach einem bestimmten Schema ab. Aus *Client*-Sicht lässt sich der typische Kommunikationsablauf wie folgt aufschlüsseln:

- (1) Der *Client* öffnet eine Verbindung<sup>25</sup> zum *Server*.
- (2) *HTTP-Request*: Der *Client* sendet eine *Anfrage* an den *Server*.
- (3) Der *Server* erhält die *Anfrage*.
- (4) *HTTP-Response*: Der *Server* antwortet auf die *Anfrage*.
- (5) [Beenden der Verbindung]

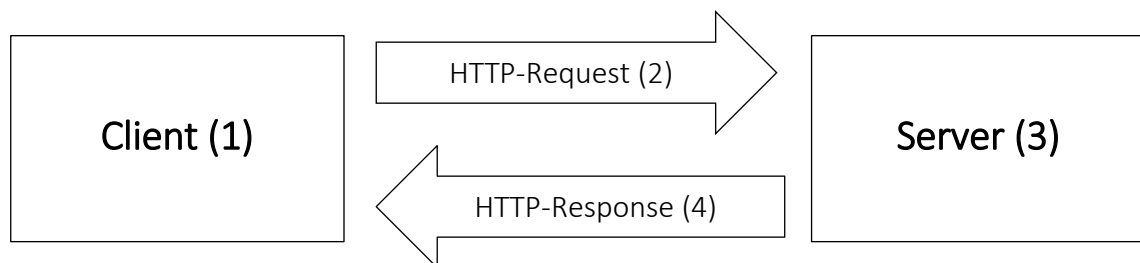


Abbildung 4: Kommunikation zwischen Client und Server gemäß HTTP (2)

### 2.2.3 HTTP-Request

Ein *Request (Anfrage)* wird in Form einer *Request-Message*<sup>26</sup> (*Anfragenachricht*) vom *Client* zum *Server* ausgeschildt. In der ersten Zeile der Nachricht sind die auf die Ressource

<sup>25</sup> vgl. <https://www.rfc-editor.org/rfc/rfc793.txt> (19.01.2016)

<sup>26</sup> vgl. <http://www.w3.org/Protocols/rfc2616/rfc2616-sec5.html> (19.01.2016)

angewendete Methode<sup>27</sup>, der Bezeichner der Ressource (*Request-URL*<sup>28</sup>) und die *Protokollversion* angeführt.

GET	/index.php	HTTP/1.1
Methode	Request-URL	Protokollversion

Abbildung 5: Ein konkretes Beispiel für eine Request-Zeile

Das *HTTP* bietet verschiedene Methoden einer *Request-Message* an. Für die automatisierte Datenextraktion sind in erster Linie folgende relevant:

*GET*<sup>29</sup>: Diese Methode sendet eine *Anfrage* nach der in der *Request-URL* angegebenen Ressource. *Client*-Daten werden als Bestandteil der *URL* der Ressource übergeben. Diese Methode wird für jene Informationen verwendet, die offenkundig in der *URL* stehen dürfen. Außerdem wird der Benutzerin beziehungsweise dem Benutzer auf diese Art und Weise die Möglichkeit geboten, *Parameter* direkt in der *Adresszeile* des *Browsers* in der *URL* auszubessern.

*POST*<sup>30</sup>: Diese Methode funktioniert wie *GET*, allerdings werden dabei *Client*-Daten nicht an die *URL* angehängt, sondern im Verborgenen übermittelt. Diese Methode ist für sensible Informationen wie zum Beispiel Passwörter wichtig.

## 2.2.4 HTTP-Response

Nach dem Erhalten und Interpretieren einer *Request-Message* antwortet der *Server* mit einem *Response* (*Antwort*) in Form einer *Response-Message*<sup>31</sup> (*Antwortnachricht*).

Die erste Zeile einer *Antwortnachricht* ist die Statuszeile, bestehend aus der *Protokollversion*, gefolgt von einem numerischen Statuscode und der ihm zugeordneten Phrase.<sup>32</sup>

HTTP/1.1	200	OK
Protokollversion	numerischer Statuscode <sup>33</sup>	zugehörige Phrase <sup>34</sup>

Abbildung 6: Ein konkretes Beispiel für eine Statuszeile

<sup>27</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-5.1.1> (19.01.2016)

<sup>28</sup> vgl. <http://tools.ietf.org/html/rfc3986#section-1.1.3> (19.01.2016)

<sup>29</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-9.3> (19.01.2016)

<sup>30</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-9.5> (19.01.2016)

<sup>31</sup> vgl. <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)

<sup>32</sup> vgl. <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)

<sup>33</sup> vgl. <http://www.iana.org/assignments/http-status-codes/http-status-codes.xhtml> (19.01.2016)

<sup>34</sup> vgl. <http://www.iana.org/assignments/http-status-codes/http-status-codes.xhtml> (19.01.2016)

### 2.2.5 HTTPS

Die verschlüsselte Variante des *HTTP* heißt *HTTPS*, wobei das „S“ für „secure“ (sicher) steht. Diese ist im *RFC 2818*<sup>35</sup> spezifiziert. Durch *HTTPS* ist eine sichere Übertragung von sensiblen Daten<sup>36</sup> wie beispielsweise Passwörtern möglich.

### 2.2.6 Session-Management

Da das *HTTP*, wie in der Einleitung bereits erwähnt, ein *zustandsloses Protokoll* ist, bezieht es keine Information aus einer Kommunikation, die bereits in der Vergangenheit stattgefunden hat.<sup>37</sup>

In diesem Zusammenhang ist der Begriff *Session (Sitzung)* zu erwähnen. Eine *Session* beschreibt einen Dialog, welcher sich über mehrere Zyklen von *Request* und *Response* erstreckt. Weiter soll innerhalb einer *Session* der Rückbezug auf die vergangene Kommunikation möglich sein.<sup>38</sup>

Zur Codierung von *Session*-Information werden verschiedene Techniken<sup>39</sup> verwendet, die das *HTTP* unterstützt:

- Die *Session*-Information wird beim *URL Rewriting* in der *URL* codiert.
- Die *Session*-Information wird in Form von *Cookies*<sup>40</sup> beim *Client* gespeichert.
- Die *Session*-Information wird in sogenannten *Hidden Fields*, unsichtbaren Formularfeldern, untergebracht.

### 2.2.7 Content-Type / MIME-Type

Unter Verwendung des *Protokolls HTTP* können verschiedene Typen von Daten übertragen werden. Die sogenannten *Content-* beziehungsweise *MIME-Types (Medientypen)*<sup>41</sup> geben an, um welche Art von Daten es sich handelt. Sie bestehen aus einem *Top-Level-Typ* und einem *Untertyp* und können durch *Parameter* weiter spezifiziert werden.

---

<sup>35</sup> vgl. <https://tools.ietf.org/html/rfc2818> (19.01.2016)

<sup>36</sup> vgl. <https://tools.ietf.org/html/rfc2818#section-1> (19.01.2016)

<sup>37</sup> vgl. <http://tools.ietf.org/html/rfc7230#section-2.3> (19.01.2016)

<sup>38</sup> vgl. <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 40 (19.01.2016)

<sup>39</sup> vgl. <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 42 (19.01.2016)

<sup>40</sup> vgl. <https://tools.ietf.org/html/rfc6265> (19.01.2016)

<sup>41</sup> vgl. <http://www.w3.org/2001/tag/2002/0129-mime> (19.01.2016)



Folgende Tabelle veranschaulicht einige Arten der übertragenen Daten:

Top-Level-Typ	Untertyp	Beschreibung
text	<i>plain</i> <i>enriched</i>	unformatierter <i>ASCII-Text</i> <i>ASCII-Text</i> mit einfachen Formatierungen
image	gif jpeg	Standbild im <i>GIF-Format</i> Standbild im <i>JPEG-Format</i>
audio	basic	Klangdaten
application	postscript	druckbares Dokument im <i>PostScript-Format</i>
...	...	...

Tabelle 1: Auszug aus einer Übersicht der erlaubten MIME-Types<sup>42</sup>

## 2.3 Woraus Webseiten bestehen

Da gemäß dem *HTTP* verschiedene *MIME-Types* übertragen werden können, steht der eine *Webseite* anbietenden Person ein breites Feld an Möglichkeiten zur Verfügung, Daten bereitzustellen.

### 2.3.1 Hypertext Markup Language (HTML)

Sehr viele *Webseiten* sind aus *Hypertext* aufgebaut. *Hypertext* ist ein über *Links* verbundenes Netz aus Text-, Bild- und Dateneinheiten. *Hypertext* wird durch die *Hypertext Markup Language (HTML)* ausgezeichnet.

*HTML* ist eine Auszeichnungssprache zum Beschreiben der Struktur von *Webseiten*. *HTML* ermöglicht Entwicklerinnen und Entwicklern Online-Dokumente mit Überschriften, Text, Tabellen, Listen, Fotos, Formularen, Videoclips, Soundclips et cetera zu veröffentlichen. Weiter ermöglicht *HTML* durch klickbare *Links* zwischen einzelnen *Webseiten* zu navigieren.<sup>43</sup>

<sup>42</sup> vgl. <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)

<sup>43</sup> vgl. <http://www.w3.org/standards/webdesign/htmlcss#whathtml> (19.01.2016)

Ein Beispiel für ein einfaches *HTML*-Dokument, zeilenweise aufgegliedert:

1	<!DOCTYPE html PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
2	<html>
3	<head>
4	<title>Titel der Webseite</title>
5	</head>
6	<body>
7	<h1>Überschrift erster Ebene</h1>
8	<p>einfacher Text</p>
9	</body>
10	</html>

Quelltext 1: Aufbau eines einfachen HTML-Dokuments

Ein *HTML*-Dokument besteht aus drei Teilen:<sup>44</sup>

1. einer Zeile, die Informationen über die *HTML*-Version enthält (Zeile 1)
2. dem *HTML*-Kopf, der Informationen über das Dokument, sogenannte *Metadaten*, bereitstellt, die in der Regel nicht dargestellt werden (Zeilen 3 bis 5)
3. dem *HTML*-Körper, der den Inhalt enthält, welcher dargestellt wird (Zeilen 6 bis 9)

Der *HTML*-Kopf und der *HTML*-Körper sind von den *Tags* `<html>` und `</html>` umschlossen.

*Tags* zeichnen Elemente aus. Sie können Elemente beispielsweise als Überschrift oder *Link* auszeichnen. Hierbei wird zwischen öffnenden und schließenden *Tags* unterschieden:

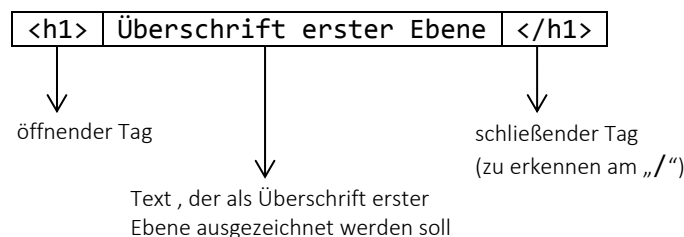


Abbildung 7: Öffnende und schließende Tags in HTML

Diese *Tags* lassen sich nach der entsprechenden *Dokumenttypdefinition* verschachteln<sup>45</sup>, dadurch wird beispielsweise die *Verlinkung* einer Überschrift ermöglicht:

```
<h1><a href="Zieladresse">Überschrift</a></h1>
```

<sup>44</sup> vgl. <http://www.w3.org/TR/html401/struct/global.html#h-7.1> (19.01.2016)

<sup>45</sup> vgl. <http://www.w3.org/TR/html401/sgml/dtd.html> (19.01.2016)

Innerhalb der *Tags* `<h1>` und `</h1>` befindet sich jener Inhalt, der als Überschrift formatiert werden soll. Innerhalb von `<a>` und `</a>` steht der *Linktext*, also jener Text, der auf dem Bildschirm dargestellt wird, wobei `href` ein *Parameter* des *Link-Tags* ist, der die *Zieladresse* angibt.

### 2.3.2 Cascading Style Sheets (CSS)

Während *HTML* die Struktur einer *Webseite* deklariert, gibt *Cascading Style Sheets (CSS)* die Darstellung einer *Webseite* an.<sup>46</sup>

Durch *CSS* können Schriftart, Größe und Farbe eines Texts verändert, sowie sonstige Formatierungen vorgenommen werden.<sup>47</sup>

*CSS* kann direkt in einem *HTML*-Dokument verwendet, oder in ein separates Dokument ausgelagert werden.<sup>48</sup>

Der Unterschied zwischen *HTML* und *CSS* ist wie folgt zu veranschaulichen:

Struktur	Darstellung
<i>HTML</i>	<i>CSS</i>
Das Wort „VwA“ ist eine Überschrift	Überschriften sind rot

Tabelle 2: HTML und CSS haben unterschiedliche Aufgaben

Sofern keine *CSS*-Informationen angegeben beziehungsweise verfügbar sind, wird eine *Webseite* wie in den Einstellungen des *Browsers* festgelegt dargestellt.<sup>49</sup>

### 2.3.3 JavaScript

*HTML* und *CSS* sind keine Programmiersprachen, *JavaScript* ist eine. *JavaScript* bietet die Möglichkeit, Webanwendungen zu realisieren. Dies ist durch *HTML* und *CSS* allein nicht möglich. Viele *Browser* haben ihre eigene *JavaScript*-Spezifikation<sup>50</sup>, daher kann es zu Kompatibilitätsproblemen kommen. Der Sprachkern wurde von der *Ecma International (Ecma)* in der *ECMA 262* als *ECMAScript* spezifiziert.<sup>51</sup>

<sup>46</sup> vgl. <http://www.w3.org/standards/webdesign/htmlcss#whatcss> (19.01.2016)

<sup>47</sup> vgl. <http://www.w3.org/TR/CSS1/#abstract> (19.01.2016)

<sup>48</sup> vgl. <http://www.w3.org/TR/CSS21/intro.html#html-tutorial> (19.01.2016)

<sup>49</sup> vgl. u. a. <https://support.mozilla.org/de/kb/Einstellungen-Fenster--Inhalts-Abschnitt> (19.01.2016)

<sup>50</sup> vgl. u. a. <https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference> (19.01.2016)

<sup>51</sup> vgl. <http://ecma-international.org/ecma-262/5.1/> (19.01.2016)

## 2.4 Manueller Abruf und manuelle Datenextraktion

Damit ist geklärt, was beim *Internetsurfen* passiert, sprich wie die Kommunikation zwischen *Client* und *Server* abläuft, wie *Webseiten* transportiert werden und woraus sie bestehen. Dies hat allerdings grundsätzlich noch nicht direkt etwas mit Automatisierung zu tun, da es beim *Internetsurfen* immer eine Person ist, welche die *URL* in die *Adresszeile* des *Browsers* eingibt oder auf ein Suchergebnis klickt. Das Abrufen einer *Webseite* geht also manuell vonstatten. Ein Beispiel hierfür ist das manuelle Abrufen von Börsenkursen oder Sportergebnissen. Eine Person besucht eine *Webseite*, um diese Information zu erhalten und bei Bedarf zu extrahieren. Die Extraktion der Daten erfolgt ebenfalls auf manuellem Wege, beispielsweise durch das Aufschreiben auf ein Blatt Papier oder durch die Eingabe in beispielsweise ein *Microsoft Excel*<sup>52</sup> Dokument.

Aufgrund des hohen Zeitaufwandes, der Fehleranfälligkeit oder anderer Faktoren kann diese Art der Datenextraktion auf Dauer auf vielen Ebenen eventuell nicht rentabel sein und den Bedarf einer Automatisierung hervorrufen.

## 2.5 Automatisierter Abruf und automatisierte Datenextraktion

Um den obig angesprochenen Prozess zu automatisieren, müssen sowohl Abruf als auch Datenextraktion automatisiert werden. Der Gesamtprozess dieser Automatisierung sieht hierbei wie folgt aus („3. Datenverarbeitung“ wird, wie in der Einleitung bereits erwähnt, im Rahmen dieser Arbeit nicht behandelt):



Abbildung 8: Gesamtprozesses der automatisierten Datenextraktion aus Webseiten

Für den automatisierten Abruf und die automatisierte Datenextraktion werden *Werkzeuge*, auch *Tools* genannt, benötigt.

Das *Tool* zum Abruf kann entweder der *Browser* selbst oder ein speziell dafür angefertigtes sein. Von den bereits angesprochenen *Browsern* gibt es zwei verschiedene Arten: *Browser* mit grafischer Benutzeroberfläche und jene, die textbasiert sind. Die in Kapitel 2.1 bereits

<sup>52</sup> vgl. <https://products.office.com/de-at/excel> (19.01.2016)

genannten sind allesamt *Browser* mit grafischer Benutzeroberfläche. *Browser* mit grafischer Benutzeroberfläche stellen eine *Webseite* grafisch dar, textbasierte in Textform. Ein Beispiel für einen textbasierten *Browser* ist lynx<sup>53</sup>.

## 2.6 Reflexion über das Thema Webseiten und daraus zu ziehende Schlüsse

Bei der Reflexion über das Thema *Webseiten* kommt man zum Schluss, dass diese auf interaktive Bedienung und nicht auf die automatische Verarbeitung, welche eine nicht interaktive Anwendung darstellt, ausgelegt sind. *Webseiten* haben Bedienelemente wie Navigationsleisten und Menüs, eine entsprechende grafische Aufbereitung und ähnliche Elemente, die mit der eigentlichen Information, welche die *Webseite* transportieren möchte, teilweise wenig zu tun haben.

*Browser* als Mittel zum *Internetsurfen* sind demnach auch primär auf interaktive Bedienung ausgelegt, weshalb der automatische Abruf der Daten ohne das Zutun einer Benutzerin beziehungsweise eines Benutzers kein vorgesehener Anwendungsfall ist.

Der textbasierte *Browser* eröffnet in diesem Zusammenhang eine Möglichkeit, da er die Daten in Textform extrahieren kann. Diese Vorgehensweise klappt allerdings bei unterschiedlichen *Webseiten* unterschiedlich gut. Wie performant diese Methode ist, hängt von der Technologie, die hinter einer *Webseite* steckt, ab. Andere *Tools* verhalten sich ähnlich. *Webseiten* sind also bezüglich ihres technologischen Hintergrunds, der für den *Client* zumindest teilweise sichtbar ist, einzuteilen.

---

<sup>53</sup> vgl. <http://lynx.invisible-island.net/> (19.01.2016)

### 3 Klassifizierung von Webseiten

Die Technologie von *Webseiten* betreffend werden verschiedene Entwicklungsparadigmen verfolgt, bezüglich welcher *Webseiten* im Folgenden klassifiziert werden.

Die in dieser Arbeit getroffene Klassifizierung sieht wie folgt aus:

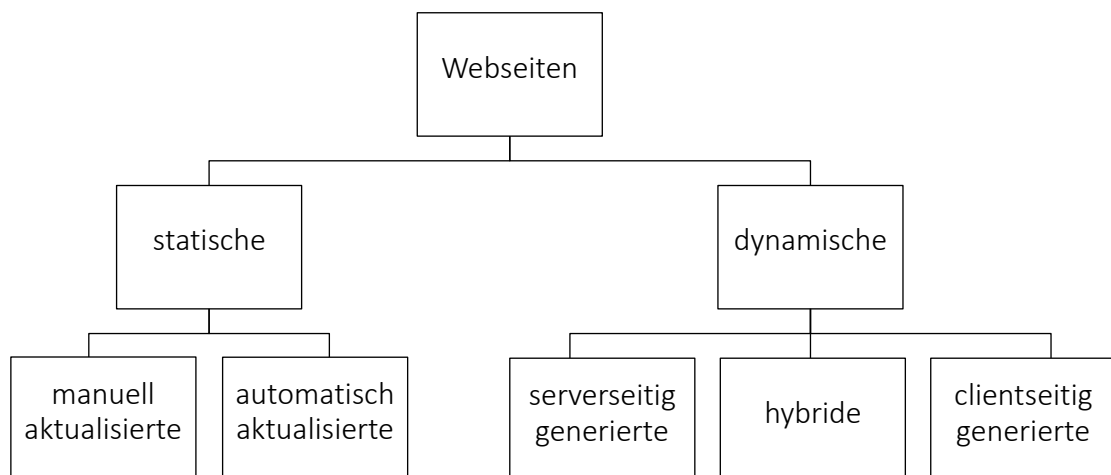


Abbildung 9: Einteilung der Arten von Webseiten

Die erste Unterscheidung, die getroffen werden muss, ist somit jene zwischen *statischen* und *dynamischen Webseiten*:

#### 3.1 Statische Webseiten

*Statische Webseiten* zeichnen sich dadurch aus, dass der *Client* beim Abruf exakt die gleiche Datei erhält, die auf dem *Server* gespeichert ist. Der *Client* holt also statische Dateien ab.

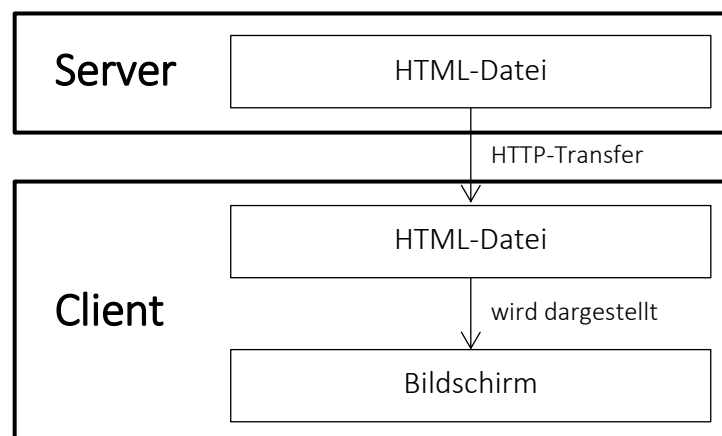


Abbildung 10: Statische Webseiten

Bei *statischen Webseiten* wird zwischen zwei Arten differenziert:

### 3.1.1 Manuell aktualisierte statische Webseiten

Bei *manuell aktualisierten statischen Webseiten* liegen die Dateien auf dem *Server* und werden von der sie anbietenden Person manuell aktualisiert.

