

CENTRUM FÜR INFORMATIONS- UND SPRACHVERARBEITUNG STUDIENGANG COMPUTERLINGUISTIK



## Bachelorarbeit

im Studiengang Computerlinguistik an der Ludwig- Maximilians- Universität München Fakultät für Sprach- und Literaturwissenschaften

# Entwicklung eines Rankingverfahrens der Suchtreffer für die FinderApp WiTTFind im Nachlass Ludwig Wittgensteins

vorgelegt von Florian Babl

Betreuer: Dr. Maximilian Hadersbeck Prüfer: Dr. Maximilian Hadersbeck Bearbeitungszeitraum: 1. April - 11. Juni 2019

### Selbstständigkeitserklärung

Ich erkläre hiermit, dass ich die vorliegende Arbeit selb-

ständig angefertigt, alle Zitate als solche kenntlich gemacht sowie alle benutzten Quellen und Hilfsmittel angegeben ha-
be.
München, den 11. Juni 2019
Florian Babl

### **Abstract**

Ziel der vorliegenden Bachelorarbeit ist es, für die FinderApp WiTTFind ein Rankingverfahren zu entwickeln, das dazu dienen soll, Mehrwortanfragen zu verarbeiten und die Ergebnisse nach ihrer Relevanz zu sortieren. Dabei wurde die schon existierende Suche der WiTTFind-Suche so optimiert, dass zu jedem Satz und damit auch Absatz, ein Score berechnet wird. Der Score dient als eine Aussage darüber, wie relevant der Suchtreffer für den Suchenden ist. Zusätzlich wurde im Rahmen dieser Arbeit den Nutzern der Suche ermöglicht relevante Ergebnisse zu einer Suchanfrage zu liefern. Die Ergebnisse müssen dabei nicht wortwörtlich und hintereinander im Wittgenstein Nachlass vorkommen. Für das Rankingverfahren wird eine eigens erstellte Formel für die Berechnung der Scores vorgestellt. Sie beinhaltet mehrere selbst gewählte und implementiere Features, wie Distanz zwischen Matches, Lemmata und Matchgenauigkeit. Ein wesentlicher Bestandteil des Rankings sind lineare Belohnungen und Bestrafungen. Am Ende der Arbeit werden die Ergebnisse vorgestellt, die durch die neue Suche erreicht werden. Teil der Arbeit ist es auch, einen unkomplizierten Einstieg die Programmierung des wf zu ermöglichen. Zusätzlich wird kurz die Entstehung und Funktionalität der neuen Wittgenstein Advanced Search Tools Dokumentation erklärt. Die Arbeit ist für jeden interessant, der sich einen Einblick in die Erstellung eines Rankingverfahrens verschaffen möchte und ein besseres Verständnis der WiTTFind Suche erstrebt.

# Inhaltsverzeichnis

Αŀ	strac	t	ı
1	Einle	eitung	3
2	2.1 2.2 2.3 2.4	Motivation	5 5 5 6
3	3.1	r J J	7 7 9 9
	3.3 3.4 3.5 3.6 3.7 3.8	Message.hpp/.cpp	10 11 12 12 13 13
4		U .	15
	4.1		15
	4.2	0	15
	4.3	4.3.1 Erste Überlegungen und Probleme	15 15 16
	4.4	Die Features für das Ranking	16 16 17 18 21
	4 5		$\frac{22}{23}$
	$4.5 \\ 4.6$	0	23 24
	4.7		$\frac{24}{25}$
5	Erge	bnisse des Rankingverfahrens	27
	5.1	1	27
			27
			28
			29
			30
			30
	5.2	Fazit	31

6	Aus	olick	33
	6.1	Verbesserungsmöglichkeiten	33
		6.1.1 Einzelne Matches	
		6.1.2 Frequenzlisten	
		6.1.3 Eingliederung in die Website	
		6.1.4 Index aufbauen für Suche	
		6.1.5 Rankingstatistik	
	6.2	Schluss	
7	Anh	ang	37
		Code	37
		7.1.1 Searcher.cpp	
		7.1.2 reply.json	
		7.1.3 RankedSearch.cpp/.hpp	
	7.2	Suchergebnisse	
	7.3	READMEs aus Gitlab	
		7.3.1 Sphinx	
		7.3.2 wf	
Lit	teratı	rverzeichnis	55
Αŀ	bildu	ngsverzeichnis	57
Inl	halt d	es beigelegten USB-Sticks	50