Ein einfaches Beispiel hierfür ist die persönliche *Webseite* einer Person, oftmals als *Homepage* bezeichnet. Persönliche *Homepages* zielen meist darauf ab, die sie betreibende Person zu präsentieren. Die Beschreibung von Hobbys, persönlichen Vorlieben, sowie persönliche Daten, wie der eigene Geburtstag oder die E-Mail-Adresse, aber auch Neuigkeiten aus dem Leben der Person sind Beispiele für Inhalte dieser Art. Bei besagten Inhalten muss die die Webseite anbietende Person selbst über die Aktualisierung entscheiden.

### 3.1.2 Automatisch aktualisierte statische Webseiten

Jenes Merkmal, das die Art der *automatisch aktualisierten statischen Webseiten* auszeichnet, ist, dass die statischen Dateien automatisch aktualisiert werden.

Es gibt neben den obig bereits angesprochenen, über einen längeren Zeitraum gültigen Inhalten auch solche, die nur über einen kürzeren Zeitraum hinweg gültig sind. Das heutige, durch Tag, Monat und Jahr angegebene Datum ist ein Beispiel für einen derartigen Inhalt. Es ist in der Regel für 24 Stunden gültig. Ein weiteres Beispiel für Inhalte dieser Art ist die Wartezeit auf der Anzeigetafel einer Wiener Busstation, deren Daten nur eine Gültigkeitsdauer im Minutenbereich aufweisen.

Eine Möglichkeit, derartige Inhalte in eine *Webseite* einzuspeisen, besteht darin, die sie beinhaltenden Dateien in passenden Zeitabständen automatisch zu aktualisieren.

Dies kann im Fall des Datums beispielsweise durch ein Programm geschehen, welches die das Datum beinhaltende Datei automatisch aktualisiert und diese um 00:00 Uhr jede Nacht auf den *Server* hochlädt oder diese direkt auf dem *Server* manipuliert.

Die Wahl des Aktualisierungszeitpunkts hängt also vom Inhalt ab.

## 3.2 Dynamische Webseiten

*Dynamische Webseiten* hingegen zeichnen sich dadurch aus, dass der *Client* nicht exakt die gleiche Datei erhält, die auf dem *Server* gespeichert ist, sondern eine Datei, welche erst beim Abruf aus jenen Anweisungen, die in der abgerufenen Datei festgelegt sind, generiert wird.

Hierbei wird zwischen zwei unterschiedlichen Paradigmen differenziert: der *serverseitigen* sowie der *clientseitigen* Generierung. Diese werden im Folgenden beschrieben.

Festzuhalten ist, dass die beiden angesprochenen Paradigmen einander nicht gegenseitig ausschließen. Eine *serverseitig generierte dynamische Webseite* kann *clientseitig* generierte Elemente enthalten und umgekehrt.

### 3.2.1 Serverseitig generierte dynamische Webseiten

*Serverseitig generierte dynamische Webseiten* basieren zumeist<sup>54</sup> auf der Programmiersprache *PHP Hypertext Processor (PHP)*<sup>55</sup>, oder auch auf *Active Server Pages .NET (ASP.NET)*<sup>56</sup> oder *JavaServer Pages (JSP)*<sup>57</sup>. Der grundsätzliche Ablauf des Abrufs einer *serverseitig generierten Webseite* lässt sich wie folgt beschreiben:

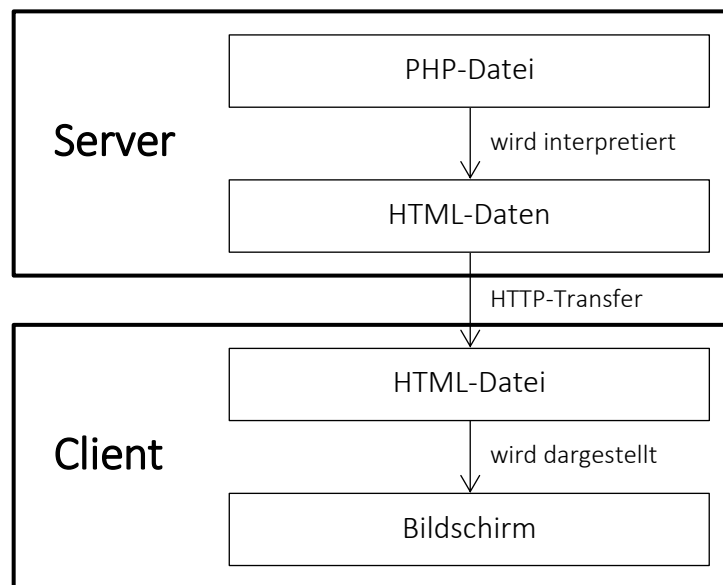


Abbildung 11: Serverseitig generierte dynamische Webseiten

Auf dem *Server* befindet sich eine Datei, die mit Instruktionen versehen ist, denen zufolge beim Abruf eine Datei generiert und an den *Client* geliefert wird. Ein Beispiel hierfür ist die bereits in Kapitel 3.1 angesprochene Ausgabe des heutigen Datums. Bei *serverseitig generierten dynamischen Webseiten* steht in der abzurufenden Datei ein *Befehl*, der veranlasst, dass vom *Server* oder einer anderen Quelle das heutige Datum abgeholt und dann beim *Client* abgeliefert wird.

<sup>54</sup> vgl. <http://redmonk.com/sogrady/2013/07/25/language-rankings-6-13/> (19.01.2016)

<sup>55</sup> vgl. <http://php.net/> (19.01.2016)

<sup>56</sup> vgl. <http://www.asp.net/> (19.01.2016)

<sup>57</sup> vgl. <http://www.oracle.com/technetwork/java/javaee/jsp/index.html> (19.01.2016)



Der *Quelltext* bleibt für den *Client* im Verborgenen, da er auf dem *Server* ausgeführt wird und lediglich die Ausgabe dieser Ausführung an den *Client* geschickt wird.<sup>58</sup>

Anzumerken ist, dass die Interpretation von im *Quelltext* enthaltenen *Befehlen browserunabhängig* vonstattengeht, da sie von der auf dem *Server* installierten Software abhängt.

### 3.2.2 Clientseitig generierte dynamische Webseiten

*Clientseitig generierte dynamische Webseiten* basieren zumeist<sup>59</sup> auf der Sprache *JavaScript*<sup>60</sup>. Der grundsätzliche Ablauf bei der *clientseitigen* Generierung lässt sich wie im Folgenden dargelegt beschreiben:

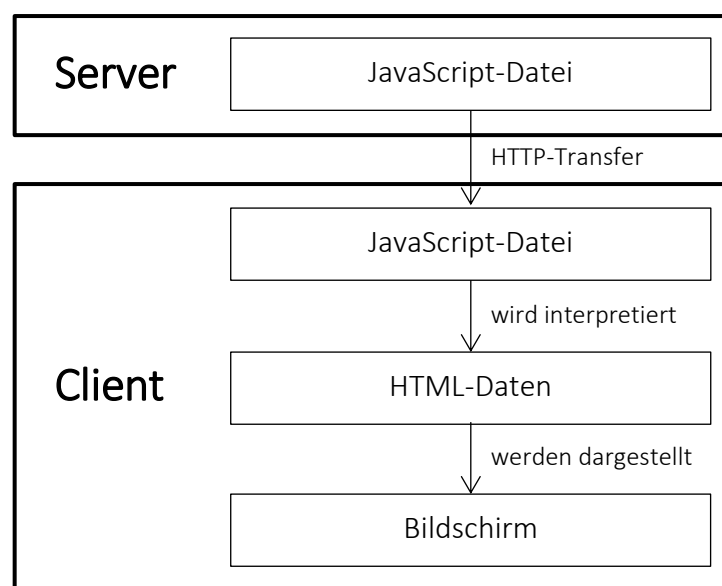


Abbildung 12: Clientseitig generierte dynamische Webseiten

Auf dem *Server* befindet sich eine Datei, die steuert, wie die *Webseite* schlussendlich aussieht. Zum Beispiel kann der *Server* dem *Client* mitteilen, dass er eine Überschrift erstellen soll, die einen bestimmten Text oder auch einen bestimmten Button beinhaltet. Die Generierung wird hierbei zwar durch den *Server* gesteuert, sie erfolgt aber erst beim *Client*. Der Client erhält also keine fertige Webseite, sie wird erst von diesem gefertigt. Hierfür muss dem *Client* der *Quelltext* bekannt sein, da der *Browser* und nicht der *Server* die Instanz ist, welche die darin enthaltenen *Befehle* ausführt. Der *Quelltext* kann von einer Person ausge-

<sup>58</sup> vgl. <http://php.net/manual/de/intro-what-is.php> (19.01.2016)

<sup>59</sup> vgl. <http://redmonk.com/sograde/2013/07/25/language-rankings-6-13/> (19.01.2016)

<sup>60</sup> vgl. <http://www.ecma-international.org/publications/files/ECMA-ST/Ecma-262.pdf> (19.01.2016) und <https://developer.mozilla.org/en-US/docs/Web/JavaScript> (19.01.2016)

lesen werden, da die ausführende Instanz von ihr selbst kontrolliert wird, was beim *Server* nicht der Fall ist.

Anzumerken ist, dass die Interpretation von *JavaScript-Befehlen browserabhängig* vonstattengeht, da manche *Browser* eine eigene Spezifikation von *JavaScript* verwenden.<sup>61</sup> Ob sich diese an die standardisierte *ECMA-International-Spezifikation*<sup>62</sup> hält, liegt im Ermessen der Entwicklerinnen und Entwickler des *Browsers*.

### 3.2.3 Hybride dynamische Webseiten

*Hybride dynamische Webseiten* sind eine Mischform aus *serverseitig generierten dynamischen Webseiten* und *clientseitig generierten dynamischen Webseiten*. Ein beispielhafter Ablauf einer hybriden Generierung lässt sich wie folgt darstellen:

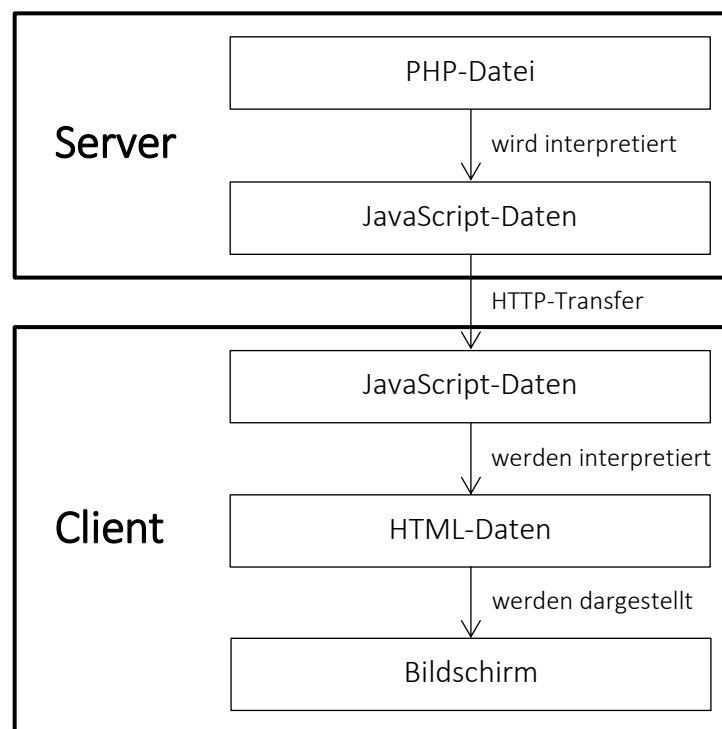


Abbildung 13: Hybride dynamische Webseiten

*Hybride dynamische Webseiten* sind *Webseiten*, die sowohl Elemente einer *serverseitig generierten dynamischen Webseite* als auch welche einer *clientseitig generierten dynamischen Webseite* enthalten. Die Generierung wird auch hierbei durch den *Server* gesteuert. Sie

<sup>61</sup> vgl. <https://developer.mozilla.org/en-US/docs/Web/JavaScript> (19.01.2016)

<sup>62</sup> vgl. <http://www.ecma-international.org/publications/files/ECMA-ST/Ecma-262.pdf> (19.01.2016)

erfolgt teilweise auf dem *Server*, teilweise beim *Client*. Die Zugänglichkeit des *Quelltexts* ist elementabhängig.

### 3.2.4 Serverseitige Generierung versus clientseitige Generierung

*Serverseitig generierte dynamische Webseiten* zeichnen sich zusätzlich zu den bereits angesprochenen Merkmalen dadurch aus, dass die Rechenlast hauptsächlich beim *Server* liegt. Je komplexer die auszuführenden Operationen sind, desto ausgelasteter und folglich desto weniger performant ist der *Server*.

Bezüglich der *clientseitig generierten dynamischen Webseiten* ist hingegen zu bemerken, dass die Rechenlast größtenteils beim *Client* liegt. Bei schwächeren Computern wie zum Beispiel Smartphones kann dies zu starken Einbußen in Leistung und vor allem im Komfort der Bedienung führen.

### 3.2.5 Login-Bereiche

Ergänzend ist zu bemerken, dass geschützte Bereiche einer *Webseite*, sogenannte *Login-Bereiche*, eine *serverseitige* Verarbeitung erfordern.

Beispiele für *Webseiten* mit *Login-Bereichen* sind [elischa.brg19.at](http://elischa.brg19.at), Twitter, Facebook, YouTube und etliche mehr.

Das *Login* erfolgt dabei oftmals durch die Eingabe eines Account-Namens und des dazugehörigen Passworts.

Ein solcher *Login-Bereich* stellt besondere Anforderungen an ein Automatisierungstool.

## 3.3 Relevanz der Art einer Webseite im Kontext dieser Arbeit

Bei Wahl eines geeigneten *Tools* muss, neben den obig genannten besonderen Anforderungen, vor allem auf die Art einer *Webseite* Rücksicht genommen werden.

Die im weiteren Verlauf der Arbeit genauer besprochenen *Tools* haben jeweils spezifische Anwendungsgebiete und eignen sich daher unterschiedlich gut für die verschiedenen Arten von *Webseiten*.

Für die Anwenderin beziehungsweise den Anwender ist es daher wichtig, bestmöglich über die *Webseite*, auf die besagte *Tools* angewendet werden sollen, Bescheid zu wissen, da dadurch eine möglichst genau auf die betreffende *Webseite* zugeschnittene Anwendung erfolgen kann.

## 4 Tools zum Abruf

Wie de facto alles in der Informatik, gehorcht auch die automatisierte Datenextraktion dem Grundprinzip „Eingabe → Verarbeitung → Ausgabe“, respektive einem linearen Ablauf. Darauf basierend lässt sich eine Einteilung in Kategorien vornehmen, in welcher „Tools zum Abruf“ die Eingabe repräsentiert, da durch den Abruf jene Daten gewonnen werden, welche in die für die Verarbeitung zuständigen Programme eingespeist werden. Anzumerken ist, dass die Klassifizierung hinsichtlich der Grundfunktionalität der einzelnen *Tools* getroffen worden ist.

### 4.1 wget

Das erste besprochene *Tool* zum Abruf ist **wget**<sup>63</sup>. Dieses ist ein unter der *GPL-Lizenz*<sup>64</sup> veröffentlichtes Programm zum Download von Informationen aus dem Internet. Es unterstützt die *Protokolle HTTP, HTTPS* und *FTP* und ist in der Programmiersprache *C*<sup>65</sup> geschrieben.

Bei **wget** handelt es sich um ein nicht-interaktives Programm. Daher kann es leicht in *Skripten* oder *Cronjobs* aufgerufen werden. *Cronjobs* sind automatisiert zeitgesteuert ausgeführte wiederkehrende Aufgaben.<sup>66</sup> Neben dem Download einzelner Dateien unterstützt **wget** auch das Speichern ganzer *Webseiten* mitsamt den Bildern, beispielsweise zum Zweck des Offline-Lesens.<sup>67</sup> Es bietet ferner die Möglichkeit der Wiederaufnahme von unterbrochenen Downloads. **wget** eignet sich erfahrungsgemäß ausgesprochen gut zum automatisierten Download von *Webseiten*.

Bei **wget** werden die Daten, die übertragen werden sollen, in Form von *URLs* angegeben.

Die allgemeine *Syntax* sieht wie folgt aus:<sup>68</sup>

```
wget [Option] ... [URL] ...
```

---

<sup>63</sup> vgl. <https://www.gnu.org/software/wget/> (19.01.2016)

<sup>64</sup> vgl. <http://www.gnu.org/licenses/gpl-3.0.de.html> (19.01.2016)

<sup>65</sup> vgl. u. a. <http://csapp.cs.cmu.edu/3e/docs/chistory.html> (19.01.2016)

<sup>66</sup> vgl. u. a. <http://pubs.opengroup.org/onlinepubs/9699919799/utilities/crontab.html> (19.01.2016)

<sup>67</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Overview> (19.01.2016)

<sup>68</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Invoking> (19.01.2016)

#### 4.1.1 Optionen von wget

wget bietet der Benutzerin beziehungsweise dem Benutzer verschiedene *Optionen*. Wichtige davon werden im Folgenden erläutert.

Option	Beschreibung
-i <i>datei</i> --input-file= <i>datei</i>	Diese <i>Option</i> veranlasst, dass die <i>URLs</i> der Dateien, die übertragen werden sollen, aus der angegebenen Textdatei namens <i>datei</i> gelesen werden. <sup>69</sup>
-b --background	Diese <i>Option</i> bewirkt, dass der Download als Hintergrundprozess durchgeführt wird.  Wenn nicht anders angegeben, werden alle Statusmeldungen in die Datei <b>wget-log</b> geschrieben, die im <i>Arbeitsverzeichnis</i> erstellt wird. <sup>70</sup>
-B <i>url</i> --base= <i>url</i>	Diese <i>Option</i> hat zur Folge, dass allen relativen <i>Links</i> innerhalb der Datei, die die Benutzerin beziehungsweise der Benutzer durch -i angegeben hat, die <i>Basissadresse url</i> vorangestellt wird. <sup>71</sup>
-c --continue	Diese <i>Option</i> veranlasst, dass ein unterbrochener Download wiederaufgenommen wird.  Wenn sich die Datei allerdings in der Zwischenzeit auf dem <i>Server</i> geändert hat, ist die heruntergeladene Datei eventuell fehlerhaft. <sup>72</sup>
-F --force-html	Diese <i>Option</i> bewirkt, dass die über die <i>Option</i> -i angegebene Datei als Datei im <i>HTML</i> -Format interpretiert wird, und dass alle Dateien, auf welche <i>Links</i> der <i>HTML</i> -Datei verweisen, geladen werden. <sup>73</sup>
--limit-rate= <i>n</i>	Durch diese <i>Option</i> wird die Download-Menge pro Sekunde auf die durch <i>n</i> angegebene Bytezahl limitiert.  Die Suffixe <i>k</i> und <i>m</i> bezeichnen hierbei Kilo- beziehungsweise Megabytes. <sup>74</sup>

<sup>69</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options> (19.01.2016)

<sup>70</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Basic-Startup-Options> (19.01.2016)

<sup>71</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options> (19.01.2016)

<sup>72</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Download-Options> (19.01.2016)

<sup>73</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options> (19.01.2016)

<sup>74</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Download-Options> (19.01.2016)

-N	Diese <i>Option</i> veranlasst, dass eine Datei nur dann heruntergeladen wird, wenn diese neuer ist als eine gleichnamige Datei, die bereits auf dem Computer gespeichert ist. <sup>75</sup>
--output-document= <i>datei</i>	Durch diese <i>Option</i> wird der heruntergeladenen Datei der Name <i>datei</i> gegeben.
-P <i>verzeichnis</i>	Diese <i>Option</i> sorgt dafür, dass alle heruntergeladenen Dateien im <i>Verzeichnis verzeichnis</i> gespeichert werden. <sup>76</sup>
-q --quiet	Ist diese <i>Option</i> aktiviert, werden keine Statusmeldungen ausgegeben. <sup>77</sup>
<b>Optionen für rekursives Herunterladen</b>	
-E --html-extension	Diese <i>Option</i> sorgt dafür, dass die Dateienendung <i>.html</i> bei allen heruntergeladenen Dateien des <i>MIME-Types application/xhtml+xml</i> oder <i>text/html</i> angefügt wird, sofern deren Name nicht mit <i>.html</i> beziehungsweise <i>.htm</i> endet. <sup>78</sup>
-H --span-hosts	Durch diese <i>Option</i> werden auch <i>Links</i> , die auf andere <i>Webseiten</i> verweisen, verfolgt. <sup>79</sup>
-k --convert-links	Durch diese <i>Option</i> werden bei den heruntergeladenen <i>HTML</i> -Dateien die externen <i>Links</i> so geändert, dass sie auf die lokalen Dateien verweisen. <sup>80</sup>
-r --recursive	Diese <i>Option</i> veranlasst, dass <i>rekursiv heruntergeladen</i> wird. <sup>81</sup>
-p --page-requisites	Durch diese <i>Option</i> werden alle Dateien heruntergeladen, welche zum Betrachten der Ausgangsdatei erforderlich sind. <sup>82</sup>

Tabelle 3: Optionen von **wget**

<sup>75</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Download-Options> (19.01.2016)

<sup>76</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Directory-Options> (19.01.2016)

<sup>77</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options> (19.01.2016)

<sup>78</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#HTTP-Options> (19.01.2016)

<sup>79</sup> vgl. [https://www.gnu.org/software/wget/manual/wget.html#Recursive-Accept\\_002fReject-Options](https://www.gnu.org/software/wget/manual/wget.html#Recursive-Accept_002fReject-Options) (19.02.2016)

<sup>80</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Retrieval-Options> (19.01.2016)

<sup>81</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Retrieval-Options> (19.01.2016)

<sup>82</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Retrieval-Options> (19.02.2016)

### 4.1.2 Beispieleingaben zu **wget**

Im Folgenden werden einige Beispieleingaben zur Verwendung von **wget** erläutert.

**Beispiel 1:** `wget http://www.brg19.at/public.php/2/86`

Diese Eingabe lädt die unter der *Adresse* `http://www.brg19.at/public.php/2/86` befindliche *Webseite* herunter. Am 22. Jänner 2016 war dies das Schulprofil des BRG19.

**Beispiel 2:** `wget -c http://www.brg19.at/public.php/2/86`

Wurde der Download aus irgendeinem Grund abgebrochen, kann er über das Aktivieren der *Option* `-c` wiederaufgenommen werden.

**Beispiel 3:** `wget -b http://www.brg19.at/public.php/2/86`

Diese *Option* lässt den Download im Hintergrund ablaufen. Das bietet sich vor allem bei großen Dateien an, die mehr Downloadzeit in Anspruch nehmen.

**Beispiel 4:** `wget -k http://www.brg19.at/public.php/2/86`

Alle externen *Links* werden so verändert, dass sie auf lokale Dateien verweisen. Dies bietet sich an, wenn man eine *Webseite* offline lesen möchte.

### 4.1.3 Anwendung von **wget** auf die verschiedenen Arten von Webseiten

**wget** kann zur Gänze auf die Art der statischen *Webseiten* angewendet werden. Schwierigkeiten treten jedoch unter anderem bei *Webseiten*, welche mit einem *Login-Bereich* geschützt sind, auf, da **wget** diese nicht downloaden kann.

## 4.2 **curl**

`curl`<sup>83</sup> ist ein weiteres Programm zum Download von *Webseiten*. Es ist ein *Open Source Kommandozeilenprogramm*, mit dem man Dateien in der *URL-Syntax* von und zu einem *Server* übertragen kann. Es unterstützt eine Vielzahl an *Protokollen*, unter anderem auch das *HTTP*.<sup>84</sup> Auch `curl` kann in *Skripten* aufgerufen werden.

Die allgemeine *Syntax* sieht wie folgt aus:<sup>85</sup>

```
curl [Option] ... [URL] ...
```

---

<sup>83</sup> vgl. <http://curl.haxx.se/> (19.01.2016)

<sup>84</sup> vgl. <http://curl.haxx.se/>, curl is an open source command line tool and library (19.01.2016)

<sup>85</sup> vgl. <http://curl.haxx.se/docs/manpage.html>, SYNOPSIS (19.01.2016)

### 4.2.1 Optionen von `curl`

`curl` bietet der Benutzerin beziehungsweise dem Benutzer verschiedene *Optionen*. Die wichtigsten werden im Folgenden erläutert.<sup>86</sup>

Option	Beschreibung
<code>--limit-rate n</code>	Durch diese <i>Option</i> wird die Download-Menge pro Sekunde auf die durch <i>n</i> angegebene Bytezahl limitiert.  Die Suffixe <i>k</i> und <i>m</i> bezeichnen hierbei Kilo- beziehungsweise Megabytes.
<code>-o datei</code>	Diese <i>Option</i> veranlasst, dass die heruntergeladenen Daten in einer Datei namens <i>datei</i> gespeichert werden, anstatt sie über die <i>Standardausgabe</i> auszugeben.
<code>-r n1-n2</code>	Durch diese <i>Option</i> wird der durch <i>n1</i> und <i>n2</i> angegebene Bytebereich der Datei übertragen.

Tabelle 4: Optionen von `curl`

### 4.2.2 Beispieleingaben zu `curl`

Im Folgenden werden einige Beispieleingaben zur Verwendung von `curl` erläutert.

**Beispiel 1:** `curl http://www.brg19.at/public.php/2/86`

Diese Eingabe lädt die unter der *Adresse* `http://www.brg19.at/public.php/2/86` befindliche *Webseite* herunter. Am 22. Jänner 2016 war dies das Schulprofil des BRG19.

**Beispiel 2:** `curl -r 2012-2016 http://www.brg19.at/public.php/27/2`

Diese Eingabe lädt den Inhalt der Datei, der sich im Bytebereich 2012 bis 2016 befindet, herunter. Am 22. Jänner 2016 war das das Wort „Stand“.

### 4.2.3 Anwendung von `curl` auf die verschiedenen Arten von Webseiten

Die Anwendung auf die verschiedenen Arten von *Webseiten* betreffend gilt für `curl` das selbe wie für `wget`.

`curl` ist ebenso wie `wget` ein eigenständiges *Tool*, das ohne einen *Browser* auskommt.

<sup>86</sup> vgl. <http://curl.haxx.se/docs/manpage.html>, OPTIONS (19.01.2016)



## 4.3 Selenium

*Selenium* hingegen kommt nicht ohne einen *Browser* aus. Es ist nämlich eine Ansammlung von *Tools*, die sich zum Ziel gesetzt hat, den Test von webbasierten Anwendungen zu automatisieren, es ist folglich eine Automatisierung von *Browsern*.<sup>87</sup>

### 4.3.1 Selenium IDE

In dieser Arbeit wird lediglich ein *Tool* aus der erwähnten Ansammlung genauer besprochen, *Selenium IDE*. Dieses ist ein *Addon* für *Mozilla Firefox*. Es ermöglicht, einzelne Abläufe von Aktionen im *Browser* aufzunehmen und diese später wiederholt abzuspielen. Außerdem bietet es die Möglichkeit, diese Aufnahme in der *Syntax* diverser Programmiersprachen zu exportieren und somit die Automatisierung in *Skripten* zu ermöglichen.<sup>88</sup>

Allerdings unterstützt *Selenium IDE* einige *Befehle* aus anderen *Selenium-Tools* nicht.<sup>89</sup> Hierbei kann man sich Abhilfe schaffen, indem man nach dem Export die *Skripte* dahingehend editiert, dass die fehlenden Funktionen entsprechend von Hand ergänzt werden.

Die Benutzeroberfläche von *Selenium IDE* sieht wie folgt aus:

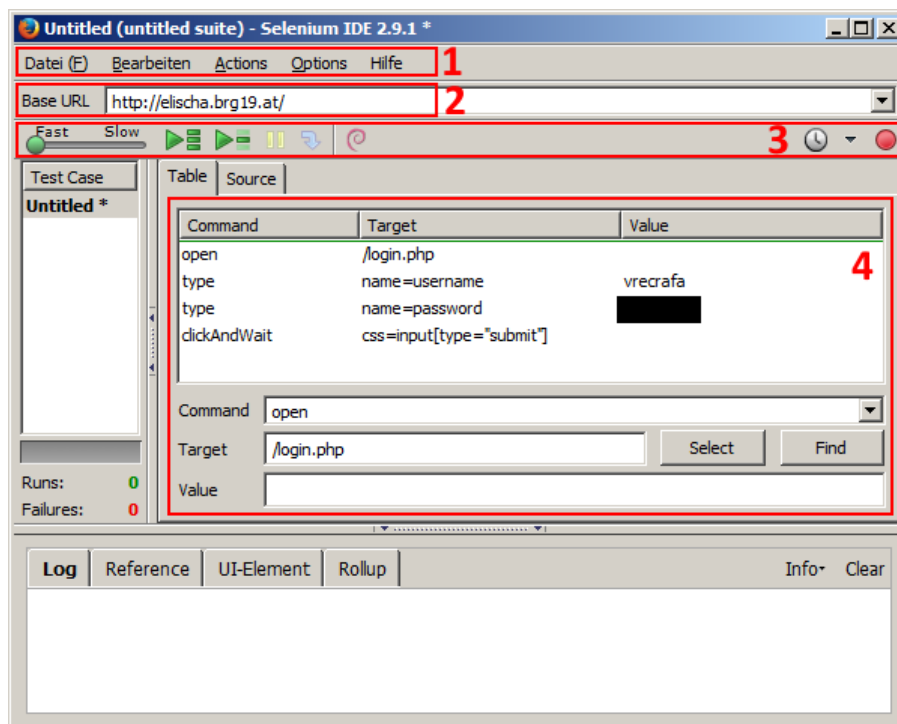


Abbildung 14: Benutzeroberfläche von Selenium IDE

<sup>87</sup> vgl. <http://docs.seleniumhq.org/> (19.01.2016)

<sup>88</sup> vgl. [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp) (19.01.2016)

<sup>89</sup> vgl. <http://www.seleniumhq.org/docs/> (19.01.2016)

- 1 **Menu Bar:** Diese beinhaltet diverse *Programmoptionen*. Im Menüpunkt „Datei“ wird unter anderem die bereits erwähnte *Option* angeboten, Aufnahmen in der *Syntax* verschiedener Programmiersprachen zu exportieren.<sup>90</sup>
- 2 **Base URL:** Diese ist die *URL*, von welcher ausgehend mit der Ausführung der *Befehlskette* begonnen werden soll.
- 3 **Toolbar:** Die *Toolbar* bietet verschiedene *Optionen* an. Man kann eine vollständige Aufnahme oder auch einzelne Teile davon abspielen. Ebenfalls wichtig ist der *Record-Button* auf der rechten Seite, der besagte Aufnahmen ermöglicht.<sup>91</sup>
- 4 **Test Case Pane:** In diesem Fenster werden die einzelnen *Befehle* aufgelistet. Es können weitere von Hand ergänzt beziehungsweise bestehende bearbeitet werden.<sup>92</sup>

In diesem Fall ist in Abbildung 14 bereits das *Einloggen* des Users „vrecrafa“ in das Portal elischa.brg19.at aufgezeichnet worden.

#### 4.3.2 Befehle von Selenium

Folgende Tabelle beinhaltet drei wichtige *Selenium-Befehle*:

Befehl	Beschreibung
open ( <i>URL</i> )	Durch Eingabe dieses <i>Befehls</i> wird die <i>URL URL</i> geöffnet.
type ( <i>locator, value</i> )	Dieser <i>Befehl</i> ermöglicht die Eingabe eines Wertes ( <i>value</i> ) in ein Feld ( <i>locator</i> ).
clickAndWait	Durch Eingabe dieses <i>Befehls</i> wird <i>Selenium</i> mitgeteilt, dass nach dem Klicken auf ein Element (beispielsweise ein Button) eine Interaktion mit dem <i>Server</i> stattfindet.

Tabelle 5: Befehle von Selenium (angewendet in Selenium IDE)<sup>93</sup>

#### 4.3.3 Anwendung von Selenium auf die verschiedenen Arten von Webseiten

*Selenium* ist das mächtigste der hier besprochenen *Tools* zum Abruf. Es ist de facto uneingeschränkt auf die verschiedenen Arten von *Webseiten* anwendbar. Da es *Browser* automatisiert, stellen für *Selenium* auch *Login-Bereiche* kein Problem dar, mit *JavaScript* generierte *Webseiten* ebenso wenig.

<sup>90</sup> vgl. [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp), IDE Features, Menu Bar (19.01.2016)

<sup>91</sup> vgl. [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp), IDE Features, Toolbar (19.01.2016)

<sup>92</sup> vgl. [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp), IDE Features, Test Case Pane (19.01.2016)

<sup>93</sup> vgl. <http://release.seleniumhq.org/selenium-core/1.0.1/reference.html> (19.01.2016)

## 5 Tools zur Verarbeitung

Nach dem Download der jeweiligen *Webseite* liegt der Inhalt, abhängig vom verwendeten *Tool*, im *HTML*-Format vor. Ziel der Weiterverarbeitung ist es, aus diesem die gewünschten Daten zu extrahieren. Zur Weiterverarbeitung sind ebenfalls *Tools* notwendig, welche im Folgenden besprochen werden. Die folgenden *Tools* sind allesamt in *Skripten* aufrufbar und eignen sich somit sehr gut für automatisierte Anwendungen.

### 5.1 GNU Coreutils und ähnliche

In diesem Zusammenhang ist das für diese Arbeit wichtigste Paket *GNU Coreutils*<sup>94</sup>. Dieses ist eine Ansammlung von diversen Programmen zur nicht-interaktiven Textverarbeitung. In dieser Arbeit werden einzelne Programme aus diesem Paket inklusive ihrer *Eingabeparameter* besprochen. Einige Programme, die in ihrer Anwendung ähnlich sind, jedoch nicht zu diesem Paket gehören, werden im Folgenden ebenfalls besprochen.

#### 5.1.1 `cat`<sup>95</sup>

`cat` gibt den Inhalt einer Textdatei aus. Es bietet außerdem die Möglichkeit, mehrere einzelne Dateien zu einer größeren Datei zusammenzufügen.

Die allgemeine *Syntax* lautet wie folgt:

```
cat [Option] ... [Datei] ...
```

`cat` verfügt über mehrere *Optionen*, einige wichtige davon sind:

Option	Beschreibung
-s	Durch diese <i>Option</i> werden mehrere Leerzeilen auf eine einzige Leerzeile reduziert.
-T	Diese <i>Option</i> veranlasst, dass <i>Tabulatoren</i> in der Form <code>^I</code> angezeigt werden.
-v	Diese <i>Option</i> veranlasst, dass Zeichen, die nicht druckbar sind, in einer speziellen Schreibweise angezeigt werden.

Tabelle 6: Optionen von `cat`

<sup>94</sup> vgl. <http://www.gnu.org/software/coreutils/coreutils.html> (22.01.2016)

<sup>95</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#cat-invocation> (19.01.2016)

**Beispiel:** `cat testdatei.txt`

Diese Eingabe gibt den Inhalt der Textdatei `testdatei.txt` auf dem Bildschirm aus.

Die Datei `testdatei.txt` wurde zur Veranschaulichung der Funktionalität der einzelnen *Tools* erstellt. Ihr Inhalt und damit die Ausgabe des obig angeführten *Befehls* sieht wie folgt aus:

1	Vorwissenschaftliche Arbeit 2017
2	Vorwissenschaftliche Arbeit 2015
3	Vorwissenschaftliche Arbeit 2018
4	Vorwissenschaftliche Arbeit 2016
5	Wissenschaftliche Diplomarbeit 2020
6	Matura 2015
7	Volksschulabschluss 2000
8	Referat am 22. Oktober 2015
9	Vorstellung der neuen 1. Klasse um 13:00 Uhr

Quelltext 2: Ausgabe von `cat`

Die obig abgebildete Datei dient im weiteren Verlauf als Muster, auf welches die einzelnen *Tools* angewendet werden. Sie wird nicht für jedes Beispiel erneut abgebildet.

### 5.1.2 `cut`<sup>96</sup>

`cut` extrahiert zeilenweise aus einem Text jene Spalten, die durch die jeweiligen *Optionen* angegeben werden.

Die allgemeine *Syntax* des Programms lautet:

```
cut [Option] ... [Datei] ...
```

`cut` weist verschiedene *Optionen* auf, folgende ist wichtig:

Option	Beschreibung
<code>-c Liste</code> <code>--characters=Liste</code>	Bei dieser <i>Option</i> werden Zeichen abgezählt. Der <i>Parameter Liste</i> gibt an, welche und wie viele Zeichen ausgegeben werden.

Tabelle 7: Eine wichtige Option von `cut`

**Beispiel:** `cut -c 1-19 testdatei.txt`

Diese Eingabe gibt die Zeichen Nummer 1 bis 19 jeder Zeile der angeführten Datei aus.

<sup>96</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#cut-invocation> (19.01.2016)

Die Ausgabe sieht so aus:

```

1
1234567890123456789

```

1	Vorwissenschaftlich
2	Vorwissenschaftlich
3	Vorwissenschaftlich
4	Vorwissenschaftlich
5	Wissenschaftliche D
6	Matura 2015
7	Volksschulabschluss
8	Referat am 22. Okto
9	Vorstellung der neu

Quelltext 3: Ausgabe von `cut`

### 5.1.3 `grep`<sup>97</sup>

`grep` durchsucht Dateien nach bestimmten Suchmustern und zeigt die Zeilen an, in welchen diese gefunden wurden. `grep` ist zwar ein *Tool* des GNU-Projekts, gehört aber nicht zum Paket der *GNU Coreutils*.

Die allgemeine *Syntax* lautet wie folgt:

```
grep [Option] ... Suchmuster [Datei] ...
```

`grep` stellt eine Vielzahl an *Optionen* bereit, folgende werden besprochen:

Option	Beschreibung
<code>-n</code>	Diese <i>Option</i> veranlasst <code>grep</code> dazu, nicht nur die Zeile mit dem gefundenen Textstück anzuzeigen, sondern auch die <i>n</i> vorangehenden und nachfolgenden Zeilen.
<code>-i</code>	Mit dieser <i>Option</i> wird bei der Suche nicht zwischen Groß- und Kleinschreibung unterschieden.
<code>-l</code>	Diese <i>Option</i> gibt die Dateinamen der Dateien aus, in denen das eingegebene Suchmuster gefunden wurde.
<code>-v</code>	Diese <i>Option</i> hat eine inverse Wirkung. Das heißt, dass sie alle Zeilen ausgibt, die dem Suchmuster nicht entsprechen.

Tabelle 8: Optionen von `grep`

<sup>97</sup> vgl. <http://www.gnu.org/software/grep/manual/grep.html> (19.01.2016)

**Beispiel: `grep wissenschaft testdatei.txt`**

Diese Eingabe gibt diejenigen Zeilen aus, in welchen „wissenschaft“ gefunden wird. Das Wort wird farblich hervorgehoben und **fett** gedruckt. Die Groß- und Kleinschreibung wird hierbei berücksichtigt.

Die Ausgabe sieht also wie folgt aus:

1	<b>Vorwissenschaftliche</b> Arbeit 2017
2	<b>Vorwissenschaftliche</b> Arbeit 2015
3	<b>Vorwissenschaftliche</b> Arbeit 2018
4	<b>Vorwissenschaftliche</b> Arbeit 2016

Quelltext 4: Ausgabe von `grep`

#### 5.1.4 `head`<sup>98</sup> und `tail`<sup>99</sup>

`head` beziehungsweise `tail` gibt die ersten beziehungsweise die letzten zehn Zeilen einer Datei auf dem Bildschirm aus, sofern kein *Parameter* übermittelt wird.

Über die *Option* `-n x` lässt sich die Zahl von zehn auf `x` Zeilen verändern.

Für `head` gilt außerdem: Falls dieser Zahl das Zeichen „-“ vorangestellt ist, werden nicht die ersten `x` Zeilen, sondern die letzten `x` Zeilen ausgegeben.

Für `tail` gilt hingegen: Falls dieser Zahl das Zeichen „+“ vorangestellt ist, werden alle Zeilen ab inklusive der angegebenen Zeile ausgegeben.

Die allgemeine *Syntax* lautet:

```
head/tail [Option] ... [Datei] ...
```

**Beispiel: `head -n -7 testdatei.txt`**

Diese Eingabe gibt alle Zeilen von `testdatei.txt` außer die letzten sieben aus:

1	<b>Vorwissenschaftliche</b> Arbeit 2017
2	<b>Vorwissenschaftliche</b> Arbeit 2015

Quelltext 5: Ausgabe von `head`

---

<sup>98</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#head-invocation> (19.01.2016)

<sup>99</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#tail-invocation> (19.01.2016)

Die Ausgabe aller Zeilen außer diesen beiden, sprich die inverse Ausgabe zu der eben genannten, kann durch den *Befehl* `tail -n +3 testdatei.txt` erzielt werden. `tail` gibt dabei alle Zeilen ab inklusive der dritten aus.

### 5.1.5 sed

`sed` steht für *Stream Editor*.<sup>100</sup> Mit ihm können *Text-Datenströme* editiert werden. Der Begriff „*Datenstrom*“ bezeichnet in der Informatik einen kontinuierlichen Fluss von Datensätzen.<sup>101</sup> `sed` verarbeitet Datensätze fortlaufend, sobald jeweils ein neuer Datensatz eingetroffen ist.<sup>102</sup> `sed` ist zwar ein *Tool* des GNU-Projekts, gehört aber nicht zum Paket der *GNU Coreutils*.

Die allgemeine *Syntax* von `sed` sieht so aus:<sup>103</sup>

```
sed [Option] ... [Script] [Datei] ...
```

Da `sed` in der Anwendung wesentlich komplizierter als die bereits besprochenen *Tools* ist, wurde die Erklärung auf jene Grundlagen beschränkt, die für die automatisierte Datenextraktion im Kontext dieser Arbeit erforderlich sind.

Die *Eingabeparameter* werden bei `sed` durch sogenannte *Delimiter* (Trennzeichen) getrennt, welche für die Anwendung passend zu wählen sind.

**Beispiel:** `sed 's/Vorwissenschaftliche/Wissenschaftliche/g' testdatei.txt`

Durch diese Eingabe werden alle Vorkommnisse von „Vorwissenschaftliche“ in „Wissenschaftliche“ geändert. Zu berücksichtigen ist, dass beispielsweise „Vorwissenschaftliches“ in „Wissenschaftliches“ geändert werden würde. Dabei ist irrelevant, wie viele Zeichen auf „Vorwissenschaftliche“ folgen. Sie sind von der Änderung nicht betroffen. Die Groß- und Kleinschreibung wird hierbei berücksichtigt, also wäre „vorwissenschaftliche“ nicht betroffen.