# 1 Einleitung

Als Nutzer einer Websuche ist man es durch Google gewöhnt, einfach mehrere Begriffe hintereinander zu nennen und als Folge relevante Suchergebnisse geliefert werden. Das Ergebnis ist dabei immer stichhaltig und im Bezug zur Frage. Dadurch verändert sich in der Konsequenz die Einstellung in Bezug auf andere Suchmaschinen. Die WiTTFind Suchmaschine durchsucht einen Korpus und soll nicht die Funktion einer Suchanfrage wie Google erkennen können. Das ist auch gar nicht ihr Zweck. Denn sie soll Textstellen aus dem Nachlass von Ludwig Wittgenstein finden, die relevant für die Suchanfrage sind. Bei einer Suchanfrage für zwei oder mehr Begriffe, möchte der Nutzer Textstellen, in denen möglichst alle Begriffe vorkommen. Jedoch ist die Suche momentan limitiert darauf Begriffe zu finden, die direkt aufeinanderfolgen. Das Problem lässt sich am besten anhand eines Beispiels erklären. In der WiTTFind-Suchanfrage "Gedanken Welt" werden keine Ergebnisse zurückgeliefert, die für den User möglichst relevant sein können. Da die Wortkombination "Gedanken Welt" im Korpus an keiner Stelle direkt hintereinander auftritt, würde die aktuelle Suchmethode keine relevanten Ergebnisse ausgeben. Dennoch gibt es Sätze in Wittgensteins Nachlass, in denen "Gedanken" und "Welt" in einem Satz vorkommen, jedoch nicht unmittelbar. Die Problemstellung ist es dementsprechend eine Methode zu entwickeln, mit der diese direkte Reihenfolge nicht mehr von Nöten ist, um Treffer zu finden. Nichtsdestotrotz soll dabei gewährleistet sein, dass Begriffe, die aufeinander folgen als Erstes in den Ergebnissen auftauchen. Für so ein Ranking muss also jeder Treffer mit einem Score ausgestattet werden, der beschreibt, wie relevant ein Treffer ist.

Ein Patent aus dem Jahr 2003 von Google Developer Krishna Bharat liefert eine Möglichkeit, einen solchen Score zu berechnen. Die Umstände und Ziele von Google sind nicht mit dieser Arbeit zu vergleichen, jedoch wird klar, dass auch im Patent von einem Ranking nach Relevanz die Rede ist. Es wird außerdem jedem Hit, d. h. jeder Website, ein Rank zugeordnet und nach dem höchsten Rank sortiert [4][5 ff.]. Für diese Arbeit soll der Fokus darauf liegen, eine Formel mit Features zu entwickeln, die für die Zwecke von WiTTFind relevant sind. Die Ranks sollen dann den Hits zugewiesen werden bevor sie anschließend sortiert werden.

# 2 Erstellung der WAST Dokumentation mit Sphinx

#### 2.1 Motivation

Das Wittgenstein Advanced Search Tools (WAST), Projekt läuft seit mehreren Jahren unter der Leitung von Dr. Max Hadersbeck und Alois Pichler. Obwohl das Projekt in der Vergangenheit bereits im Rahmen diverser Bachelorarbeiten um einige Features erweitert und verfeinert wurde, konnte die vollständige Dokumentation des Projekts bisher nocht nicht abgeschlossen werden. Wie den letzten Aufzeichnungen von Dr. Maximilian Hadersbeck und Daniel Bruder zu entnehmen ist, war die letzte Aktualisierung auf den neuesten Stand im Januar 2014 [3, 1]. Eine umfassende Dokumentation ist jedoch essenziell für eine schnelle und erfolgreiche Einarbeitung in das gesamte bzw. eines der mittlerweile zahlreichen Teilprojekte.

Für jede der angesprochenen Bachelorarbeiten und für jeden Bestandteil von WAST wird im entsprechenden Gitlab Ordner ein neuer Branch bzw. ein neues Projekt erstellt. Hierbei ist der Ablauf fest vorgeschireben und bewirkt, dass für jedes Teilprojekt eine READ-ME.md Datei erstellt wird. Der Aufbau von Gitlab sorgt jedoch dafür, dass die Suche nach speziellen Begriffen keine Ergebnisse liefert. Deshalb ist die Erstellung einer eigenen Dokumentation, die alle Projekte umspannt, vonnöten. Von zentraler Bedeutung ist dabei, dass die entstehende Übersicht online auf ansprechende Art und Weise eingebunden werden kann. Darüber hinaus soll die Dokumentation auch einfach zu durchsuchen sein. Wie ein Vorlesungsskript sollte man nach bestimmten Begriffen suchen können und alle Ergebnisse angezeigt bekommen. Die Bachelor und Master Studenten finden dadurch Zeit sparender und effektiver einen Zugang zum WAST Projekt.