`s` steht hierbei für „substitute“ (ersetze) und ersetzt Zeichen(*ketten*) durch andere. `g` steht hierbei für „global“, was bedeutet, dass die Eingabe nicht nur auf das erste Vorkommnis

---

<sup>100</sup> vgl. <https://www.gnu.org/software/sed/> (19.01.2016)

<sup>101</sup> vgl. [http://www.its.bldrdoc.gov/fs-1037/dir-010/\\_1451.htm](http://www.its.bldrdoc.gov/fs-1037/dir-010/_1451.htm) (19.01.2016)

<sup>102</sup> vgl. <https://www.gnu.org/software/sed/manual/sed.html#Introduction> (19.01.2016)

<sup>103</sup> vgl. <https://www.gnu.org/software/sed/manual/sed.html#Invoking-sed> (19.01.2016)

innerhalb jedes Datensatzes, sondern auf alle Vorkommnisse innerhalb jedes Datensatzes angewendet werden soll.<sup>104</sup>

Die Ausgabe sieht daher wie folgt aus:

10	Wissenschaftliche Arbeit 2017
11	Wissenschaftliche Arbeit 2015
12	Wissenschaftliche Arbeit 2018
13	Wissenschaftliche Arbeit 2016
14	Wissenschaftliche Diplomarbeit 2020
15	Matura 2015
16	Volksschulabschluss 2000
17	Referat am 22. Oktober 2015
18	Vorstellung der neuen 1. Klasse um 13:00 Uhr

Quelltext 6: Ausgabe von `sed`

Bei der obig angesprochenen Eingabe sind die „/“-Zeichen die *Delimiter*. Ein Beispiel für eine Eingabe, in welcher die „/“-Zeichen nicht als *Delimiter* fungieren können, ist, wenn man exemplarisch in einem *HTML*-Dokument alle schließenden *HTML-Tags* in eine andere *Zeichenkette* ändern will. Ein schließender *HTML-Tag*, hier am Beispiel einer Überschrift erster Ebene, sieht wie folgt aus:

```
</h1>
```

Die Verwendung der „/“-Zeichen als *Delimiter* ist in diesem Fall deshalb nicht möglich, da die Eingabe *syntaktisch* falsch wäre, da sie vier als *Delimiter* zu erkennende Zeichen beinhaltet. Ein konkretes Beispiel:

```
sed 's/</h1>/VwA/g'
```

In einem solchen Fall muss ein anderes Zeichen als *Delimiter* verwendet werden. Eine korrekte Eingabe wäre daher:

```
sed 's#</h1>#VwA#g'
```

### 5.1.6 `sort`<sup>105</sup>

`sort` sortiert angegebene Dateien zeilenweise und gibt das Ergebnis davon auf dem Bildschirm aus.

<sup>104</sup> vgl. <https://www.gnu.org/software/sed/manual/sed.html#index-Global-substitution-108> (19.01.2016)

<sup>105</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#sort-invocation> (19.01.2016)



Die allgemeine *Syntax* von `sort` lautet:

```
sort [Option] ... [Datei] ...
```

Einige wichtige *Optionen* von `sort` sind:

Option	Beschreibung
-c	Durch diese <i>Option</i> wird geprüft, ob die Datei bereits sortiert ist oder nicht.
-f	Durch diese <i>Option</i> werden Groß- und Kleinbuchstaben gleichwertig behandelt.
-k <i>n1</i> , <i>n2</i>	Durch diese <i>Option</i> werden nur die Zeichen zwischen der <i>n1</i> -ten und der <i>n2</i> -ten Spalte berücksichtigt. Die Spaltennummerierung beginnt hierbei mit „1“. Spalten werden in der Regel durch Leerzeichen oder <i>Tabulatoren</i> getrennt.
-r	Durch diese <i>Option</i> wird in umgekehrter Reihenfolge sortiert.
-tz	Mit dieser <i>Option</i> wird das Trennzeichen zwischen zwei Spalten angegeben.

Tabelle 9: Optionen von `sort`

**Beispiel:** `sort testdatei.txt`

Diese Eingabe sortiert die Datei `testdatei.txt` zeilenweise, beginnend mit dem ersten Zeichen der jeweiligen Zeile.

Die Ausgabe sieht wie folgt aus:

1	Matura 2015
2	Referat am 22. Oktober 2015
3	Volksschulabschluss 2000
4	Vorstellung der neuen 1. Klasse um 13:00 Uhr
5	Vorwissenschaftliche Arbeit 2015
6	Vorwissenschaftliche Arbeit 2016
7	Vorwissenschaftliche Arbeit 2017
8	Vorwissenschaftliche Arbeit 2018
9	Wissenschaftliche Diplomarbeit 2020

Quelltext 7: Ausgabe von `sort`

### 5.1.7 `tr`<sup>106</sup>

`tr` ersetzt oder löscht Zeichen einer Eingabe und gibt das Ergebnis über die *Standardausgabe* aus.

Die allgemeine *Syntax* von `tr` lautet:

```
tr [Option] ... Zeichenfolge1 [Zeichenfolge2]
```

Eine wichtige *Option* von `tr` ist `-d`. Diese löscht die angegebenen Zeichen. In diesem Fall muss keine zweite *Zeichenfolge* mitgeliefert werden.

**Beispiel:** `tr 'V' 'U' < testdatei.txt`

Durch diese Eingabe wird jedes „V“ in der Datei `testdatei.txt` durch „U“ ersetzt. Die Groß- und Kleinschreibung wird hierbei berücksichtigt. Bei `tr` kann die Eingabedatei nicht als *Parameter* mitgeliefert werden. Das Zeichen „<“ bedeutet *befehlsunabhängig*, dass die Eingabe aus der jeweiligen Datei bezogen werden soll.

Die Ausgabe sieht so aus:

1	Uorwissenschaftliche Arbeit 2017
2	Uorwissenschaftliche Arbeit 2015
3	Uorwissenschaftliche Arbeit 2018
4	Uorwissenschaftliche Arbeit 2016
5	Wissenschaftliche Diplomarbeit 2020
6	Matura 2015
7	Uolkssschulabschluss 2000
8	Referat am 22. Oktober 2015
9	Uorstellung der neuen 1. Klasse um 13:00 Uhr

Quelltext 8: Ausgabe von `tr`

### 5.1.8 Abschlussbemerkungen zu den GNU Coreutils

Die *GNU Coreutils* sind ein sehr mächtiges *Werkzeug* im Umgang mit Text, mit welchem gerade durch die Automatisierung in *Skripten* viele Anwendungsfälle abgedeckt werden können, so auch die automatisierte Datenextraktion aus *Webseiten*. Allerdings stoßen die *GNU Coreutils* dann an ihre Grenzen, wenn sich die Struktur einer *Webseite* ändert beziehungsweise Unregelmäßigkeiten in dieser auftreten. Diese Grenzen werden im Folgenden erläutert. Die *Tools* `grep` und `sed` gehören zwar nicht zum Paket der *GNU Coreutils*, die hier getroffenen Bemerkungen gelten für sie trotzdem.

<sup>106</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#tr-invocation> (19.01.2016)

## 5.2 Parser

Das in Kapitel 5.1.2 zu `cut` angeführte Beispiel ist ein Hinweis darauf, dass die *GNU Coreutils* gewisse Einschränkungen mit sich bringen.

Versucht man mit `cut` beispielsweise gewisse Zeichen jeder Zeile zu extrahieren, sind die erhaltenen Daten eventuell schon unbrauchbar, wenn nur ein Leerzeichen zu viel gesetzt wird, obwohl die überflüssigen Leerzeichen in *HTML* keine Änderung in der Darstellung bewirken.

Ein Beispiel hierfür ist das manuelle Eintragen von Daten über Online-Formulare. Angenommen, jede Lehrperson der Schule trägt ihre Wahlmodulnummer manuell in eine Tabelle ein, dann liefert `cut`, sofern aus Versehen nur ein Leerzeichen zwischen den Ziffern gesetzt wird, ein falsches Ergebnis zurück.

Eine komplexere Struktur kann beispielsweise auch die Sprechstundenliste einer Schule haben. Angenommen, die Sprechstunde einer Lehrperson entfällt jeden ersten Donnerstag im Monat aufgrund von Fortbildung, so ist die beispielhafte Struktur

```
Vorname Nachname  
Sprechstunde
```

bereits unregelmäßig, da sie bei besagter Lehrperson wie folgt aussähe:

```
Vorname Nachname  
Sprechstunde  
Der erste Donnerstag jedes Monats entfällt aufgrund von Fortbildung!
```

Besagte Lehrperson benötigt drei Zeilen, während die restlichen Lehrpersonen lediglich zwei Zeilen benötigen. Will man nun beispielsweise die Sprechstunden herausfiltern und extrahiert mit einem Programm jede zweite Zeile, da in dieser üblicherweise die Sprechstunde steht, so erhält man ab der Zeile nach der Anmerkung die falschen Daten, da sich die Struktur verschoben hat und die Folgezeile keine Sprechstunde mehr ist, wie sie es eigentlich sein sollte, sondern durch die Verschiebung einen Vornamen und einen Nachnamen beinhaltet.

Aus *Webseiten*, die eine sehr komplexe Struktur haben, ist es also schwierig, mithilfe der *GNU Coreutils* zuverlässig Daten zu extrahieren. In diesem Fall kommen sogenannte *Parser* zur Anwendung, um Abhilfe zu schaffen.

Ein *Parser* ist ein Programm, das eine beliebige Eingabe zerlegt und in ein für die Weiterverarbeitung brauchbares Format umwandelt.

Allgemein wird ein *Parser* verwendet, um einen Text in eine neue Struktur zu übersetzen. Dabei müssen *lexikalische* und *grammatische Analysen* durchgeführt werden.

Bei einer *lexikalischen Analyse* wird *plain text* in Folgen von logisch zusammengehörigen Einheiten (beispielsweise in Folgen von Zeichen), sogenannte *Tokens*, zerlegt. Bei der *grammatischen Analyse* werden Eingabedaten auf *syntaktische* Regeln überprüft, an die sie sich halten.<sup>107</sup> Folgendes Beispiel veranschaulicht dies:

Ein *Parser* würde, um bei der obig bereits angesprochenen bei den Sprechstunden auftretenden Problematik zu bleiben, die Struktur

```
Vorname Nachname  
Sprechstunde
```

auch in

```
Vorname Nachname  
Sprechstunde  
Der erste Donnerstag jedes Monats entfällt aufgrund von Fortbildung!
```

erkennen und feststellen, dass die dritte Zeile dieser Struktur nicht entspricht. Die Struktur der zu extrahierenden Daten könnte dem *Parser*, selbstverständlich nicht wie hier in deutscher Sprache, sondern in einer Programmiersprache, wie folgt beschrieben werden:

- Zeilen, die einen Vornamen und einen Nachnamen beinhalten, bestehen aus genau zwei Wörtern. Diese sind durch ein Leerzeichen getrennt.
- Sprechstunden sind im Format Wochentag, Stunde (Uhrzeit) beschrieben.

Die dritte Zeile der unregelmäßigen Struktur entspricht dieser Beschreibung nicht und würde daher nicht erkannt werden.

---

<sup>107</sup> vgl. <http://www.iue.tuwien.ac.at/phd/demel/node107.html> (19.01.2016)

Bei der Datenextraktion aus *Webseiten* ist zumeist eine *HTML*-Interpretation notwendig. Es wird folglich ein *HTML-Parser* benötigt. Ein Beispiel hierfür ist der in der Programmiersprache *Python*<sup>108</sup> geschriebenen `html.parser`.<sup>109</sup>

Reicht selbst dieser nicht aus, um brauchbare Daten zu erhalten, besteht die Möglichkeit, selbst einen auf die jeweilige Anwendung zugeschnittenen *Parser* zu schreiben.

Diese Unternehmung wird durch Programme wie *Bison*<sup>110</sup> und *Flex*<sup>111</sup> unterstützt. Diese werden im Rahmen dieser Arbeit jedoch nicht besprochen.

---

<sup>108</sup> vgl. <https://www.python.org/> (19.01.2016)

<sup>109</sup> vgl. <https://docs.python.org/2/library/htmlparser.html> (19.01.2016)

<sup>110</sup> vgl. <https://www.gnu.org/software/bison/> (19.01.2016)

<sup>111</sup> vgl. <http://flex.sourceforge.net/> (19.01.2016)

## 6 Fallbeispiele

Nach der theoretischen Analyse der einzelnen *Tools* folgt in diesem Teil die praktische Anwendung jener. Die hier angeführten Fallbeispiele wurden speziell für die besprochenen *Tools* erarbeitet und werden im Folgenden Schritt für Schritt erklärt.

### 6.1 Extraktion der Sprechstundenliste des BRG19

Das erste Fallbeispiel ist die Extraktion der Sprechstundenliste des BRG19 aus dessen *Homepage*. Ziel ist es, diese als Datei im Format *Comma-separated values (CSV)* zu erhalten. Das *CSV*-Format wird dafür verwendet, um Dateien zwischen diversen Tabellenkalkulationsprogrammen (wie zum Beispiel *Microsoft Excel*<sup>112</sup>) auszutauschen und zu konvertieren. Es ist im *RFC 4180*<sup>113</sup> grundlegend beschrieben.

Die Sprechstundenliste war am 20. Jänner 2016 unter

```
http://www.brg19.at/public.php/27/2
```

abzurufen. Der erste Schritt ist, diese *Webseite* im *HTML*-Format herunterzuladen, um sie bearbeiten zu können. Diese Aufgabe übernimmt in diesem Fall **wget**, sie könnte allerdings auch in äquivalenter Art und Weise mit **curl** umgesetzt werden.

Folgende Eingabe in ein *Terminal* erfolgt:

```
wget --output-document=sprechstundenliste.html http://www.brg19.at/public.php/27/2
```

Die Datei liegt nun im *Arbeitsverzeichnis*.

Durch das Mitliefern des entsprechenden *Parameters* wurde der heruntergeladenen Datei ein aussagekräftiger Dateiname, in diesem Fall der Name **sprechstundenliste.html** zugewiesen.

Aufgrund der Länge der Datei ist diese auf dem beiliegenden Speichermedium vorhanden. Der *Quelltext* kann im *Browser Mozilla Firefox* durch folgende Eingabe in die *Adresszeile* eingesehen werden:

```
view-source:http://www.brg19.at/public.php/27/2
```

Die Sprechstundenliste liegt in Tabellenform im *HTML*-Format vor.

---

<sup>112</sup> vgl. <https://products.office.com/de-at/excel> (19.01.2016)

<sup>113</sup> vgl. <https://tools.ietf.org/html/rfc4180#page-2> (19.01.2016)

Dabei sieht die Tabellenzeile einer Lehrperson im *Browser* dargestellt wie folgt aus:

(Die Spalten beinhalten chronologisch: Name, Email: @brg19.at, Wochentag, Zeit, Ort)

SCHÖBEL Dipl.-Ing. Christian	scoe	Dienstag	7. Stunde (13:50 - 14:40)	Informatik Sammlung
------------------------------	------	----------	------------------------------	---------------------

Abbildung 15: Tabellenzeile einer Lehrperson

Im *HTML*-Format sieht diese Zeile so aus:

1	<tr bgcolor="#e0e0ff">
2	<td>SCHÖBEL Dipl.-Ing. Christian</td>
3	<td>scoe</td>
4	<td>Dienstag</td>
5	<td>7. Stunde  (13:50 - 14:40)</td>
6	<td>Informatik Sammlung</td>
7	</tr>

Quelltext 9: Aufbau der Tabellenzeile einer Lehrperson

Der *HTML-Tag* <tr> beschreibt eine Zeile. <td> beschreibt die Spalten, die sich in dieser Zeile befinden.

Ziel ist es, die extrahierten Daten zur Weiterverarbeitung in ein für die Anwendungszwecke geeignetes Format umzuwandeln, in diesem Fall in das bereits erwähnte *CSV*-Format.

Die extrahierten Daten werden über das *Terminal* mithilfe der *GNU Coreutils* und *sed* weiterverarbeitet. Die endgültige *Befehlskette*, die wegen der Übersichtlichkeit in Zeilen unterteilt wurde und im Folgenden zeilenweise erklärt wird, sieht wie folgt aus:

1	cat sprechstundenliste.html
2	tr '\n' ' '
3	sed 's#</tr>#\n#g'
4	cut -c 32-
5	sed 's#</td>##g'
6	sed 's#<td>#;#g'
7	sed 's#\t##g'
8	tail -n +7
9	head -n -2
10	sed 's# ##g'
11	sed 's#<td colspan="2">#;#g'
12	sed 's# ;#;#g'
13	> sprechstundenliste.csv

Quelltext 10: Befehlskette zur Extraktion der Sprechstundenliste

Gearbeitet wird hierbei mit sogenannten *Pipelines (Pipes)*. Diese ermöglichen der Anwenderin beziehungsweise dem Anwender die Ausgabe eines Programms in die Eingabe eines anderen Programms umzuleiten. Die *Pipelines* werden durch das Zeichen „|“ ausgezeichnet. Sie sind im zuvor abgebildeten *Quelltext* am Ende der jeweiligen Zeile angeführt.

**Zeile 1:** `cat sprechstundenliste.html |`

`cat` würde hierbei die Datei auf dem Bildschirm als Text ausgeben. Diese Ausgabe wird durch eine *Pipeline* (`|`) zum nächsten *Befehl* umgeleitet. Diese Umleitung erfolgt bei den weiteren *Befehlen* ebenfalls, wird im Folgenden aber nicht mehr gesondert erwähnt.

**Zeile 2:** `tr '\n' ' '`

`tr` ersetzt alle `'\n'`, also alle Zeilenumbrüche, durch ein Leerzeichen, da zwischen den zwei Hochkommata ein Leerzeichen () steht.

**Zeile 3:** `sed 's#</tr>#\n#g'`

`sed` ersetzt alle schließenden `</tr>`-Tags (`</tr>`) durch einen Zeilenumbruch (`\n`). Dies hat zur Folge, dass nun jede Lehrperson eine eigene Zeile hat.

**Zeile 4:** `cut -c 32-`

`cut` schneidet in jeder Zeile den Text ab inklusive dem 32. Zeichen aus. Diese sind speziell für diese Anwendung abgezählt und somit beginnt eine neue Zeile genau mit dem Namen der Lehrperson.

**Zeile 5:** `sed 's#</td>##g'`

`sed` ersetzt nun alle schließenden `</td>`-Tags (`</td>`) durch einen Leerstring, da zwischen den zwei *Delimitern* kein Zeichen angeführt ist. Dadurch werden diese *Tags* entfernt.

**Zeile 6:** `sed 's#<td>#;#g'`

Nun werden alle öffnenden `<td>`-Tags (`<td>`) von `sed` durch ein Semikolon (`;`), in diesem Fall das Trennzeichen der *CSV*-Datei, ersetzt.

**Zeile 7:** `sed 's#\t##g'`

Von `sed` werden alle Tabstopps (`\t`) entfernt.

**Zeile 8:** `tail -n +7`

`tail` gibt den gesamten Text ab inklusive der siebten Zeile aus.

**Zeile 9:** `head -n -2`

`head` gibt den gesamten Text exklusive der letzten zwei Zeilen aus.



Durch die Verwendung von **tail** und **head** werden durch die Extraktion überflüssig gewordene Zeilen, wie zum Beispiel jene, die den *HTML*-Kopf beinhalten, entfernt. Ab beziehungsweise bis zu welcher Zeile ausgegeben werden soll, ist abgezählt.

**Zeile 10:** `sed 's#<br />##g'`

`sed` entfernt nun noch die durch die Extraktion überflüssig gewordenen `<br />`-Tags.

**Zeile 11:** `sed 's#<td colspan="2">#;#;#g'`

`sed` ersetzt den nun überflüssigen Tag `<td colspan="2">` durch zwei Semikolons. Dies ist für den Erhalt der Tabellenstruktur erforderlich, da die Sprechstunde der Frau Direktorin einen Ausnahmefall darstellt, da für diese kein fixer Termin festgelegt ist.

**Zeile 12:** `sed 's# ;#;#g'`

`sed` sorgt dafür, dass jedes Leerzeichen, das vor einem Semikolon steht, entfernt wird.

**Zeile 13:** `> sprechstundenliste.csv`

Das Zeichen „>“ leitet die *Standardausgabe* statt auf den Bildschirm in die Datei `sprechstundenliste.csv` um, die im *Arbeitsverzeichnis* angelegt wird. Die Zeile einer Lehrperson sieht nun wie folgt aus:

SCHÖBEL Dipl.-Ing. Christian;scoe;Dienstag;7. Stunde (13:50 - 14:40);Informatik Sammlung
--

Quelltext 11: Tabellenzeile einer Lehrperson (`sprechstundenliste.csv`)

Diese Struktur entspricht exakt der bereits in Abbildung 15 dargestellten, jedoch ist sie nun nicht im *HTML*-, sondern im *CSV*-Format abgebildet. Die angewendete *Befehlskette* kann über ein *Skript* vollständig automatisiert werden.

Eine etwas andere Vorgangsweise ist beim Download mit dem textbasierten *Browser lynx* und der darauffolgenden Weiterverarbeitung bei der Extraktion einer Liste der Lehrpersonen aus der gleichen Quelldatei zu verzeichnen. Der *Quelltext* hierzu ist im Anhang (A.II) angeführt.

## 6.2 Extraktion der eigenen Modulliste von elischa.brg19.at

Das zweite Fallbeispiel befasst sich mit der Extraktion der eigenen Modulliste aus dem Portal „Elischa“ (elischa.brg19.at) des BRG19. Diese erweist sich als etwas komplizierter, da die eigene Modulliste erst nach dem *Einloggen* in das Portal verfügbar ist. Für das automatisierte *Einloggen* und den Download der Modulliste wird das *Tool Selenium* zur Anwendung gebracht.

Folgende *Befehlsabfolge* ist hierfür erforderlich:

Befehl	Ziel	Wert
open	/login.php	
type	name=username	vrecrafa
type	name=password	1234strengGEHEIM#
clickAndWait	css=input[type="submit"]	
open	/pupils/elias/elias-angemeldet.php	
...	...	...

Tabelle 10: Befehle zur Extraktion der Modulliste im HTML-Format

Zunächst wird unter elischa.brg19.at die Datei `login.php` geöffnet. Dieser werden der Account-Name, in diesem Fall „vrecrafa“, und das Passwort übergeben. Danach wird der Button mit der Aufschrift „anmelden!“ geklickt, um *eingeloggt* zu werden.

### ELISCHA Elektronische Lösung für Intelligente SCHul Administration

Abbildung 16: Diese Eingaben werden von Selenium automatisiert

Nach dem erfolgreichen *Einloggen* wird die unter

`http://elischa.brg19.at/pupils/elias/elias-angemeldet.php`

abzurufende Datei geöffnet. Diese beinhaltet die Modulliste des dazugehörigen Accounts, in diesem Fall jene von „vrecrafa“.

Für die Speicherung der Datei ist *Selenium IDE* in diesem Fall nicht geeignet. Um die Datei dennoch zu speichern, wird die bereits erwähnte Möglichkeit des Exportierens in der *Syntax* einer Programmiersprache, in diesem Fall *Python*<sup>114</sup>, verwendet. Das vollständige *Python-Skript* ist im Anhang (A.III) angeführt. Jene Passagen, die nicht von *Selenium IDE* generiert wurden, sind **fett** gedruckt. Im Folgenden werden sie kurz besprochen.

27	<b># Pfadangabe - Wo soll Datei gespeichert werden?</b>
28	<b>f = codecs.open("C:/frei/" + "modulliste" + ".html", "w", "utf-8")</b>
29	<b>f.write(driver.page_source)</b>
30	<b>f.close()</b>

Quelltext 12: Befehle zur Speicherung des Quelltexts der Modulliste

- Zeile 27 ist ein Kommentar, dieser wird bei der Programmausführung nicht berücksichtigt.<sup>115</sup>
- Zeile 28 öffnet beziehungsweise erstellt die Datei, in welcher der *Quelltext* der Modulliste gespeichert werden soll. Die Datei soll in diesem Fall unter `C:/frei/` gespeichert werden (*Pfadangabe* hier mit *Slashes*), den Namen `modulliste` tragen und die Endung `.html` haben. Weiter soll die Datei mit Schreibrechten geöffnet werden, sowie in *UTF-8*<sup>116</sup> codiert sein.<sup>117</sup>
- Zeile 29 sorgt dafür, dass der *Quelltext* der Modulliste in diese Datei geschrieben wird.<sup>118</sup> Dabei kommt der *Selenium WebDriver*<sup>119</sup>, ein *Tool* aus der Ansammlung der *Selenium-Tools*, zur Anwendung.
- Zeile 30 veranlasst das Schließen der geöffneten Datei.<sup>120</sup>

<sup>114</sup> vgl. <https://www.python.org/> (19.01.2016)

<sup>115</sup> vgl. <https://www.python.org/dev/peps/pep-0008/#inline-comments> (19.01.2016)

<sup>116</sup> vgl. <https://tools.ietf.org/html/rfc3629> (19.01.2016)

<sup>117</sup> vgl. <https://docs.python.org/2/library/codecs.html> (19.01.2016)

<sup>118</sup> vgl. <https://docs.python.org/2/tutorial/inputoutput.html> (19.01.2016)

<sup>119</sup> vgl. <http://www.seleniumhq.org/projects/webdriver/> (19.01.2016)

<sup>120</sup> vgl. <https://docs.python.org/2/tutorial/inputoutput.html> (19.01.2016)

Die Datei ist nach erfolgreichem Download im von der Benutzerin beziehungsweise dem Benutzer gewählten *Verzeichnis* abgelegt. Aufgrund der Länge der Datei wird sinnvollerweise lediglich mit Ausschnitten gearbeitet. Die Datei ist auf dem beiliegenden Speichermedium vorhanden.

Auch die Modulliste liegt in Tabellenform im *HTML*-Format vor. Sie sieht im *Browser* dargestellt wie folgt aus:

#	Modul-ID	Fächer	Matura	Titel	LehrerInnen	Zeit	Optionen/Status
21	15225	M		Kompetent in Mathematik	DMA	Mi, 14:40 - 16:20	
20	15144	SQ	SQ	Selbstmanagement	POB	Mi, 14:40 - 16:20	
19	15131	PE	PE	Psychotherapien	POB	Mo, 14:40 - 16:20	
18	15123	INF		Projektarbeit aus Informatik	SCOE	Di, 16:20 - 18:00	
17	15120	INF	INF	Informatik 1: PASCAL für Anfänger	SCOE	Di, 14:40 - 16:20	
16	14229	M		Mathematik Kompetent	DMA,FRA,DOV	Di, 16:20 - 18:00	
15	14222	INF	INF	PASCAL für Fortgeschrittene	SCOE	Di, 14:40 - 16:20	
14	14220	H	H	Lebensformen im Mittelalter	JLK	Mo, 16:20 - 18:00	
13	14208	BIO	BIO	Meeresbiologie	SCL	Mo, 14:40 - 16:20	
12	14129	L	L	Rom als Keimzelle Europas	DOB	Mo, 14:40 - 16:20	
11	14127	INF	INF	Computernetzwerke	SCOE	Di, 16:20 - 18:00	
10	14126	INF	INF	Einführung in Linux	SCOE	Di, 14:40 - 16:20	
9	14125	INF	INF	Projektarbeit aus Informatik	SCOE	Do, 16:20 - 18:00	
8	13251	SP	SP	Spanisch 2		nach Vereinbarung	
7	13247	SQ	SQ	Konflikte verstehen und bearbeiten	LED	Mi, 14:40 - 16:20	
6	13220	INF	INF	Webscripting und Datenbanken	SCOE	Do, 16:20 - 18:00	
5	13219	INF	INF	Betriebssysteme	SCOE	Do, 14:40 - 16:20	
4	13158	SP	SP	Spanisch 1		nach Vereinbarung	
3	13128	INF	INF	Maschinennahe Programmierung	SCOE	Do, 14:40 - 16:20	
2	13127	INF	INF	Microcontroller	SCOE	Do, 16:20 - 18:00	
1	13108	CH	CH	Chemische Übungen	STN	Mo, 14:40 - 16:20	

Abbildung 17: Modulliste

Die Tabellenzeile eines Moduls sieht im *HTML*-Format so aus:

1	<tr class="hi1">
2	<td class="num">10</td>
3	<td class="mid">14126</td>
4	<td class="fks">INF</td>
5	<td class="fks">INF</td>
6	<td>Einführung in Linux</td>
7	<td>SCOE</td>
8	<td class="tright"><span class="nobr">Di, 14:40 - 16:20</span></td>
9	<td><span class="nobr"></span></td>
10	</tr>

Quelltext 13: Aufbau der Tabellenzeile eines Moduls

Das Ziel ist auch hierbei, die extrahierten Daten in ein an die Anwendungszwecke angepasstes Format zur Weiterverarbeitung umzuwandeln. Auch in diesem Fall ist das, wie schon im ersten Fallbeispiel, das bereits erwähnte *CSV*-Format.

Auch in diesem Fallbeispiel werden die extrahierten Daten mithilfe der *GNU Coreutils* weiterverarbeitet.

Erneut die endgültige *Befehlskette* in durch Zeilen separierter Form:

1	cat modulliste.html
2	tr '\n' ' '
3	sed 's#</tr>#\n#g'
4	sed 's#<td[^>]*>##'
5	sed 's#<td[^>]*>#;#g'
6	sed 's#[^>]*>##g'
7	sed 's\t##g'
8	cut -c 3-
9	sed 's# ;#;#g'
10	tail -n +2
11	> modulliste.csv

Quelltext 14: Befehlskette zur Extraktion der Modulliste

Im Folgenden werden nun allerdings nur mehr jene Zeilen besprochen, die im ersten Fallbeispiel noch nicht in dieser Form erwähnt wurden.

Folgende Zeilen werden nicht mehr erklärt:

- Zeile 1, da sie bis auf den Dateinamen mit Zeile 1 aus Fallbeispiel 1 ident ist
- Zeilen 2, 3 und 7, da sie ident mit selbigen aus Fallbeispiel 1 sind
- Zeile 9, da sie ident mit Zeile 12 aus Fallbeispiel 1 ist

Folgende Zeilen werden gesondert erklärt:

**Zeile 4:** `sed 's#<td[^>]*>##'`

`sed` entfernt den ersten `<td>`-Tag, mitsamt dessen *Parametern*.

Da *HTML-Tags* verschieden aussehen, sich aber alle an eine bestimmte Logik halten, wird hierbei mit *Regular Expressions*<sup>121</sup> (*Reguläre Ausdrücke*) gearbeitet.

<sup>121</sup> vgl. u. a. [http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd\\_chap09.html](http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd_chap09.html) (19.01.2016)

Folgendes Beispiel veranschaulicht die Vorgehensweise bei der Arbeit mit *Regular Expressions*:

```
<h1>
```

```
<p>
```

Die beiden angeführten *HTML-Tags* haben folgende Gemeinsamkeiten:

- Sie sind beide öffnende *Tags*, das heißt, sie haben beide kein „/“.
- Sie beginnen mit <.
- Sie enden mit >.

Eine *Regular Expression* beschreibt diese Gemeinsamkeiten. Dies passiert auch im obigen *Befehl*.

Folgende Beschreibung trifft der Ausdruck `<td[^>]*>` bei der Verwendung mit `sed`:<sup>122</sup>

- Der *HTML-Tag* wird mit < geöffnet, danach folgt `td`.
- `[^>]` wählt alle Zeichen aus, die kein > sind. Das ist wichtig, da weitere *Parameter* vorhanden sein könnten.
- \* bedeutet, dass die ausgewählten Zeichen beliebig oft vorkommen dürfen.
- Der *HTML-Tag* wird mit > geschlossen.

In jeder Zeile wird die erste *Zeichenkette*, die dieser Beschreibung entspricht, entfernt, alle weiteren jedoch nicht, da die *Option* für `g` (für „global“) nicht aktiviert wurde.

**Zeile 5:** `sed 's#<td[^>]*>#;#g'`

`sed` ersetzt nun global alle weiteren `<td>`-*Tags* durch Semikolons, in diesem Fall das Trennzeichen der *CSV-Datei*.

**Zeile 6:** `sed 's#[^>]*>##g'`

`sed` entfernt global alle weiteren *HTML-Tags*. Der einzige Unterschied verglichen mit der Beschreibung in Zeile 4 liegt darin, dass nach dem Zeichen „<“ nicht auch noch die *Zeichenfolge* „`td`“ stehen muss, sondern beliebige Zeichen stehen dürfen.

---

<sup>122</sup> vgl. [https://www.gnu.org/software/sed/manual/html\\_node/Regular-Expressions.html](https://www.gnu.org/software/sed/manual/html_node/Regular-Expressions.html) (20.01.2016)

**Zeile 8:** `cut -c 3-`

`cut` schneidet den Text ab inklusive dem dritten Zeichen aus. Diese sind speziell für diesen Fall abgezählt.

**Zeile 10:** `tail -n +2`

`tail` gibt den gesamten Text ab inklusive der zweiten Zeile aus.

**Zeile 11:** `> modulliste.csv`

Das Ergebnis der *Befehlskette* wird in der Datei `modulliste.csv` gespeichert:

```
10;14126;INF;INF;Einführung in Linux;SCOE;Di, 14:40 - 16:20;
```

Quelltext 15: Tabellenzeile eines Moduls (`modulliste.csv`)

Auch hier liegt nun die gleiche Struktur vor, die in Abbildung 17 bereits dargestellt worden ist. Die verwendete *Befehlskette* kann auch hier über ein *Skript* vollständig automatisiert werden.

## 7 Bewertungen und Schlussfolgerungen

In dieser Arbeit wurde der gesamte Prozess einer automatisierten Datenextraktion aus *Webseiten* abgehandelt. Daraus lassen sich einige Schlussfolgerungen ziehen, die im folgenden Kapitel abgehandelt werden.

### 7.1 Kenntnissen über das Internetsurfen und den Aufbau von Webseiten

Grundsätzlich ist es wichtig, über die Vorgänge, die beim *Internetsurfen* ablaufen, Bescheid zu wissen, um den Prozess der automatisierten Datenextraktion nachvollziehen zu können. Ist kein Wissen über die theoretische Funktionsweise des *HTTP* oder den grundlegenden Aufbau einer *Webseite* vorhanden, ist es de facto unmöglich, einen Automatisierungsprozess zu planen, da der Aufbau einer *Webseite* und der Transport dieser unabdingbar für die Datenextraktion aus dieser sind. Die wohl wichtigste Grundlage ist in diesem Zusammenhang die Tatsache, dass die Daten, welche man aus einer *Webseite* extrahieren will, in der Regel im *HTML*-Format vorliegen. Insofern sind detaillierte Kenntnisse über *HTML* eine essenzielle Voraussetzung für eine automatisierte Datenextraktion aus *Webseiten*. Die einzig noch sinnvoll erscheinende Möglichkeit, *HTML* zu umgehen, ist jene, die *Webseite* nicht mithilfe von `wget` oder `curl` herunterzuladen, sondern einen textbasierten *Browser* wie beispielsweise `lynx` zu verwenden, da dieser die Seite als Text ohne die *HTML*-Struktur abliefern kann.

### 7.2 Die Art einer Webseite

Neben dem Aufbau einer *Webseite* ist auch deren damit in Verbindung stehende Art für die automatisierte Datenextraktion von großer Bedeutung. Diese Arbeit nahm empirisch eine eigene Klassifizierung von *Webseiten* vor. In diesem Zusammenhang ist zu sagen, dass die eine automatisierte Datenextraktion planende Person über die grundlegende Funktionsweise und die Art der jeweiligen *Webseite* Bescheid wissen muss. Mutmaßungen über die Art einer *Webseite* lassen sich in jedem Fall leicht anstellen. Beispielsweise ist anzunehmen, dass *Webseiten*, die Informationen wie das heutige Datum oder Ähnliches beinhalten, generiert werden. Ob diese Elemente *serverseitig* oder *clientseitig* generiert werden, lässt sich leicht feststellen, indem man sich den erhaltenen *Quelltext* ansieht. Sind darin keine *clientseitig* generierten Elemente enthalten, ist festzustellen, dass die *Webseite* statisch ist oder *serverseitig* generiert worden ist. Sofern die *Webseite* keinen *Login-Bereich* hat, ist eine weitere Unterscheidung zwischen den beiden genannten Arten für die Anwenderin beziehungsweise den Anwender sogar irrelevant, da sich die besprochenen *Tools* für



beide Arten gleich verhalten. Die Extraktion von Daten aus *Webseiten*, die durch einen *Login-Bereich* geschützt sind, lässt sich hingegen sinnvollerweise nur mit *Selenium* realisieren.

### 7.3 Sinnvolle Verwendung geeigneter Tools

In diesem Zusammenhang ist generell zu sagen, dass die wohl wichtigste Entscheidung im Prozess der automatisierten Datenextraktion jene über die verwendeten *Tools* ist. Eine unbedachte Wahl kann hierbei auf der einen Seite einen enormen Mehraufwand und auf der anderen Seite eine geringere Zuverlässigkeit, was die Korrektheit der extrahierten Daten betrifft, bedeuten. Insofern ist es notwendig, grundlegende Überlegungen über die einzelnen *Tools* anzustellen und zu ermitteln, ob sie wirklich die Anforderungen für den entsprechenden Anwendungsbereich erfüllen. Weiter ist hierbei zu bemerken, dass die *GNU Coreutils* grundsätzlich mit wenig Aufwand sehr schnell brauchbare Ergebnisse liefern. Bei einer komplexeren Struktur der Quelle stoßen sie allerdings schnell an ihre Grenzen. Abhilfe kann man sich vorerst noch schaffen, indem man mit *Regular Expressions* arbeitet. Diese gehen aber ohnehin schon in Richtung *Parser*, da sie ebenfalls gewisse Vorkommnisse allgemein beschreiben. Oftmals gibt es aber wohl kaum einen anderen Ausweg, als einen *Parser* zu verwenden beziehungsweise sogar eigenhändig zu schreiben.

### 7.4 Abschlussbewertung

Abschließend lässt sich feststellen, dass die automatisierte Datenextraktion aus *Webseiten* in vielen Fällen erstrebenswert scheint, da sie eine große Zeitersparnis bedeutet und vor allem auch die Fehleranfälligkeit reduziert. Es ist jedoch zu bemerken, dass ein großes Maß an Vorkenntnissen benötigt wird, um eine zuverlässige Anwendung zu gewährleisten. Es kommt dabei darauf an, inwiefern der dafür erforderliche Aufwand in einem Verhältnis zum daraus resultierenden Nutzen steht. Eine Automatisierung macht nur dann Sinn, wenn dadurch merklich Zeit eingespart werden kann und die extrahierten Daten regelmäßig in einer aktuellen Form benötigt werden. Ansonsten steht der Aufwand wohl kaum im Verhältnis zum tatsächlichen Nutzen. Die Automatisierung einer einmaligen Extraktion von Daten aus einer bestimmten *Webseite* scheint ohne fundierte Vorkenntnisse kaum sinnvoll. Für eine wiederkehrende Extraktion ist es allerdings sehr wohl empfehlenswert, sich entsprechende Kenntnisse anzueignen und eine Automatisierung anzustreben.

## A Anhang

A.I	Abstract (deutschsprachige Fassung) .....	58
A.II	Fallbeispiel 1, Daten im Textformat von <b>lynx</b> abgeholt .....	59
A.III	Python-Skript zur Speicherung der Modulliste .....	59

### A.I Abstract (deutschsprachige Fassung)

Die automatisierte Datenextraktion aus *Webseiten* hat das Ziel der Zeitersparnis und der Verringerung der Fehlerquote. Diese Arbeit befasst sich mit einem derartigen Automatisierungsprozess.

Begonnen wird mit der Beschreibung der Vorgänge beim *Internetsurfen*. Dabei wird auf das *Hypertext Transfer Protocol (HTTP)* und speziell auf dessen Funktionseinheiten *Client* und *Server* eingegangen. Weiter wird empirisch eine Klassifizierung von *Webseiten* vorgenommen. In weiterer Folge werden diverse *Tools*, welche für die Automatisierung erforderlich sind, besprochen. Diese werden in die Kategorien „Tools zum Abruf“ und „Tools zur Verarbeitung“ eingeteilt.

Nach der theoretischen Beschreibung der einzelnen *Tools* folgt die Analyse eines Automatisierungsprozesses auf Basis zweier Fallbeispiele. Dabei werden die Sprechstundenliste des BRG19 sowie eine individuelle Liste an gewählten Modulen aus den im Vorfeld heruntergeladenen Dateien im Format *Hypertext Markup Language (HTML)* extrahiert und in Dateien im Format *Comma-separated values (CSV)* umgewandelt.

Abschließend werden Bewertungen und Schlussfolgerungen den Automatisierungsprozess betreffend angestellt, die analysieren, welche Kenntnisse für eine derartige Automatisierung erforderlich sind und inwiefern eine Automatisierung sinnvoll beziehungsweise nicht sinnvoll ist.

## A.II Fallbeispiel 1, Daten im Textformat von **lynx** abgeholt

```
lynx --dump http://www.brg19.at/public.php/27/2 > lynxdump.txt
```

1	cat lynxdump.txt
2	cut -c 4-
3	head -n -17
4	tail -n +15
5	sort
6	sed 's/\s[a-z]/\n/'
7	grep -v '^[a-z]'
8	tail -n +71
9	> lehrpersonenliste_lynx.txt

Quelltext 16: Befehlskette zur Extraktion einer Liste der Lehrpersonen (**lynx**)

Die Datei `lynxdump.txt` ist auf dem beiliegenden Speichermedium vorhanden.

## A.III Python-Skript zur Speicherung der Modulliste

1	# -*- coding: utf-8 -*-
2	from selenium import webdriver
3	from selenium.webdriver.common.by import By
4	from selenium.webdriver.common.keys import Keys
5	from selenium.webdriver.support.ui import Select
6	from selenium.common.exceptions import NoSuchElementException
7	from selenium.common.exceptions import NoAlertPresentException
8	import unittest, time, re, codecs
9	
10	class Webdriver(unittest.TestCase):
11	def setUp(self):
12	self.driver = webdriver.Firefox()
13	self.driver.implicitly_wait(30)
14	self.base_url = "http://elischa.brg19.at/"
15	self.verificationErrors = []
16	self.accept_next_alert = True
17	
18	def test_webdriver(self):
19	driver = self.driver
20	driver.get(self.base_url + "/login.php")
21	driver.find_element_by_name("username").clear()
22	driver.find_element_by_name("username").send_keys("vrecrafa")
23	driver.find_element_by_name("password").clear()
24	driver.find_element_by_name("password").send_keys("██████████")
25	driver.find_element_by_css_selector("input[type='submit']").click()
26	driver.get(self.base_url + "/pupils/elias/elias-angemeldet.php")
27	# Pfadangabe - Wo soll Datei gespeichert werden?
28	f = codecs.open("C:/frei/" + "modulliste" + ".html", "w", "utf-8")
29	f.write(driver.page_source)
30	f.close()
31	
32	def is_element_present(self, how, what):

33	try: self.driver.find_element(by=how, value=what)
34	except NoSuchElementException as e: return False
35	return True
36	def is_alert_present(self):
37	try: self.driver.switch_to_alert()
38	except NoAlertPresentException as e: return False
39	return True
40	def close_alert_and_get_its_text(self):
41	try:
42	alert = self.driver.switch_to_alert()
43	alert_text = alert.text
44	if self.accept_next_alert:
45	alert.accept()
46	else:
47	alert.dismiss()
48	return alert_text
49	finally: self.accept_next_alert = True
50	
51	def tearDown(self):
52	finally: self.accept_next_alert = True
53	
54	self.driver.quit()
55	self.assertEqual([], self verificationErrors)
56	
57	if __name__ == "__main__":
58	unittest.main()

Quelltext 17: Python-Skript zur Speicherung der Modulliste

## B Glossar

Im folgenden Glossar sind alle im Text gekennzeichneten Begriffe alphabetisch aufsteigend sortiert. Auf der linken Seite ist jeweils der Begriff, auf der rechten Seite die Erklärung zu finden. Die Quellen wurden als Fußnoten angefügt und sind auch im Literaturverzeichnis (C Literaturverzeichnis) angeführt. Die Voranstellung von „Abk.“ kennzeichnet Abkürzungen als solche. Mehrfachbedeutungen wurden durch Voranstellung von „hier:“ ausgewiesen, worauf die im Kontext dieser Arbeit relevante Bedeutung folgt. Verweise wurden durch die Voranstellung von „siehe“ gekennzeichnet. Referenzen auf das Glossar wurden nach dem bisher gebrauchten Schema auf der rechten, nicht aber auf der linken Seite gesetzt.