### 2.2 Sphinx Tool

Zur Erstellung der Dokumentation wurde das Sphinx Tool ausgewählt. Sphinx ist ein Werkzeug zur Softwaredokumentation, das erstmals 2008 aufgetreten ist. Sphinx nutzt zum Generieren einer Website oder einer PDF Datei reStructuredText oder Markdown Dateien. Dabei ist Sphinx vor allem für die gesamte Dokumentation des Python-Projekts bekannt geworden, welche komplett mit Sphinx realisiert wurde [11]. Ferner wurde das Programm zur Dokumentation anderer Python Module genutzt und sogar zum Schreiben von Büchern genutzt [9]. Sphinx erweist sich dank seiner weitreichenden Möglichkeiten als noch hilfreicher als Latex, um eine Dokumentation zu erstellen und ermöglicht es, die bestehenden Teil-Dokumentationen des WAST Projekts zusammenzufassen. Alle in 2.1 erwähnten Ziele können folglich mit Sphinx erreicht werden.

#### 2.3 Funktionalität

Mithilfe von Sphinx und der Hilfe der READMEs aus wf, sis3, witt-data uvm. kann eine Website erstellt werden, die es möglich macht, alle READMEs auf einmal zu durchsuchen und damit die wichtigsten Treffer anzuzeigen. Die einzelnen README Dateien können im Verzeichnis "WAST-DOC/source" gefunden werden. An den nachfolgenden Ablauf ist sich zu halten, sobald der Dokumentation neue Markdown Dateien hinzugefügt oder die Ansicht der Website abgeändert werden soll:

Zuerst muss wird der Branch "CAST/infrastructure/wast/tree/WAST\_Dokumentation/WAST-DOKU" aus Gitlab geclont, Sphinx, recommonmark und Sphinx-Markdown-Tables mit "pip" oder "pip3" installiert. Danach wird in den entsprechenden Ordner gewechselt und make html aufgerufen. Nachdem nun alle Markdown Dokumente verarbeitet wurden, kann jetzt die File "index.html" aud dem "build" Ordner mit Firefox geöffnet werden (Seite 52). Die Datei config.py musste um das Modul CommonMarkParser und sphinx\_markdown\_tables erweitert werden, damit Markdown Dateien auch als solche geparst werden können und Tabellen erkannt werden. Für eine verbesserte Ansicht der Website sind zusätzlich die Überschriften aus den Markdown Dateien angepasst worden, denn der CommonMarktParser macht aus allen Level eins Überschriften einen Unterpunkt.

```
1
   from recommonmark.parser import CommonMarkParser
2
   source_parsers = {
3
4
       '. md': CommonMarkParser,
5
6
   source_suffix = ['.rst', '.md']
7
8
9
10
   extensions = [
11
        'sphinx_markdown_tables',
12
13
```

### 2.4 Ergebnis

Zusätzlich zu den bereits genannten Features, bietet Sphinx außerdem die Möglichkeit, das Design der Website frei zu wählen. Die fertige Internetseite ist im nachfolgenden Bild zu sehen. Neben der Möglichkeit, alle links aufgelisteten Teilprojekte zu durchsuchen, existieren auch noch Unterpunkte für die einzelnen Stichpunkte innerhalb des Projekts. Für die grafische Oberfläche sorgt "sphinx\_rtd\_theme". Ein anderes Design-Template kann in der Datei config.py eingestellt werden.

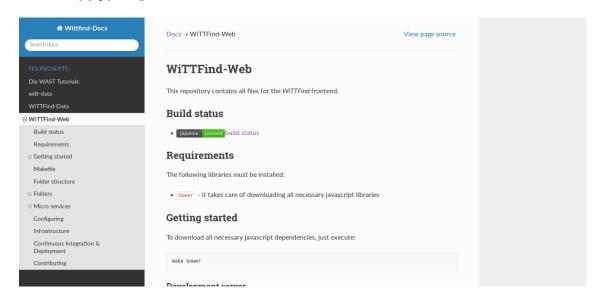


Abbildung 2.1: Ansicht der neuen WAST Dokumentationswebsite

In den nächsten Schritten muss die Dokumentation nun online angeboten werden. Des Weiteren wäre es sinnvoll, ein Skript zu schreiben, das einmal im Monat ausgewählte README Dateien aus dem WAST-Gitlab liest, aktualisiert und dadurch die Dokumentation stets aktuell hält.