In einem Glossar ist es teilweise notwendig, Begriffe wörtlich oder zumindest beinahe wörtlich zu übernehmen. Da aber jeweils nur die im Kontext relevante Bedeutung des Begriffs, teilweise sogar in leicht abgewandelter Form, übernommen wurde, wurde auch hierbei mit der in der Arbeit gebräuchlichen Zitierweise gearbeitet.

Abk. ASCII	siehe <i>American Standard Code for Information Interchange</i>
Abk. ASP.NET	siehe <i>Active Server Pages .NET</i>
Abk. CSS	siehe <i>Cascading Style Sheets</i>
Abk. CSV	siehe <i>Comma-separated values</i>
Abk. FTP	siehe <i>File Transfer Protocol</i>
Abk. GPL	siehe <i>GNU General Public License</i>
Abk. HTML	siehe <i>Hypertext Markup Language</i>
Abk. HTTP	siehe <i>Hypertext Transfer Protocol</i>
Abk. IETF	siehe <i>Internet Engineering Task Force</i>
Abk. IMAP	siehe <i>Internet Message Access Protocol</i>
Abk. JSP	siehe <i>JavaServer Pages</i>
Abk. LDAP	siehe <i>Lightweight Directory Access Protocol</i>
Abk. NTP	siehe <i>Network Time Protocol</i>
Abk. PHP	rekursives Akronym für <i>PHP Hypertext Processor</i> ; siehe <i>PHP Hypertext Processor</i>
Abk. POP	siehe <i>Post Office Protocol</i>
Abk. RFC	siehe <i>Request for Comments</i>
Abk. SMTP	siehe <i>Simple Mail Transfer Protocol</i>
Abk. URL	siehe <i>Uniform Resource Locator</i>
Abk. W3C	siehe <i>World Wide Web Consortium</i>
Active Server Pages .NET	Software zur Entwicklung <i>dynamischer Webseiten</i> <sup>123</sup>
Addon	Hilfsprogramm, mit welchem ein Anwendungsprogramm erweitert wird <sup>124</sup>

<sup>123</sup> vgl. <http://www.asp.net/> (19.01.2016)

<sup>124</sup> vgl. [http://www.duden.de/rechtschreibung/Add\\_on](http://www.duden.de/rechtschreibung/Add_on) (19.01.2016)

<b>Adresse</b>	hier: siehe <i>Internetadresse</i>
<b>Adresszeile</b>	Bedienelement eines <i>Browsers</i> , in welches unter anderem die <i>URL</i> einer <i>Webseite</i> eingegeben werden kann <sup>125</sup>
<b>American Standard Code for Information Interchange</b>	genormter Code für Ziffern, Buchstaben und Sonderzeichen <sup>126</sup>
<b>Anfrage</b>	hier: siehe <i>Request</i>
<b>Anfragenachricht</b>	siehe <i>Request-Message</i>
<b>Antwort</b>	hier: siehe <i>Request</i>
<b>Arbeitsverzeichnis</b>	siehe <i>Verzeichnis</i> ; <i>Verzeichnis</i> , auf dem gearbeitet wird beziehungsweise auf dessen Dateien Operationen angewendet werden sollen <sup>127</sup>
<b>ASCII-Text</b>	Text, der <i>ASCII</i> codiert ist <sup>128</sup>
<b>automatisch aktualisierte statische Webseite</b>	<i>Webseite</i> , die automatisch aktualisiert wird; Wahl des Aktualisierungszeitpunkts hängt vom Inhalt der <i>Webseite</i> ab (Teil der im Rahmen dieser Arbeit getroffenen Klassifizierung von <i>Webseiten</i> )
<b>Bash</b>	Shell (= Benutzerschnittstelle, die eingegebene Befehle interpretiert) des GNU-Projekts <sup>129</sup>
<b>Befehl</b>	hier: Anweisung zur Ausführung einer bestimmten Operation an den Computer <sup>130</sup>
<b>Befehlskette</b>	Zusammenkettung von <i>Befehlen</i> ; siehe <i>Befehl</i>
<b>Bison</b>	hier: <i>Parser</i> -Generator des GNU-Projekts <sup>131</sup>
<b>Browser</b>	Programm, mit welchem <i>Webseiten</i> abgerufen werden können; kann grafische Benutzeroberfläche haben oder textbasiert sein <sup>132</sup>
<b>C</b>	Programmiersprache <sup>133</sup>
<b>Cascading Style Sheets</b>	Sprache zum Festlegen der Darstellung von <i>HTML</i> -Dokumenten <sup>134</sup>
<b>cat</b>	hier: Programm aus der Ansammlung der <i>GNU Coreutils</i> zur Ausgabe des Inhalts einer Textdatei <sup>135</sup>
<b>Client</b>	hier: initiiert Verbindungen zu <i>Servern</i> ; ist in der Regel ein <i>Browser</i> <sup>136</sup>
<b>clientseitig generierte dynamische Webseite</b>	<i>Webseite</i> , die beim Abruf vom <i>Client</i> generiert wird; Generierung wird vom <i>Server</i> gesteuert, erfolgt aber beim <i>Client</i> ; clientseitig ausgeführter <i>Quelltext</i> ist für den <i>Client</i> sichtbar (Teil der im Rahmen dieser Arbeit getroffenen Klassifizierung von <i>Webseiten</i> )
<b>Comma-separated values</b>	Format, das verwendet wird, um Dateien zwischen diversen Tabellenkalkulationsprogrammen auszutauschen und zu konvertieren <sup>137</sup>
<b>Content-Type</b>	siehe <i>MIME-Type</i>

<sup>125</sup> vgl. u. a. <https://support.mozilla.org/de/kb/autovervollstaendigung-adresseleiste> (19.01.2016)

<sup>126</sup> vgl. <http://tools.ietf.org/html/rfc20> (22.01.2016)

<sup>127</sup> vgl. u. a. <https://www.dict.cc/?s=Arbeitsverzeichnis> (22.01.2016)

<sup>128</sup> vgl. <http://tools.ietf.org/html/rfc20> (22.01.2016)

<sup>129</sup> vgl. <https://www.gnu.org/software/bash/> (19.01.2016)

<sup>130</sup> vgl. <http://www.duden.de/rechtschreibung/Befehl#Bedeutung1b> (22.01.2016)

<sup>131</sup> vgl. <https://www.gnu.org/software/bison/> (19.01.2016)

<sup>132</sup> vgl. <http://www.duden.de/rechtschreibung/Browser> (22.01.2016)

<sup>133</sup> vgl. u. a. <http://csapp.cs.cmu.edu/3e/docs/chistory.html> (19.01.2016)

<sup>134</sup> vgl. <http://www.w3.org/standards/webdesign/htmlcss#whatcss> (19.01.2016)

<sup>135</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#cat-invocation> (19.01.2016)

<sup>136</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-1.3> (19.01.2016)

<sup>137</sup> vgl. <https://tools.ietf.org/html/rfc4180#page-2> (19.01.2016)

<b>Cookie</b>	Möglichkeit, <i>Session</i> -Information zu codieren; ist ein Datensatz, mit welchem Benutzer einer <i>Webseite</i> identifiziert werden kann <sup>138</sup>
<b>Cronjob</b>	automatisiert zeitgesteuert ausgeführte wiederkehrende Aufgabe <sup>139</sup>
<b>curl</b>	<i>Kommandozeilenprogramm</i> zum Download von <i>Webseiten</i> <sup>140</sup>
<b>cut</b>	Programm aus der Ansammlung der <i>GNU Coreutils</i> zur zeilenweisen Extraktion von Spalten aus einer Eingabe <sup>141</sup>
<b>Dateipfad</b>	Ort, an welchem sich eine Datei befindet <sup>142</sup>
<b>Datenstrom</b>	kontinuierlicher Fluss von Datensätzen <sup>143</sup>
<b>Delimiter</b>	Trennzeichen <sup>144</sup>
<b>Dokumenttyp</b>	Klasse von ähnlichen Dokumenten <sup>145</sup>
<b>dynamische Webseite</b>	<i>Webseite</i> , die beim Abruf generiert wird (Teil der im Rahmen dieser Arbeit getroffenen Klassifizierung von <i>Webseiten</i> )
<b>Ecma (International)</b>	(European Computer Manufacturers Association); Organisation zur Normung von Informations- und Kommunikationssystemen und Unterhaltungselektronik <sup>146</sup>
<b>ECMAScript</b>	Sprachkern von <i>JavaScript</i> <sup>147</sup>
<b>Eingabeparameter</b>	siehe <i>Parameter</i>
<b>Einloggen</b>	siehe <i>Login</i>
<b>enriched text</b>	<i>ASCII-Text</i> mit einfachen Formatierungen <sup>148</sup>
<b>File Transfer Protocol</b>	<i>Protokoll</i> , das hauptsächlich für den Transfer von Dateien verwendet wird <sup>149</sup>
<b>Flex</b>	hier: (Beispiel für ein) Programm, das eine <i>lexikalische Analyse</i> durchführt <sup>150</sup>
<b>GET</b>	Methode einer <i>Request-Message</i> <sup>151</sup>
<b>GIF-Format</b>	(Graphics Interchange Format); Grafikformat für Rasterbilder <sup>152</sup>
<b>GNU Coreutils</b>	Ansammlung von diversen Programmen zur Textverarbeitung <sup>153</sup>
<b>GNU General Public License</b>	Softwarelizenz, die einer Person ermöglicht, Software auszuführen, zu studieren, zu ändern und zu verbreiten (kopieren) <sup>154</sup>
<b>Google Chrome</b>	<i>Browser</i> von Google Inc. <sup>155</sup>

<sup>138</sup> vgl. <https://tools.ietf.org/html/rfc6265> (19.01.2016) und

<http://www.duden.de/rechtschreibung/Cookie> (22.01.2016)

<sup>139</sup> vgl. u. a. <http://pubs.opengroup.org/onlinepubs/9699919799/utilities/crontab.html> (19.01.2016)

<sup>140</sup> vgl. <http://curl.haxx.se/> (19.01.2016)

<sup>141</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#cut-invocation> (19.01.2016)

<sup>142</sup> vgl. <http://www.linfo.org/path.html> (22.01.2016)

<sup>143</sup> vgl. [http://www.its.bldrdoc.gov/fs-1037/dir-010/\\_1451.htm](http://www.its.bldrdoc.gov/fs-1037/dir-010/_1451.htm) (19.01.2016)

<sup>144</sup> vgl. [http://www.its.bldrdoc.gov/fs-1037/dir-011/\\_1544.htm](http://www.its.bldrdoc.gov/fs-1037/dir-011/_1544.htm) (22.01.2016)

<sup>145</sup> vgl. <http://www.duden.de/rechtschreibung/Dokumenttyp> (22.01.2016)

<sup>146</sup> vgl. <http://www.ecma-international.org/> (22.01.2016)

<sup>147</sup> vgl. <http://ecma-international.org/ecma-262/5.1/> (19.01.2016)

<sup>148</sup> vgl. <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)

<sup>149</sup> vgl. <https://tools.ietf.org/html/rfc354> (19.01.2016)

<sup>150</sup> vgl. <http://flex.sourceforge.net/> (19.01.2016)

<sup>151</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-9.3> (19.01.2016)

<sup>152</sup> vgl. <https://www.w3.org/Graphics/GIF/spec-gif89a.txt> (23.01.2016)

<sup>153</sup> vgl. <http://www.gnu.org/software/coreutils/coreutils.html> (22.01.2016)

<sup>154</sup> vgl. u. a. <http://www.gnu.org/licenses/gpl-3.0.html>

<sup>155</sup> vgl. <https://www.google.de/chrome/browser/desktop/index.html> (19.01.2016)

<b>Google-Suche</b>	Suche im Internet mithilfe der Suchmaschine Google von Google Inc. <sup>156</sup>
<b>grammatische Analyse</b>	auch <i>syntaktische Analyse</i> ; Überprüfung von Eingabe auf <i>syntaktische</i> Regeln, an die sie sich halten <sup>157</sup>
<b>grep</b>	Programm des GNU-Projekts, das Dateien nach bestimmten Suchmustern durchsucht <sup>158</sup>
<b>head</b>	Programm, das die ersten zehn (oder eine andere Anzahl an) Zeilen einer Textdatei auf dem Bildschirm ausgibt <sup>159</sup>
<b>Hidden Field</b>	deutsch: unsichtbares Formularfeld; Möglichkeit, <i>Session</i> -Information unterzubringen <sup>160</sup>
<b>Homepage</b>	beispielsweise die persönliche <i>Webseite</i> einer Person <sup>161</sup>
<b>html.parser</b>	in der Programmiersprache <i>Python</i> geschriebener <i>HTML-Parser</i> <sup>162</sup>
<b>HTML-Parser</b>	<i>Parser</i> für Daten im <i>HTML</i> -Format; siehe <i>HTML</i> ; siehe <i>Parser</i>
<b>HTML-Tag</b>	siehe <i>Tag</i> (Anmerkung: Definition nur auf <i>HTML</i> bezogen)
<b>HTTPS</b>	verschlüsselte Variante des <i>HTTP</i> ; „S“ steht für „secure“ (sicher) <sup>163</sup>
<b>hybride dynamische Webseite</b>	<i>Webseite</i> , die sowohl Elemente einer <i>clientseitig generierten dynamischen Webseite</i> als auch Elemente einer <i>serverseitig generierten dynamischen Webseite</i> beinhaltet; Zugänglichkeit des <i>Quelltexts</i> für den <i>Client</i> ist elementabhängig; Generierung erfolgt teilweise auf dem <i>Server</i> , teilweise beim <i>Client</i> (Teil der im Rahmen dieser Arbeit getroffenen Klassifizierung von <i>Webseiten</i> )
<b>Hypertext</b>	über <i>Links</i> verbundenes Netz aus Bildern und Text, in dem sich Benutzerinnen und Benutzer frei bewegen können <sup>164</sup>
<b>Hypertext Markup Language</b>	Sprache zum Beschreiben der Struktur von aus <i>Hypertext</i> bestehenden <i>Webseiten</i> <sup>165</sup>
<b>Hypertext Transfer Protocol</b>	<i>zustandsloses Protokoll</i> zur Übertragung von Daten über ein <i>Rechnernetz</i> <sup>166</sup>
<b>Internet Engineering Task Force</b>	Organisation, welche mit der Weiterentwicklung des Internets befasst ist <sup>167</sup>
<b>Internet Message Access Protocol</b>	<i>Protokoll</i> , gemäß welchem E-Mails abgerufen werden <sup>168</sup>
<b>Internetadresse</b>	aus standardisierten Zeichen bestehende Angabe, unter der jemand im Internet erreichbar ist <sup>169</sup>
<b>Internetsurfen</b>	Tätigkeit der wahllosen oder gezielten Suche von Informationen im Internet <sup>170</sup>

<sup>156</sup> vgl. <https://www.google.com> (19.01.2016)

<sup>157</sup> vgl. <http://www.iue.tuwien.ac.at/phd/demel/node107.html> (19.01.2016)

<sup>158</sup> vgl. <http://www.gnu.org/software/grep/manual/grep.html> (19.01.2016)

<sup>159</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#head-invocation> (19.01.2016)

<sup>160</sup> vgl. <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 42 (19.01.2016)

<sup>161</sup> vgl. <http://www.duden.de/rechtschreibung/Homepage> (22.01.2016)

<sup>162</sup> vgl. <https://docs.python.org/2/library/htmlparser.html> (19.01.2016)

<sup>163</sup> vgl. <https://tools.ietf.org/html/rfc2818> (19.01.2016)

<sup>164</sup> vgl. <http://www.duden.de/rechtschreibung/Hypertext> (22.01.2016)

<sup>165</sup> vgl. <http://www.w3.org/standards/webdesign/htmlcss#whathtml> (19.01.2016)

<sup>166</sup> vgl. <https://tools.ietf.org/html/rfc2616>, Abstract (19.01.2016)

<sup>167</sup> vgl. <https://www.ietf.org/> (19.01.2016)

<sup>168</sup> vgl. <https://tools.ietf.org/html/rfc3501> (19.01.2016)

<sup>169</sup> vgl. <http://www.duden.de/rechtschreibung/Internetadresse> (22.01.2016)

<sup>170</sup> vgl. <http://www.duden.de/rechtschreibung/Internetsurfen> (22.01.2016)



<b>JavaScript</b>	Programmiersprache, mit der Webanwendungen realisiert werden können <sup>171</sup>
<b>JavaServer Pages</b>	von Sun Microsystems Inc. entwickelte Programmiersprache zur Realisierung von Webanwendungen <sup>172</sup>
<b>JPEG-Format</b>	Grafikformat für Rasterbilder <sup>173</sup>
<b>Kommandozeile</b>	auf dem Bildschirm gekennzeichnete Stelle, an der man <i>Befehle</i> (= Kommandos) zum Starten von Programmen eingeben kann <sup>174</sup>
<b>Kommandozeilenprogramm</b>	Programm ohne grafische Benutzeroberfläche, welches über die <i>Kommandozeile</i> bedient wird <sup>175</sup>
<b>lexikalische Analyse</b>	Zerlegung von <i>plain text</i> in <i>Tokens</i> <sup>176</sup>
<b>Lightweight Directory Access Protocol</b>	<i>Protokoll</i> zur Abfrage und Änderung von Informationen in einem <i>Verzeichnisdienst</i> <sup>177</sup>
<b>Link</b>	Verknüpfung mit einer anderen Datei oder einer anderen Stelle in der gleichen Datei <sup>178</sup>
<b>Login</b>	deutsch: <i>Einloggen</i> <sup>179</sup>
<b>Login-Bereich</b>	Bereich, der ein <i>Login</i> erfordert, um Zugriff zu gewähren <sup>180</sup>
<b>Lynx</b>	textbasierter <i>Browser</i> <sup>181</sup>
<b>manuell aktualisierte statische Webseite</b>	<i>Webseite</i> , die manuell verändert wird (Teil der im Rahmen dieser Arbeit getroffenen Klassifizierung von <i>Webseiten</i> )
<b>Medientyp</b>	siehe <i>MIME-Type</i>
<b>Metadaten</b>	Daten, die Informationen über andere Daten enthalten <sup>182</sup>
<b>Microsoft Edge</b>	<i>Browser</i> der Microsoft Corporation <sup>183</sup>
<b>Microsoft Excel</b>	Tabellenkalkulationsprogramm der Microsoft Corporation <sup>184</sup>
<b>Microsoft Windows XP</b>	Betriebssystem der Microsoft Corporation aus dem Jahr 2001 <sup>185</sup>
<b>MIME-Type</b>	gibt (bei der Übertragung gemäß <i>HTTP</i> ) an, um welche Art von Daten es sich handelt <sup>186</sup>
<b>Mozilla Firefox</b>	<i>Browser</i> der Mozilla Foundation <sup>187</sup>
<b>Network Time Protocol</b>	<i>Protokoll</i> , das für die Synchronisation von Uhren in Computersystemen von Bedeutung ist <sup>188</sup>

<sup>171</sup> vgl. <http://www.ecma-international.org/publications/files/ECMA-ST/Ecma-262.pdf> (19.01.2016) und <https://developer.mozilla.org/en-US/docs/Web/JavaScript> (19.01.2016)

<sup>172</sup> vgl. <http://www.oracle.com/technetwork/java/javaee/jsp/index.html> (19.01.2016)

<sup>173</sup> vgl. <http://jpeg.org/> (23.01.2016)

<sup>174</sup> vgl. <http://www.duden.de/rechtschreibung/Kommandozeile> (22.01.2016)

<sup>175</sup> vgl. <http://www.duden.de/rechtschreibung/Kommandozeile> (22.01.2016)

<sup>176</sup> vgl. u. a. <https://www.ibm.com/developerworks/community/blogs/nlp/entry/tokenization?lang=en> (22.01.2016)

<sup>177</sup> vgl. <https://tools.ietf.org/html/rfc4511> (19.01.2016)

<sup>178</sup> vgl. <http://www.duden.de/rechtschreibung/Link> (22.01.2016)

<sup>179</sup> vgl. [http://www.duden.de/rechtschreibung/Log\\_in](http://www.duden.de/rechtschreibung/Log_in) (22.01.2016)

<sup>180</sup> vgl. [http://www.duden.de/rechtschreibung/Log\\_in](http://www.duden.de/rechtschreibung/Log_in) (22.01.2016)

<sup>181</sup> vgl. <http://lynx.invisible-island.net/> (19.01.2016)

<sup>182</sup> vgl. <https://wiki.selfhtml.org/wiki/HTML/Kopfdaten/meta> (22.01.2016)

<sup>183</sup> vgl. <https://www.microsoft.com/de-at/windows/microsoft-edge> (19.01.2016)

<sup>184</sup> vgl. <https://products.office.com/de-at/excel> (19.01.2016)

<sup>185</sup> vgl. <http://windows.microsoft.com/de-de/windows/windows-help#windows=windows-xp> (21.01.2016)

<sup>186</sup> vgl. <http://www.w3.org/2001/tag/2002/0129-mime> (19.01.2016)