### 3 wf: Aktueller Stand

Derzeit stehen bereits viele Suchmöglichkeiten in wf zur Verfügung. So kann zum Beispiel nach einem oder mehreren Begriffen gesucht werden. Die Suche ermöglicht es dem Nutzer, mit regulären Ausdrücken zu arbeiten oder auch nach Wortart-Tags zu suchen. Auch existiert eine Methode, die Partikelverben finden kann. Des Weiteren gibt es unterschiedliche Arten der Suche. Beispielsweise wird genau der String  $\alpha$  gefunden, wenn dieser in Hochkommata gesetzt wird. Der Code ist zwar außerordentlich strukturiert und leserlich, jedoch fehlt es an Kommentaren oder gar einer gesamten Dokumentation. Weshalb ein guter Einstieg in den Code des Programms zu finden, erschwert ist. Im Folgenden sollen die wichtigsten Funktionen von wf erläutert werden. Dies soll einen schnellen Einstieg in das Thema ermöglichen und ein paar Begriffe und Termini erläutern, die in den anschließenden Kapiteln häufiger benutzt werden. Die genaue Erläuterung der Funktionsweise der Suche ist dabei jedoch nicht Teil dieser Arbeit. Vielmehr wird die Funktionalität von wf erklärt, damit wichtige Schnittpunkte in Zukunft schneller erkannt werden.

#### 3.1 wf: Stuktur

#### 3.1.1 Aufbau

In der unteren Abbildung soll der allgemeine Aufbau von wf deutlich werden. Der *wfa* Server übergibt die Suchanfrage an wf, welcher dann die Suche durchführt. Die Suche findet in *Search.hpp* statt und wird dort sehr komplex. Für die vorliegende Arbeit sind hauptsächlich die Dateien von 3.2 bis 3.7 von Interesse.

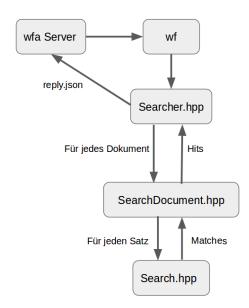


Abbildung 3.1: wf Schaubild

Wf ist eine in C++ programmierte Suchmaschine, die vom wfa Server aufgerufen wird und ein Json Objekt zurückliefert, das Informationen zu den gefundenen Treffern bzgl. der Suchanfrage beinhaltet (Siehe 3.2). Jedoch kann wf auch direkt von der Konsole gestartet werden. Dazu wird der  $wf\_client$  und  $wf\_server$  genutzt.

Sobald wf installiert ist, kann die Suche ausgeführt werden. Allerdings kann wf\_client kann nur aufgerufen werden, wenn der zuständige wf\_server bereits läuft. Auf Grund dessen wird dieser zuerst gestartet:

In files.txt stehen die Namen der Dateien, die nach Treffern durchsucht werden. Im Anschluss wird der  $wf\_client$  ausgeführt:

Wenn mehrere Dateien durchsucht werden sollen, können mehrere -f Ausdrücke angehängt werden. Alle anderen Optionen zur Anpassung der Suche können im Anhang auf Seite 54 gefunden werden. Sobald eine neue Stelle im Code ergänzt wurde, muss im build Ordner der Befehl make aufgerufen, sowie der Server neu gestartet werden, damit die Änderungen übernommen werden.

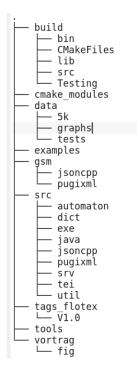


Abbildung 3.2: Verzeichnisaufbau von wf

Der Aufbau des wf Ordners ist für ein besseres Verständnis nicht zu vernachlässigen. Abbildung 3.2 zeigt die komplette Übersicht über das Verzeichnis.

Der build-Ordner entsteht erst, nachdem er mit make erstellt wurde. Aus diesem Verzeichnis heraus wird der  $wf\_client$  und  $wf\_server$  gestartet. Hier finden sich außerdem mehrere Makefiles.

Ein weiterer wichtiger Ordner ist *data*. Hier finden sich alle Text, XML, json und grf Files, die zum Testen von wf benötigt werden.

In examples sind sowohl die reply. json als auch die request json zu finden.

gsm enthält die Module pugixml und jsoncpp, das zum Bearbeiten der Json Objekte unabdingbar ist.

Der gesamte Code lässt sich in den Unterverzeichnissen von  $src^1$  finden. Dabei sind beson-

<sup>&</sup>lt;sup>1</sup>Steht für Source.

ders sind folgende Unterordner relevant:

- dict: Beinhaltet den Code für den Umgang mit dem Dictionary, dem Parser und den grammatikalischen Angaben.
- srv: Besteht aus dem Code für die Initialisierung der tatsächlichen Suche durch den Automaten.
- automaton: Der Code des Automatons<sup>2</sup>, des Handlings von Token und der Suche von Ergebnissen kann hier unter anderem recherchiert werden.
- util: Hier lassen sich einige utilitätsbezogene C++ Files finden.

#### 3.1.2 Begrifflichkeiten

Um für das Verständnis ein gewisses Hintergrundwissen zu erhalten, werden in den nachfolgenden Abschnitten immer wiederkehrende Begriffe kurz näher definiert und erläutert.

- Unter Hits bzw. einem Hit wird der Treffer auf Satzebene verstanden. In Abschnitt 3.2 ist der Aufbau eines Hit-Arrays veranschaulicht.
- Für jeden Hit gibt es einen oder mehrere Matches. Jedes Match besteht aus dem Token selber und liefert einige wichtige zusätzliche Informationen für diesen Treffer.
- Ein Query Token ist ein Token, das in die Suchanfrage eingegeben wird. Query-String oder Eingabe-String steht stellvertretend für die gesamte Suchanfrage.
- Wenn von einem Rank die Rede ist, so ist das gleichbedeutend mit einem Score oder einem Rang. In dieser Arbeit wird mit dem Begriff "Rank" der Integerwert beschrieben, der für jeden Hit berechnet werden soll.
- In jeder der nachfolgenden und einigen hier nicht erwähnten Dateien, wurden im Code Kommentare hinzugefügt, um den Code verständlicher zu machen.
- Der neu eingeführte Key s rank hält später den berechneten Rank vom Satz/Hit.
- Genauso enthält ab rank den Rank vom Absatz in dem der Hit steht.
- Kurz erklärt definiert ein sog. Automaton, oder auch Automat, auf natürliche Art und Weise eine Sprache und einen zugehörigen Erkennungsalgorithmus. Der Automat bekommt ein zufälliges Wort und soll bestimmen, ob es sich um ein Wort der Sprache handelt [8][17].

### 3.2 Reply.json und Request.json

Request.json stellt das Json Objekt dar, das vom wfa Server an die wf Suche weitergeleitet wird. Reply.json repräsentiert eines der wichtigsten Objekte aus wf. Es liefert die Suchergebnisse im Json Dateiformat zurück. Eine Beispieldatei ist auf Seite 39 zu finden. Jedes Key Value Paar steht für ein wichtiges Attribut, das die Suchergebnisse genauer definiert. Die für die Bachelorarbeit relevanten Keys werden nachfolgend erläutert. Die restlichen Key Value Paare sind im Anhang auf Seite 39 zu finden.

• "buildertype": Gibt an, wie der Automat später zusammengesetzt werden soll. Simple-Builder baut den Automaten so, dass zu einer Suchanfrage nur nach den Lemmata des Query Tokens gesucht wird. Hingegen findet InverseBuilder aber auch inverse Lemmata, d.h. wenn nach einem schon flektierten Verb gesucht wird, können jetzt das Lemma und andere Flexionen gefunden werden.

<sup>&</sup>lt;sup>2</sup>Gleichbedeutend mit Automat

- "hits": Ist ein Array, das aus mehreren einzelnen Hits besteht. Jeder Hit hat folgende Keys:
  - "abnr": Gibt die Absatznummer im Dokument an.
  - "matches": Ist ein Array, das mehrere einzelne Matches enthalten kann. Jedes Match hat eigene Anmerkungen. Jedes Token hat eine spezifische ID, die in "id" übergeben wird. Der Tokenname wird in "token" gespeichert und die Position im Satz wird in "pos" abgelegt. In "gc" ist der GrammaticalCode abgelegt. Das sind grammatische Kennzeichen, die beschreiben, in welchen Flexionen und anderen Abwandlungen das Wort vorkommt.
  - "n": Jeder Satz hat ein eigenes "n". Für das Beispiel "Ts-213,iii-r[8]\_1" lässt sich sagen, dass der Hit auf der Seite "Ts-213,iii-r" im ersten Satz des achten Absatzes vorkam.
  - "satznr": Für jedes Dokument sind alle Sätze durchnummeriert.
- "query": Der Value beinhaltet den Suchstring.