<sup>187</sup> vgl. <https://www.mozilla.org/de/firefox/new/> (19.01.2016)

<sup>188</sup> vgl. <https://tools.ietf.org/html/rfc5905> (19.01.2016)

<b>Open Source (Open-Source-Software)</b>	Software, deren Quellcode frei zugänglich ist <sup>189</sup>
<b>Opera</b>	<i>Browser</i> von Opera Software <sup>190</sup>
<b>Option</b>	siehe <i>Parameter</i>
<b>Parameter</b>	für einen Aufruf gültige Einstellungsmöglichkeiten eines Computerprogramms; <i>HTML</i> : genauere Spezifikation von <i>Tags</i> <sup>191</sup>
<b>Parser</b>	Programm, das eine beliebige Eingabe zerlegt und in ein für die Weiterverarbeitung brauchbares Format umwandelt <sup>192</sup>
<b>Pfad</b>	siehe <i>Dateipfad</i>
<b>PHP Hypertext Processor</b>	Programmiersprache, mit der Webanwendungen realisiert werden können <sup>193</sup>
<b>Pipe</b>	siehe <i>Pipeline</i>
<b>Pipeline</b>	auch: <i>Pipe</i> ; ermöglicht Umleitung der Ausgabe eines Programms in die Eingabe eines anderen Programms <sup>194</sup>
<b>plain text</b>	deutsch: einfacher, schlichter Text; Daten, die direkt in Text umgesetzt werden können <sup>195</sup>
<b>POST</b>	Methode einer <i>Request-Message</i> <sup>196</sup>
<b>Post Office Protocol</b>	<i>Protokoll</i> , gemäß welchem E-Mails abgerufen werden <sup>197</sup>
<b>PostScript-Format</b>	Sprache für Vektorgrafiken, die beschreibt wie Seiten auf Geräten (wie zum Beispiel Druckern) aussehen sollen <sup>198</sup>
<b>Protokoll</b>	hier: Kommunikationsvorschrift; schreibt vor, wie eine Kommunikation abzulaufen hat <sup>199</sup>
<b>Python</b>	hier: Programmiersprache <sup>200</sup>
<b>Quelltext</b>	in einer Programmiersprache geschriebene Abfolge von Programmanweisungen, die vom Menschen gelesen, aber erst nach einer elektronischen Übersetzung vom Computer verarbeitet werden können <sup>201</sup>
<b>Rechnernetz</b>	Vernetzung mehrerer voneinander unabhängiger Rechner, die den Datenaustausch zwischen besagten Rechnern ermöglicht <sup>202</sup>
<b>Regular Expression</b>	Ausdruck zur kontextunabhängigen Beschreibung von <i>Zeichenketten</i> <sup>203</sup>
<b>Regulärer Ausdruck</b>	siehe <i>Regular Expression</i>

<sup>189</sup> vgl. [http://www.duden.de/rechtschreibung/Open\\_Source\\_Software](http://www.duden.de/rechtschreibung/Open_Source_Software) (22.01.2016)

<sup>190</sup> vgl. <http://www.opera.com/de> (19.01.2016)

<sup>191</sup> vgl. u. a. <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2005/n1905.pdf>, Seite 16, Kap. 1.3.1 argument (22.01.2016) und <http://www.w3.org/standards/webdesign/htmlcss#whathtml> (19.01.2016)

<sup>192</sup> vgl. [https://www2.informatik.uni-erlangen.de/EN/teaching/SS2012/HalloWelt/PA\\_2012.pdf](https://www2.informatik.uni-erlangen.de/EN/teaching/SS2012/HalloWelt/PA_2012.pdf), Seite 3 (22.01.2016)

<sup>193</sup> vgl. <http://php.net/> (19.01.2016)

<sup>194</sup> vgl. u. a. <http://www.linfo.org/pipe.html> (22.01.2016)

<sup>195</sup> vgl. <https://tools.ietf.org/html/rfc2046#section-4.1.3> (19.01.2016)

<sup>196</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-9.5> (19.01.2016)

<sup>197</sup> vgl. <https://tools.ietf.org/html/rfc1939> (19.01.2016)

<sup>198</sup> vgl. <http://www.adobe.com/devnet/postscript.html> (23.01.2016)

<sup>199</sup> vgl. <http://www.duden.de/rechtschreibung/Protokoll#Bedeutung3a> (22.01.2016)

<sup>200</sup> vgl. <https://www.python.org/> (21.01.2016)

<sup>201</sup> vgl. <http://www.duden.de/rechtschreibung/Quellcode> (22.01.2016)

<sup>202</sup> vgl. <http://www.duden.de/rechtschreibung/Netzwerk#b2-Bedeutung-3> (22.01.2016)

<sup>203</sup> vgl. u. a. [http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd\\_chap09.html](http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd_chap09.html) (19.01.2016)

<b>rekursives Herunterladen</b>	bewirkt, dass alle Elemente die von der angegebenen <i>URL</i> <i>verlinkt</i> werden, ebenfalls heruntergeladen werden <sup>204</sup>
<b>Request</b>	deutsch: <i>Anfrage</i> ; wird in Form einer <i>Request-Message</i> ( <i>Anfragenachricht</i> ) vom <i>Client</i> zum <i>Server</i> geschickt <sup>205</sup>
<b>Request for Comments</b>	<i>RFCs</i> sind eine Reihe technischer und organisatorischer Dokumente das Internet betreffend <sup>206</sup>
<b>Request-Message</b>	beschreibt <i>Request</i> des <i>Clients</i> an den <i>Server</i> ; hält sich an bestimmtes Schema <sup>207</sup> (siehe Kapitel 2.2.3)
<b>Request-URL</b>	siehe <i>Request</i> ; siehe <i>URL</i>
<b>Response</b>	deutsch: <i>Antwort</i> ; wird in Form einer <i>Response-Message</i> ( <i>Antwortnachricht</i> ) vom <i>Server</i> an den <i>Client</i> geschickt <sup>208</sup>
<b>Response-Message</b>	deutsch: <i>Antwortnachricht</i> ; beschreibt <i>Antwort</i> des <i>Servers</i> an den <i>Client</i> ; hält sich an bestimmtes Schema <sup>209</sup> (siehe Kapitel 2.2.4)
<b>Safari</b>	<i>Browser</i> von Apple Inc. <sup>210</sup>
<b>sed</b>	( <i>Stream Editor</i> ); Programm des GNU-Projekts, mit dem Text- <i>Datenströme</i> editiert werden <sup>211</sup>
<b>Selenium</b>	Ansammlung von <i>Tools</i> zur Automatisierung von <i>Browsern</i> <sup>212</sup>
<b>Selenium IDE</b>	<i>Tool</i> aus der Ansammlung der <i>Selenium-Tools</i> <sup>213</sup>
<b>Selenium WebDriver</b>	<i>Tool</i> aus der Ansammlung der <i>Selenium-Tools</i> <sup>214</sup>
<b>Server</b>	hier: wartet auf Verbindungsanfragen von <i>Clients</i> ; antwortet auf diese durch Liefern der gewünschten Ressource oder Statusinformation <sup>215</sup>
<b>serverseitig generierte dynamische Webseite</b>	<i>Webseite</i> , die beim Abruf vom Server generiert wird; ausgeführter <i>Quelltext</i> ist für den <i>Client</i> nicht sichtbar; <i>Client</i> erhält lediglich die Ausgabe dieser Ausführung (Teil der im Rahmen dieser Arbeit getroffenen Klassifizierung von <i>Webseiten</i> )
<b>Session</b>	deutsch: <i>Sitzung</i> ; Dialog, welcher sich über mehrere Zyklen von <i>Request</i> und <i>Response</i> erstreckt; innerhalb einer <i>Session</i> soll Rückbezug auf vergangene Kommunikation möglich sein <sup>216</sup>
<b>Simple Mail Transfer Protocol</b>	<i>Protokoll</i> , gemäß welchem E-Mails im Internet übertragen werden <sup>217</sup>
<b>Sitzung</b>	hier: siehe <i>Session</i>
<b>Skript</b>	hier: Liste von <i>Bash-Befehlen</i> , die der Reihe nach ausgeführt werden sollen; kleines Computerprogramm <sup>218</sup>

<sup>204</sup> vgl. <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Download> (22.01.2016)

<sup>205</sup> vgl. <http://www.w3.org/Protocols/rfc2616/rfc2616-sec5.html> (19.01.2016)

<sup>206</sup> vgl. <http://www.rfc-editor.org/> (22.01.2016)

<sup>207</sup> vgl. <http://www.w3.org/Protocols/rfc2616/rfc2616-sec5.html> (19.01.2016)

<sup>208</sup> vgl. <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)

<sup>209</sup> vgl. <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)

<sup>210</sup> vgl. <http://www.apple.com/at/safari/> (19.01.2016)

<sup>211</sup> vgl. <https://www.gnu.org/software/sed/> (19.01.2016)

<sup>212</sup> vgl. <http://docs.seleniumhq.org/> (19.01.2016)

<sup>213</sup> vgl. [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp) (19.01.2016)

<sup>214</sup> vgl. [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp) (19.01.2016)

<sup>215</sup> vgl. <http://tools.ietf.org/html/rfc2616#section-1.3> (19.01.2016)

<sup>216</sup> vgl. <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 40 (19.01.2016)

<sup>217</sup> vgl. <https://tools.ietf.org/html/rfc5321> (19.01.2016)

<sup>218</sup> vgl. u. a. <http://www.duden.de/rechtschreibung/Skript#Bedeutung4> (22.01.2016)

<b>Slash</b>	deutsch: Schrägstrich; von rechts oben nach links unten verlaufender Strich; Zeichen: „/“ <sup>219</sup> ( <i>ASCII</i> Code 47)
<b>sort</b>	Programm aus der Ansammlung der <i>GNU Coreutils</i> zur zeilenweisen Sortierung von Dateien <sup>220</sup>
<b>Standardausgabe</b>	Programm kann über diese Daten ausgeben; normalerweise mit dem Bildschirm verbunden. <sup>221</sup>
<b>statische Webseite</b>	<i>Webseite</i> , die statisch ist und nicht beim Abruf generiert wird (Teil der im Rahmen dieser Arbeit getroffenen Klassifizierung von <i>Webseiten</i> )
<b>syntaktische Analyse</b>	siehe <i>grammatische Analyse</i>
<b>Syntax</b>	hier: korrekte Verknüpfung sprachlicher Einheiten <sup>222</sup>
<b>Tabulator</b>	Durch Drücken der Tabulatortaste [↵] eingefügtes Zeichen <sup>223</sup> ( <i>ASCII</i> Code 9)
<b>Tag</b>	Markierungselement von Beschreibungssprachen (hier: <i>HTML</i> ) zur Strukturierung der Dokumente <sup>224</sup>
<b>tail</b>	Programm, das die letzten zehn (oder eine andere Anzahl an) Zeilen einer Textdatei auf dem Bildschirm ausgibt <sup>225</sup>
<b>Terminal</b>	siehe <i>Kommandozeile</i>
<b>Token</b>	Folge von logisch zusammengehörigen Einheiten (beispielsweise Folge von Zeichen) <sup>226</sup>
<b>Tool</b>	deutsch: <i>Werkzeug</i> ; Computerprogramm (in der Regel von geringerem Umfang), das zusätzliche Aufgaben für ein bestimmtes Betriebssystem übernimmt <sup>227</sup>
<b>Top-Level-Typ</b>	übergeordneter Bestandteil eines <i>MIME-Types</i> ; <sup>228</sup> Beispiele: text, image, audio, application
<b>tr</b>	Programm aus der Ansammlung der <i>GNU Coreutils</i> zum Ersetzen oder Löschen von Zeichen einer Eingabe <sup>229</sup>
<b>Ubuntu Linux 14.04.3 LTS</b>	Auf dem Linux-Kernel (Linux-Kern) basierendes Betriebssystem von Canonical Ltd. aus dem Jahr 2014 <sup>230</sup>
<b>Uniform Resource Locator</b>	Schema für die Angabe einer <i>Internetadresse</i> ; gibt an, welche Daten von welchem Computer unter Verwendung des angegebenen <i>Protokolls</i> übertragen werden sollen <sup>231</sup>
<b>Untertyp</b>	untergeordneter Bestandteil eines <i>MIME-Types</i> ; <sup>232</sup> Beispiele für <i>Untertypen</i> des <i>Top-Level-Typs</i> „image“: gif, jpeg

<sup>219</sup> vgl. <http://www.duden.de/rechtschreibung/Schraegstrich#b2-Bedeutung-b> (22.01.2016)

<sup>220</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#sort-invocation> (19.01.2016)

<sup>221</sup> vgl. [http://www.linfo.org/standard\\_output.html](http://www.linfo.org/standard_output.html) (29.01.2016)

<sup>222</sup> vgl. <http://www.duden.de/rechtschreibung/Syntax> (22.01.2016)

<sup>223</sup> vgl. u. a. <http://www-01.ibm.com/software/globalization/topics/keyboards/physical.html>

<sup>224</sup> vgl. [http://www.duden.de/rechtschreibung/Tag\\_Strukturelement\\_Markierung](http://www.duden.de/rechtschreibung/Tag_Strukturelement_Markierung) (22.01.2016)

<sup>225</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#tail-invocation> (19.01.2016)

<sup>226</sup> vgl. <http://www.duden.de/rechtschreibung/Token> (22.01.2016) und

<https://www.ibm.com/developerworks/community/blogs/nlp/entry/tokenization?lang=en> (22.01.2016)

<sup>227</sup> vgl. <http://www.duden.de/rechtschreibung/Tool> (22.01.2016)

<sup>228</sup> vgl. <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)

<sup>229</sup> vgl. <https://www.gnu.org/software/coreutils/manual/coreutils.html#tr-invocation> (19.01.2016)

<sup>230</sup> vgl. <http://releases.ubuntu.com/14.04/> (21.01.2016)

<sup>231</sup> vgl. <https://tools.ietf.org/html/rfc1738> (19.01.2016)

<sup>232</sup> vgl. <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)

<b>URL Rewriting</b>	Technik zur Codierung von <i>Session</i> -Information; <i>Session</i> -Information wird hierbei in <i>URL</i> codiert <sup>233</sup>
<b>URL-Syntax</b>	siehe <i>URL</i> ; siehe <i>Syntax</i>
<b>UTF-8</b>	genormter Code für Ziffern, Buchstaben und Sonderzeichen <sup>234</sup>
<b>Verzeichnis</b>	Teil des Dateisystems, in dem Dateien abgelegt werden <sup>235</sup>
<b>Webbrowser</b>	siehe <i>Browser</i>
<b>Webseite</b>	Teil einer <i>Website</i> <sup>236</sup>
<b>Webserver</b>	siehe <i>Server</i>
<b>Website</b>	Gesamtheit der hinter einer <i>Internetadresse</i> stehenden <i>Webseiten</i> <sup>237</sup>
<b>Werkzeug</b>	hier: siehe <i>Tool</i>
<b>wget</b>	<i>Kommandozeilenprogramm</i> des GNU-Projekts zum Download von <i>Webseiten</i> <sup>238</sup>
<b>Windows Internet Explorer</b>	<i>Browser</i> der Microsoft Corporation <sup>239</sup>
<b>World Wide Web Consortium</b>	Konsortium zur Erstellung von Web-Standards <sup>240</sup>
<b>Zeichenkette</b>	Folge von Zeichen <sup>241</sup>
<b>Zeichenfolge</b>	siehe <i>Zeichenkette</i>
<b>zustandslos</b>	hier: alle <i>Anfragen</i> als voneinander unabhängig behandelnd <sup>242</sup>

<sup>233</sup> vgl. <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 42 (19.01.2016)

<sup>234</sup> vgl. <https://tools.ietf.org/html/rfc3629> (19.01.2016)

<sup>235</sup> vgl. u. a. <http://www.duden.de/rechtschreibung/Ordner#b2-Bedeutung-3> (22.01.2016)

<sup>236</sup> vgl. <http://www.duden.de/rechtschreibung/Seite#b2-Bedeutung-11> (22.01.2016)

<sup>237</sup> vgl. <http://www.duden.de/rechtschreibung/Website> (22.01.2016)

<sup>238</sup> vgl. <https://www.gnu.org/software/wget/> (19.01.2016)

<sup>239</sup> vgl. <http://windows.microsoft.com/de-at/internet-explorer/download-ie> (19.01.2016)

<sup>240</sup> vgl. <http://www.w3.org/> (19.01.2016)

<sup>241</sup> vgl. <http://www.duden.de/rechtschreibung/Zeichenkette> (22.01.2016)

<sup>242</sup> vgl. <https://tools.ietf.org/html/rfc2616>, Abstract (19.01.2016)

## C Literaturverzeichnis

**Legende:**

Fußnotennummer      URL (Datum des letzten Zugriffs)

- 1    <https://www.gnu.org/software/bash/manual/bashref.html> (19.01.2016)
- 2    <https://www.python.org/> (21.01.2016)
- 3    <http://windows.microsoft.com/de-de/windows/windows-help#windows=windows-xp>  
(21.01.2016)
- 4    <http://releases.ubuntu.com/14.04/> (21.01.2016)
- 5    <http://www.w3.org/> (19.01.2016)
- 6    <https://www.ietf.org/> (19.01.2016)
- 7    <https://www.mozilla.org/de/firefox/new/> (19.01.2016)
- 8    <https://www.google.de/chrome/browser/desktop/index.html> (19.01.2016)
- 9    <http://windows.microsoft.com/de-at/internet-explorer/download-ie> (19.01.2016)
- 10   <https://www.microsoft.com/de-at/windows/microsoft-edge> (19.01.2016)
- 11   <http://www.apple.com/at/safari/> (19.01.2016)
- 12   <http://www.opera.com/de> (19.01.2016)
- 13   <https://www.google.com> (19.01.2016)
- 14   <https://tools.ietf.org/html/rfc1738> (19.01.2016)
- 15   <https://tools.ietf.org/html/rfc1738#section-2> (19.01.2016)
- 16   <https://tools.ietf.org/html/rfc1738#section-2.1> (19.01.2016)
- 17   <https://tools.ietf.org/html/rfc354> (19.01.2016)
- 18   <https://tools.ietf.org/html/rfc5321> (19.01.2016)
- 19   <https://tools.ietf.org/html/rfc1939> (19.01.2016)
- 20   <https://tools.ietf.org/html/rfc3501> (19.01.2016)
- 21   <https://tools.ietf.org/html/rfc4511> (19.01.2016)
- 22   <https://tools.ietf.org/html/rfc5905> (19.01.2016)
- 23   <https://tools.ietf.org/html/rfc2616>, Abstract (19.01.2016)
- 24   <http://tools.ietf.org/html/rfc2616#section-1.3> (19.01.2016)
- 25   <https://www.rfc-editor.org/rfc/rfc793.txt> (19.01.2016)
- 26   <http://www.w3.org/Protocols/rfc2616/rfc2616-sec5.html> (19.01.2016)
- 27   <http://tools.ietf.org/html/rfc2616#section-5.1.1> (19.01.2016)
- 28   <http://tools.ietf.org/html/rfc3986#section-1.1.3> (19.01.2016)
- 29   <http://tools.ietf.org/html/rfc2616#section-9.3> (19.01.2016)
- 30   <http://tools.ietf.org/html/rfc2616#section-9.5> (19.01.2016)

- 31 <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)
- 32 <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)
- 33 <http://www.iana.org/assignments/http-status-codes/http-status-codes.xhtml> (19.01.2016)
- 34 <http://www.iana.org/assignments/http-status-codes/http-status-codes.xhtml> (19.01.2016)
- 35 <https://tools.ietf.org/html/rfc2818> (19.01.2016)
- 36 <https://tools.ietf.org/html/rfc2818#section-1> (19.01.2016)
- 37 <http://tools.ietf.org/html/rfc7230#section-2.3> (19.01.2016)
- 38 <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 40 (19.01.2016)
- 39 <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 42 (19.01.2016)
- 40 <https://tools.ietf.org/html/rfc6265> (19.01.2016)
- 41 <http://www.w3.org/2001/tag/2002/0129-mime> (19.01.2016)
- 42 <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)
- 43 <http://www.w3.org/standards/webdesign/htmlcss#whathtml> (19.01.2016)
- 44 <http://www.w3.org/TR/html401/struct/global.html#h-7.1> (19.01.2016)
- 45 <http://www.w3.org/TR/html401/sgml/dtd.html> (19.01.2016)
- 46 <http://www.w3.org/standards/webdesign/htmlcss#whatcss> (19.01.2016)
- 47 <http://www.w3.org/TR/CSS1/#abstract> (19.01.2016)
- 48 <http://www.w3.org/TR/CSS21/intro.html#html-tutorial> (19.01.2016)
- 49 u. a. <https://support.mozilla.org/de/kb/Einstellungen-Fenster--Inhalts-Abschnitt> (19.01.2016)
- 50 u. a. <https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference> (19.01.2016)
- 51 <http://ecma-international.org/ecma-262/5.1/> (19.01.2016)
- 52 <https://products.office.com/de-at/excel> (19.01.2016)
- 53 <http://lynx.invisible-island.net/> (19.01.2016)
- 54 <http://redmonk.com/sogrady/2013/07/25/language-rankings-6-13/> (19.01.2016)
- 55 <http://php.net/> (19.01.2016)
- 56 <http://www.asp.net/> (19.01.2016)
- 57 <http://www.oracle.com/technetwork/java/javaee/jsp/index.html> (19.01.2016)
- 58 <http://php.net/manual/de/intro-what-is.php> (19.01.2016)
- 59 <http://redmonk.com/sogrady/2013/07/25/language-rankings-6-13/> (19.01.2016)
- 60 <http://www.ecma-international.org/publications/files/ECMA-ST/Ecma-262.pdf> (19.01.2016)  
und <https://developer.mozilla.org/en-US/docs/Web/JavaScript> (19.01.2016)
- 61 <https://developer.mozilla.org/en-US/docs/Web/JavaScript> (19.01.2016)
- 62 <http://www.ecma-international.org/publications/files/ECMA-ST/Ecma-262.pdf> (19.01.2016)
- 63 <https://www.gnu.org/software/wget/> (19.01.2016)
- 64 <http://www.gnu.org/licenses/gpl-3.0.de.html> (19.01.2016)
- 65 u. a. <http://csapp.cs.cmu.edu/3e/docs/chistory.html> (19.01.2016)