### 3.3 Searcher.hpp/.cpp

Die zwei wichtigsten Funktionen aus Searcher.hpp sind search(...) und buildAnswer(...), zu finden auf Seite 37:

• In search(...) werden Builder, Automat und Skipper initialisiert (vgl. Zeile 3, 5,6) und für jeden Pfad zu einem Dokument die SearchDocument() Klasse aufgerufen (Zeile 4). Die Ergebnisse werden dann alle in der Variable futures gespeichert und später in das SearchDocument::Message Objekt res übertragen (Zeile 9 - 14). Wie man dem Code in Zeile 1 und 11 entnehmen kann, gibt es ein request und ein res Objekt. Beide Objekte sind Teil verschiedener Klassen. Denn mit dem Message Objekt request wird sowohl die Suchanfrage übergeben, als auch das finale Json Objekt zurückgeliefert. Während mit dem SearchDocument::Message Objekt, wie in 3.5 genauer erklärt, nur die Ergebnisse abgespeichert werden. Wenn es sich nicht um eine gerankte Suche handeln soll, wird in Zeile 19 deutlich, dass im Anschluss buildAnswer(...) aufgerufen wird.

```
void wf:: Searcher:: search (Message& request, bool ranked search)
    const
2
   auto builder = request. GetBuilder(dict); // Baut Suchautomat
3
      in Message.cpp
   request. For Each Path ([\&] (const std::string& path) {
4
     auto automaton = query->GetAutomaton(*builder);
5
         SimpleQuery.hpp gibt query token zu Automaton
     auto skipper = request.GetSkipper(); //holt den Skipper (bei
6
         Satzzeichen usw.)
                                           // search bekommt die
     auto search = SearchDocument(
       std::move(automaton), std::move(skipper), files_->Get(path)
       futures.push back(std::async(std::launch::async, std::move(
           search)));
```

<sup>&</sup>lt;sup>3</sup>Steht für Results.

```
});
10
    SearchDocument::Message res; // vektor aus Hit Objekten
11
    for (auto& future : futures) { // for each Hit
12
      auto tmp = future.get(); //tmp = hits
13
      std::copy(begin(tmp), end(tmp), std::back_inserter(res));
14
1.5
       (ranked search)
16
17
18
       else buildAnswer(request, res, max);
19
20
```

• buildAnswer(...) baut aus dem Message Objekt mit dem Namen request das endgültige Json Objekt, das am Ende von wf zurückgeliefert wird. Es werden für jeden Hit mehrere Key Value Paare hinzugefügt, die den Hit genauer beschreiben sollen, (Zeile 9). Abschließend werden die gefundenen Matches noch eingefügt (Zeile 17) und die Anzahl der Hits n festgelegt (Zeile: 20).

```
void wf:: Searcher:: buildAnswer (Message& request,
                      const SearchDocument:: Message& res, size t max)
                           const {
3
      std::unordered\_map{<}std::string\ ,\ size\_t{}>\ counts\ ;
      for (const auto& r : res) { // for each Hit
        n += r.matches.size(); // n = Anzahl der Matches
         for (const auto& m : r.matches) { // alle Matches iterieren
           auto& hit = request.AppendHit()
                           . SetPath (r.path)
10
                           . SetN(r.n)
11
                           . SetF (r.f)
                           . SetAbnr(std::stoul(r.abnr))
13
                           . SetSatznr(std::stoul(r.satznr))
14
                           . Set Counter (i++);
15
16
           for (auto i = m. first; i != m. second; ++i)
17
             hit.push back(*i);
19
      request . SetNumberOfHits(n);
20
21
```

### 3.4 Message.hpp/.cpp

Message.hpp beinhaltet die drei Klassen Message, Hit und Match. Jede hat mehrere Setter/ und Getter-Funktionen, die das jeweilige Objekt bearbeiten. Die bereits in 3.2 erklärten Attribute, werden von diesen Funktionen gesetzt. Beim Erstellen eines Objekts wird immer ein Json::Value erstellt. Daraus entwickelt sich am Ende ein großes Json Objekt der Klasse Message, das auf alle anderen Json Values verlinkt. Das Objekt bildet dann die vollständige Antwort<sup>4</sup> an den wfa Server.

 $<sup>^4</sup>$ reply.json

### 3.5 SearchDocument.hpp/.cpp

In SearchDocument() wird über jeden Satz der File iteriert und für jeden Satz die Automaten Suche aufgerufen. Die gefundenen Treffer werden in matches abgespeichert. Wenn von matches die Rede ist, heißt das also nichts anderes als "Treffer in einem Satz". Dementsprechend wird für jeden Satz mit Treffern eine Hit-Struct erstellt. Die struct-Variablen werden nun von mehreren Values gefüllt. Alle Hits werden aneinander angehängt und woraus ein Vektor voller Hits entsetht, über den in buildAnswer(...) iteriert werden kann.

```
1 wf::SearchDocument::Message
2 wf::SearchDocument::operator()() const {
    for (const auto& sent : file_ .sentences) { //iterate sentences
      matches.clear();
5
      search(*automaton_, //calls automaton/Search.hpp
6
             begin (sent.seq), // seq = Vektor aus Token
             end(sent.seq), // Iter. on begin and end of sentence.
             *skip ,
             std::back inserter(matches)); // append "matches"
10
      hits.emplace back();
11
      hits.back().path = path;
12
      hits.back().n = sent.n; //siglum
13
      hits.back().f = sent.f; //facsimile
14
      hits.back().abnr = sent.abnr;
15
      hits.back().satznr = sent.satznr;
16
      hits.back().matches = matches;
17
      n += matches.size();
18
19
```

### 3.6 Dictionary.hpp/.cpp

Jeder Eintrag im Dictionary ist zu einem struct Entry verlinkt. Jeder Entry kann mit der Funktion find(...) aus der Klasse Dictionary gefunden werden. find(...) kann entweder mit der id des Entrys oder mit dem std::wstring aufgerufen werden. Das wird in 4.4.3 noch relevant sein. Für zwei Entries kann mit der Funktion isLemmaOf(...) und isInverseLemmaOf(...) überprüft werden, in welcher Relation die Entries zueinander stehen.

```
struct Entry {
                   // Jedes Wort hat einen Eintrag im Dictionary mit
      mehreren Attributen
    Entry(std::wstring s, size t i)
      : id(i)
3
        str(std::move(s))
        gc ()
        ic ()
6
        lemmata()
    bool isLemmaOf(const Entry* e) const noexcept;
9
    bool isInverseLemmaOf(const Entry*) const noexcept;
10
11
    size_t id;
12
13
    std::wstring str;
```

```
GrammaticalCode gc;
InflectionalCode ic;
std::set<const Entry*> lemmata;
};
```

### 3.7 Search.hpp

An dieser Stelle wird der Code sehr detailiert und es ist nicht Aufgabe dieser Bachelorarbeit intensiv auf den Code einzugehen. Kurz zusammengefasst gibt es zwei Funktionen: Die search(...) Funktion ruft die find(...)-Funktion (Seite 37) auf. In find(...) wird ein Stack aufgebaut, der abgearbeitet werden soll und die delta(...) Funktion übernimmt die Suche. Je nachdem, mit welchen Parametern der Automat gebaut wurde, es gibt unterschiedliche delta(...) Funktionen, die jetzt unterschiedlich suchen. Die neu gesetzten Kommentare im Code erklären die wichtigsten Stellen.

```
1 template < class It, class OutIt>
2 OutIt
3 wf::search(const Automaton&a, It b, It e, const Skip& skip,
     OutIt out) // out = matches,
4 {
          // boolean üfr ranked search
                                               //b = begin , e = end
    for (b = skip.skip(b, e); b != e; b = skip.skip(b, e)) {
5
      auto tmp = b;
6
      // Alle b abspeichern und ßanschlieend distanz messen
      if (find(a, b, e, skip)) //wenn ein match gefunden wurde.
        *out++= std::make\ pair(tmp,\ b);\ //\ tmp = iterator\ paar.
10
        ++b; // weitergehen
11
12
    return out;
13
14
```

### 3.8 wf: Problemstellung

Das Hauptproblem, das auf tieferer Ebene erkennbar wird, ist, dass die wf Suche so implementiert wurde, dass sie direkt aufeinander folgende Treffer identifizieren soll. Bei einer Sucheingabe, wie wir denken, wird nicht automatisch ein Hit gefunden, der einen oder mehrere Wildcards zwischen den Wörter zulässt, wie z.B. "wir sollen denken" oder kurz "wir \* denken". Nur mithilfe eines regulären Ausdrucks lässt sich die Suche entsprechend einstellen. Das führt jedoch zu anderen Problemen, die im nächsten Kapitel angesprochen werden.