- 66 u. a. <http://pubs.opengroup.org/onlinepubs/9699919799/utilities/crontab.html> (19.01.2016)
- 67 <https://www.gnu.org/software/wget/manual/wget.html#Overview> (19.01.2016)
- 68 <https://www.gnu.org/software/wget/manual/wget.html#Invoking> (19.01.2016)
- 69 <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options>  
(19.01.2016)
- 70 <https://www.gnu.org/software/wget/manual/wget.html#Basic-Startup-Options> (19.01.2016)
- 71 <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options>  
(19.01.2016)
- 72 <https://www.gnu.org/software/wget/manual/wget.html#Download-Options> (19.01.2016)
- 73 <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options>  
(19.01.2016)
- 74 <https://www.gnu.org/software/wget/manual/wget.html#Download-Options> (19.01.2016)
- 75 <https://www.gnu.org/software/wget/manual/wget.html#Download-Options> (19.01.2016)
- 76 <https://www.gnu.org/software/wget/manual/wget.html#Directory-Options> (19.01.2016)
- 77 <https://www.gnu.org/software/wget/manual/wget.html#Logging-and-Input-File-Options>  
(19.01.2016)
- 78 <https://www.gnu.org/software/wget/manual/wget.html#HTTP-Options> (19.01.2016)
- 79 [https://www.gnu.org/software/wget/manual/wget.html#Recursive-Accept\\_002fReject-Options](https://www.gnu.org/software/wget/manual/wget.html#Recursive-Accept_002fReject-Options)  
(19.02.2016)
- 80 <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Retrieval-Options>  
(19.01.2016)
- 81 <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Retrieval-Options>  
(19.01.2016)
- 82 <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Retrieval-Options>  
(19.02.2016)
- 83 <http://curl.haxx.se/> (19.01.2016)
- 84 <http://curl.haxx.se/>, curl is an open source command line tool and library (19.01.2016)
- 85 <http://curl.haxx.se/docs/manpage.html>, SYNOPSIS (19.01.2016)
- 86 <http://curl.haxx.se/docs/manpage.html>, OPTIONS (19.01.2016)
- 87 <http://docs.seleniumhq.org/> (19.01.2016)
- 88 [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp) (19.01.2016)
- 89 <http://www.seleniumhq.org/docs/> (19.01.2016)
- 90 [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp), IDE Features, Menu Bar (19.01.2016)
- 91 [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp), IDE Features, Toolbar (19.01.2016)
- 92 [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp), IDE Features, Test Case Pane  
(19.01.2016)
- 93 <http://release.seleniumhq.org/selenium-core/1.0.1/reference.html> (19.01.2016)
- 94 <http://www.gnu.org/software/coreutils/coreutils.html> (22.01.2016)



- 95 <https://www.gnu.org/software/coreutils/manual/coreutils.html#cat-invocation> (19.01.2016)
- 96 <https://www.gnu.org/software/coreutils/manual/coreutils.html#cut-invocation> (19.01.2016)
- 97 <http://www.gnu.org/software/grep/manual/grep.html> (19.01.2016)
- 98 <https://www.gnu.org/software/coreutils/manual/coreutils.html#head-invocation> (19.01.2016)
- 99 <https://www.gnu.org/software/coreutils/manual/coreutils.html#tail-invocation> (19.01.2016)
- 100 <https://www.gnu.org/software/sed/> (19.01.2016)
- 101 [http://www.its.bldrdoc.gov/fs-1037/dir-010/\\_1451.htm](http://www.its.bldrdoc.gov/fs-1037/dir-010/_1451.htm) (19.01.2016)
- 102 <https://www.gnu.org/software/sed/manual/sed.html#Introduction> (19.01.2016)
- 103 <https://www.gnu.org/software/sed/manual/sed.html#Invoking-sed> (19.01.2016)
- 104 <https://www.gnu.org/software/sed/manual/sed.html#index-Global-substitution-108>  
(19.01.2016)
- 105 <https://www.gnu.org/software/coreutils/manual/coreutils.html#sort-invocation> (19.01.2016)
- 106 <https://www.gnu.org/software/coreutils/manual/coreutils.html#tr-invocation> (19.01.2016)
- 107 <http://www.iue.tuwien.ac.at/phd/demel/node107.html> (19.01.2016)
- 108 <https://www.python.org/> (19.01.2016)
- 109 <https://docs.python.org/2/library/htmlparser.html> (19.01.2016)
- 110 <https://www.gnu.org/software/bison/> (19.01.2016)
- 111 <http://flex.sourceforge.net/> (19.01.2016)
- 112 <https://products.office.com/de-at/excel> (19.01.2016)
- 113 <https://tools.ietf.org/html/rfc4180#page-2> (19.01.2016)
- 114 <https://www.python.org/> (19.01.2016)
- 115 <https://www.python.org/dev/peps/pep-0008/#inline-comments> (19.01.2016)
- 116 <https://tools.ietf.org/html/rfc3629> (19.01.2016)
- 117 <https://docs.python.org/2/library/codecs.html> (19.01.2016)
- 118 <https://docs.python.org/2/tutorial/inputoutput.html> (19.01.2016)
- 119 <http://www.seleniumhq.org/projects/webdriver/> (19.01.2016)
- 120 <https://docs.python.org/2/tutorial/inputoutput.html> (19.01.2016)
- 121 u. a. [http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd\\_chap09.html](http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd_chap09.html)  
(19.01.2016)
- 122 [https://www.gnu.org/software/sed/manual/html\\_node/Regular-Expressions.html](https://www.gnu.org/software/sed/manual/html_node/Regular-Expressions.html) (20.01.2016)
- 123 <http://www.asp.net/> (19.01.2016)
- 124 [http://www.duden.de/rechtschreibung/Add\\_on](http://www.duden.de/rechtschreibung/Add_on) (19.01.2016)
- 125 u. a. <https://support.mozilla.org/de/kb/autovollstaendigung-adressleiste> (19.01.2016)
- 126 <http://tools.ietf.org/html/rfc20> (22.01.2016)
- 127 u. a. <https://www.dict.cc/?s=Arbeitsverzeichnis> (22.01.2016)
- 128 <http://tools.ietf.org/html/rfc20> (22.01.2016)
- 129 <https://www.gnu.org/software/bash/> (19.01.2016)
- 130 <http://www.duden.de/rechtschreibung/Befehl#Bedeutung1b> (22.01.2016)

- 131 <https://www.gnu.org/software/bison/> (19.01.2016)
- 132 <http://www.duden.de/rechtschreibung/Browser> (22.01.2016)
- 133 u. a. <http://csapp.cs.cmu.edu/3e/docs/chistory.html> (19.01.2016)
- 134 <http://www.w3.org/standards/webdesign/htmlcss#whatcss> (19.01.2016)
- 135 <https://www.gnu.org/software/coreutils/manual/coreutils.html#cat-invocation> (19.01.2016)
- 136 <http://tools.ietf.org/html/rfc2616#section-1.3> (19.01.2016)
- 137 <https://tools.ietf.org/html/rfc4180#page-2> (19.01.2016)
- 138 <https://tools.ietf.org/html/rfc6265> (19.01.2016) und  
<http://www.duden.de/rechtschreibung/Cookie> (22.01.2016)
- 139 u. a. <http://pubs.opengroup.org/onlinepubs/9699919799/utilities/crontab.html> (19.01.2016)
- 140 <http://curl.haxx.se/> (19.01.2016)
- 141 <https://www.gnu.org/software/coreutils/manual/coreutils.html#cut-invocation> (19.01.2016)
- 142 <http://www.linfo.org/path.html> (22.01.2016)
- 143 [http://www.its.bldrdoc.gov/fs-1037/dir-010/\\_1451.htm](http://www.its.bldrdoc.gov/fs-1037/dir-010/_1451.htm) (19.01.2016)
- 144 [http://www.its.bldrdoc.gov/fs-1037/dir-011/\\_1544.htm](http://www.its.bldrdoc.gov/fs-1037/dir-011/_1544.htm) (22.01.2016)
- 145 <http://www.duden.de/rechtschreibung/Dokumenttyp> (22.01.2016)
- 146 <http://www.ecma-international.org/> (22.01.2016)
- 147 <http://ecma-international.org/ecma-262/5.1/> (19.01.2016)
- 148 <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)
- 149 <https://tools.ietf.org/html/rfc354> (19.01.2016)
- 150 <http://flex.sourceforge.net/> (19.01.2016)
- 151 <http://tools.ietf.org/html/rfc2616#section-9.3> (19.01.2016)
- 152 <https://www.w3.org/Graphics/GIF/spec-gif89a.txt> (23.01.2016)
- 153 <http://www.gnu.org/software/coreutils/coreutils.html> (22.01.2016)
- 154 u. a. <http://www.gnu.org/licenses/gpl-3.0.html>
- 155 <https://www.google.de/chrome/browser/desktop/index.html> (19.01.2016)
- 156 <https://www.google.com> (19.01.2016)
- 157 <http://www.iue.tuwien.ac.at/phd/demel/node107.html> (19.01.2016)
- 158 <http://www.gnu.org/software/grep/manual/grep.html> (19.01.2016)
- 159 <https://www.gnu.org/software/coreutils/manual/coreutils.html#head-invocation> (19.01.2016)
- 160 <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 42 (19.01.2016)
- 161 <http://www.duden.de/rechtschreibung/Homepage> (22.01.2016)
- 162 <https://docs.python.org/2/library/htmlparser.html> (19.01.2016)
- 163 <https://tools.ietf.org/html/rfc2818> (19.01.2016)
- 164 <http://www.duden.de/rechtschreibung/Hypertext> (22.01.2016)
- 165 <http://www.w3.org/standards/webdesign/htmlcss#whathtml> (19.01.2016)
- 166 <https://tools.ietf.org/html/rfc2616>, Abstract (19.01.2016)

- 167 <https://www.ietf.org/> (19.01.2016)
- 168 <https://tools.ietf.org/html/rfc3501> (19.01.2016)
- 169 <http://www.duden.de/rechtschreibung/Internetadresse> (22.01.2016)
- 170 <http://www.duden.de/rechtschreibung/Internetsurfen> (22.01.2016)
- 171 <http://www.ecma-international.org/publications/files/ECMA-ST/Ecma-262.pdf> (19.01.2016)  
und <https://developer.mozilla.org/en-US/docs/Web/JavaScript> (19.01.2016)
- 172 <http://www.oracle.com/technetwork/java/javaee/jsp/index.html> (19.01.2016)
- 173 <http://jpeg.org/> (23.01.2016)
- 174 <http://www.duden.de/rechtschreibung/Kommandozeile> (22.01.2016)
- 175 <http://www.duden.de/rechtschreibung/Kommandozeile> (22.01.2016)
- 176 u. a. <https://www.ibm.com/developerworks/community/blogs/nlp/entry/tokenization?lang=en>  
(22.01.2016)
- 177 <https://tools.ietf.org/html/rfc4511> (19.01.2016)
- 178 <http://www.duden.de/rechtschreibung/Link> (22.01.2016)
- 179 [http://www.duden.de/rechtschreibung/Log\\_in](http://www.duden.de/rechtschreibung/Log_in) (22.01.2016)
- 180 [http://www.duden.de/rechtschreibung/Log\\_in](http://www.duden.de/rechtschreibung/Log_in) (22.01.2016)
- 181 <http://lynx.invisible-island.net/> (19.01.2016)
- 182 <https://wiki.selfhtml.org/wiki/HTML/Kopfdaten/meta> (22.01.2016)
- 183 <https://www.microsoft.com/de-at/windows/microsoft-edge> (19.01.2016)
- 184 <https://products.office.com/de-at/excel> (19.01.2016)
- 185 <http://windows.microsoft.com/de-de/windows/windows-help#windows=windows-xp>  
(21.01.2016)
- 186 <http://www.w3.org/2001/tag/2002/0129-mime> (19.01.2016)
- 187 <https://www.mozilla.org/de/firefox/new/> (19.01.2016)
- 188 <https://tools.ietf.org/html/rfc5905> (19.01.2016)
- 189 [http://www.duden.de/rechtschreibung/Open\\_Source\\_Software](http://www.duden.de/rechtschreibung/Open_Source_Software) (22.01.2016)
- 190 <http://www.opera.com/de> (19.01.2016)
- 191 u. a. <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2005/n1905.pdf>, Seite 16, Kap.  
1.3.1 argument (22.01.2016) und <http://www.w3.org/standards/webdesign/htmlcss#whathtml>  
(19.01.2016)
- 192 [https://www2.informatik.uni-erlangen.de/EN/teaching/SS2012/HalloWelt/PA\\_2012.pdf](https://www2.informatik.uni-erlangen.de/EN/teaching/SS2012/HalloWelt/PA_2012.pdf), Seite  
3 (22.01.2016)
- 193 <http://php.net/> (19.01.2016)
- 194 u. a. <http://www.linfo.org/pipe.html> (22.01.2016)
- 195 <https://tools.ietf.org/html/rfc2046#section-4.1.3> (19.01.2016)
- 196 <http://tools.ietf.org/html/rfc2616#section-9.5> (19.01.2016)
- 197 <https://tools.ietf.org/html/rfc1939> (19.01.2016)
- 198 <http://www.adobe.com/devnet/postscript.html> (23.01.2016)

- 199 <http://www.duden.de/rechtschreibung/Protokoll#Bedeutung3a> (22.01.2016)
- 200 <https://www.python.org/> (21.01.2016)
- 201 <http://www.duden.de/rechtschreibung/Quellcode> (22.01.2016)
- 202 <http://www.duden.de/rechtschreibung/Netzwerk#b2-Bedeutung-3> (22.01.2016)
- 203 u. a. [http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd\\_chap09.html](http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd_chap09.html)  
(19.01.2016)
- 204 <https://www.gnu.org/software/wget/manual/wget.html#Recursive-Download> (22.01.2016)
- 205 <http://www.w3.org/Protocols/rfc2616/rfc2616-sec5.html> (19.01.2016)
- 206 <http://www.rfc-editor.org/> (22.01.2016)
- 207 <http://www.w3.org/Protocols/rfc2616/rfc2616-sec5.html> (19.01.2016)
- 208 <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)
- 209 <http://www.w3.org/Protocols/rfc2616/rfc2616-sec6.html> (19.01.2016)
- 210 <http://www.apple.com/at/safari/> (19.01.2016)
- 211 <https://www.gnu.org/software/sed/> (19.01.2016)
- 212 <http://docs.seleniumhq.org/> (19.01.2016)
- 213 [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp) (19.01.2016)
- 214 [http://docs.seleniumhq.org/docs/02\\_selenium\\_ide.jsp](http://docs.seleniumhq.org/docs/02_selenium_ide.jsp) (19.01.2016)
- 215 <http://tools.ietf.org/html/rfc2616#section-1.3> (19.01.2016)
- 216 <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 40 (19.01.2016)
- 217 <https://tools.ietf.org/html/rfc5321> (19.01.2016)
- 218 u. a. <http://www.duden.de/rechtschreibung/Skript#Bedeutung4> (22.01.2016)
- 219 <http://www.duden.de/rechtschreibung/Schraegstrich#b2-Bedeutung-b> (22.01.2016)
- 220 <https://www.gnu.org/software/coreutils/manual/coreutils.html#sort-invocation> (19.01.2016)
- 221 [http://www.linfo.org/standard\\_output.html](http://www.linfo.org/standard_output.html) (29.01.2016)
- 222 <http://www.duden.de/rechtschreibung/Syntax> (22.01.2016)
- 223 u. a. <http://www-01.ibm.com/software/globalization/topics/keyboards/physical.html>
- 224 [http://www.duden.de/rechtschreibung/Tag\\_Strukturelement\\_Markierung](http://www.duden.de/rechtschreibung/Tag_Strukturelement_Markierung) (22.01.2016)
- 225 <https://www.gnu.org/software/coreutils/manual/coreutils.html#tail-invocation> (19.01.2016)
- 226 <http://www.duden.de/rechtschreibung/Token> (22.01.2016) und  
<https://www.ibm.com/developerworks/community/blogs/nlp/entry/tokenization?lang=en>  
(22.01.2016)
- 227 <http://www.duden.de/rechtschreibung/Tool> (22.01.2016)
- 228 <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)
- 229 <https://www.gnu.org/software/coreutils/manual/coreutils.html#tr-invocation> (19.01.2016)
- 230 <http://releases.ubuntu.com/14.04/> (21.01.2016)
- 231 <https://tools.ietf.org/html/rfc1738> (19.01.2016)
- 232 <http://www.iana.org/assignments/media-types/media-types.xhtml> (19.01.2016)

- 233 <http://www.uni-weimar.de/medien/webis/teaching/lecturenotes/web-technology/unit-de-network-protocol3.pdf>, Seite 42 (19.01.2016)
- 234 <https://tools.ietf.org/html/rfc3629> (19.01.2016)
- 235 u. a. <http://www.duden.de/rechtschreibung/Ordner#b2-Bedeutung-3> (22.01.2016)
- 236 <http://www.duden.de/rechtschreibung/Seite#b2-Bedeutung-11> (22.01.2016)
- 237 <http://www.duden.de/rechtschreibung/Website> (22.01.2016)
- 238 <https://www.gnu.org/software/wget/> (19.01.2016)
- 239 <http://windows.microsoft.com/de-at/internet-explorer/download-ie> (19.01.2016)
- 240 <http://www.w3.org/> (19.01.2016)
- 241 <http://www.duden.de/rechtschreibung/Zeichenkette> (22.01.2016)
- 242 <https://tools.ietf.org/html/rfc2616>, Abstract (19.01.2016)

## D Abbildungsverzeichnis

Abbildung 1: Eine Website besteht in der Regel aus mehreren Webseiten .....	11
Abbildung 2: Ein konkretes Beispiel für eine URL .....	12
Abbildung 3: Kommunikation zwischen Client und Server gemäß HTTP (1) .....	14
Abbildung 4: Kommunikation zwischen Client und Server gemäß HTTP (2) .....	14
Abbildung 5: Ein konkretes Beispiel für eine Request-Zeile .....	15
Abbildung 6: Ein konkretes Beispiel für eine Statuszeile .....	15
Abbildung 7: Öffnende und schließende Tags in HTML .....	18
Abbildung 8: Gesamtprozesses der automatisierten Datenextraktion aus Webseiten .....	20
Abbildung 9: Einteilung der Arten von Webseiten .....	22
Abbildung 10: Statische Webseiten.....	22
Abbildung 11: Serverseitig generierte dynamische Webseiten.....	24
Abbildung 12: Clientseitig generierte dynamische Webseiten .....	25
Abbildung 13: Hybride dynamische Webseiten.....	26
Abbildung 14: Benutzeroberfläche von Selenium IDE.....	33
Abbildung 15: Tabellenzeile einer Lehrperson .....	47
Abbildung 16: Diese Eingaben werden von Selenium automatisiert.....	50
Abbildung 17: Modulliste .....	52

## E Tabellenverzeichnis

Tabelle 1: Auszug aus einer Übersicht der erlaubten MIME-Types .....	17
Tabelle 2: HTML und CSS haben unterschiedliche Aufgaben.....	19
Tabelle 3: Optionen von <b>wget</b> .....	30
Tabelle 4: Optionen von <b>curl</b> .....	32
Tabelle 5: Befehle von Selenium (angewendet in Selenium IDE).....	34
Tabelle 6: Optionen von <b>cat</b> .....	35
Tabelle 7: Eine wichtige Option von <b>cut</b> .....	36
Tabelle 8: Optionen von <b>grep</b> .....	37
Tabelle 9: Optionen von <b>sort</b> .....	41
Tabelle 10: Befehle zur Extraktion der Modulliste im HTML-Format.....	50

## F Quelltextverzeichnis

Quelltext 1: Aufbau eines einfachen HTML-Dokuments.....	18
Quelltext 2: Ausgabe von <b>cat</b> .....	36
Quelltext 3: Ausgabe von <b>cut</b> .....	37
Quelltext 4: Ausgabe von <b>grep</b> .....	38
Quelltext 5: Ausgabe von <b>head</b> .....	38
Quelltext 6: Ausgabe von <b>sed</b> .....	40
Quelltext 7: Ausgabe von <b>sort</b> .....	41
Quelltext 8: Ausgabe von <b>tr</b> .....	42
Quelltext 9: Aufbau der Tabellenzeile einer Lehrperson.....	47
Quelltext 10: Befehlskette zur Extraktion der Sprechstundenliste.....	47
Quelltext 11: Tabellenzeile einer Lehrperson ( <b>sprechstundenliste.csv</b> ).....	49
Quelltext 12: Befehle zur Speicherung des Quelltexts der Modulliste.....	51
Quelltext 13: Aufbau der Tabellenzeile eines Moduls.....	52
Quelltext 14: Befehlskette zur Extraktion der Modulliste.....	53
Quelltext 15: Tabellenzeile eines Moduls ( <b>modulliste.csv</b> ).....	55
Quelltext 16: Befehlskette zur Extraktion einer Liste der Lehrpersonen ( <b>lynx</b> ).....	59
Quelltext 17: Python-Skript zur Speicherung der Modulliste.....	60



## G Beiliegendes Speichermedium

Folgende Dateien sind auf dem beiliegenden Speichermedium vorhanden:

- `sprechstundenliste.html`
- `modulliste.html`
- `lynxdump.txt`