# 4 Rankingverfahren

#### 4.1 Motivation

An wf wurden in den letzten Jahren keine größeren Änderungen vorgenommen. Es wurden immer wieder Features ergänzt, aber grundlegend wurde die Suche nicht mehr verändert. Während wf zwar makellos funktioniert, kann man es dennoch immer weiter verbessern. Ein Feature, welches es dem Nutzer ermöglicht, Treffer auf Satzebene zu finden, auch wenn die entsprechenden Wörter nicht direkt hintereinander stehen, wäre daher eine wertvolle Ergänzung für das Projekt. Suchende könnten leichter Abschnitte finden, die für sie relevant sind und und sie müssten nicht mehr auf schwierige reguläre Audrücke zurückgreifen, um die gewünschten Ergebnisse zu erhalten. WittFind wird außerdem einsteigerfreundlich, da die Suche dann eher der Googlesuche ähnelt. Da hauptsächlich Philosophen die Suche nutzen, kann es ihnen auch helfen, schneller zu Ergebnissen zu kommen. Das entwickelte Ranking basiert auf rationalen und logischen Annahmen für relevante Features.

### 4.2 Zielsetzung

Das Ziel dieser Arbeit ist es, dem Nutzer zu ermöglichen, Mehrwort-Anfragen in die WiTT-Find App einzugeben und dabei Ergebnisse anzuzeigen, wenn auch mehrere Worte zwischen den Suchbegriffen stehen. Bei Betrachtung des Beispiels "Welt Gedanken" liefert die WiTT-Find App nur einen Treffer zurück, bei dem die beiden Query Token direkt hintereinander vorkommen. Die Resultate sollten dann nach ihrer Relevanz sortiert werden. Die aufkommende Frage, was ein Ergebnis relevant macht und welche Features man für die Berechung verwenden kann, wird in 4.4 beantwortet. Als logische Konsequenz muss jeder Hit einen Rank zugeordnet bekommen, damit die Hits nach diesem Wert sortiert werden können. Bisher ist die Reihenfolge, in denen die Treffer ausgegeben werden, durch die Reihenfolge der Dokumente bestimmt.

Hinzu kommt, dass die Ranks einem Absatz zugewiesen werden sollen. Auf lange Sicht hin soll die Website der WiTTFind App ganze Absätze anzeigen und nicht nur einzelne Hits, deshalb reicht es nicht, die Scores nur satzweise zu berechnen.

### 4.3 Entwicklungsprozess

### 4.3.1 Erste Überlegungen und Probleme

Der erste Ansatz bestand darin, die komplette Suche von wf abzuändern und so anzupassen, dass Suchergebnisse nicht hintereinanderstehen müssen, damit sie gefunden werden. Obwohl wf einwandfrei funktioniert, ist es nicht trivial den richtigen Ansatzpunkt zu finden, um so fundamentale Änderungen durchzuführen. Wenn man so eine Veränderung durchführen will, muss gewährleistet sein, dass nicht die gesamte Suche zusammenbricht. Daher wurde der Ansatz schnell wieder verworfen, da sehr viele andere Komponenten an wf hängen, die davon abhängig sind, dass wf weiterhin die Ergebnisse liefert, die es bisher geschickt hat.

Eine weitere Überlegung bestand darin, eine neue, eigenständige Suche zu schreiben, die gesondert aufgerufen wird. Die Suche sollte die Treffer einzeln matchen. Hierbei wären neue Transitions und Skipper nötig gewesen, die Teil des Automaten sind. Auch dieser Ansatz wurde verworfen, da es an der Dokumentation des Codes und an der gesamten Verständlichkeit des Automaten scheiterte.

#### 4.3.2 Finaler Ansatz

Der finale Ansatz nimmt einen einfacheren Weg, der die anderen Komponenten von wf berücksichtigt und dennoch eigenständig aufgerufen werden kann. Mit Hilfe des logischen Operators "|" sollen die Query Token separiert werden. Mit dem "|"-Operator ist wf schon in der Lage, innerhalb von einem Satz Ergebnisse zu finden, die nicht nebeneinanderstehen müssen. Der Ausdruck  $\alpha|\beta^1$  liefert aber gleichzeitig Ergebnisse, die den linken Ausdruck und/oder den rechten Ausdruck beinhalten. [10][121]. D. h. mindestens einer der beiden Ausdrücke muss wahr sein, damit der ganze Ausdruck wahr wird. Das dadurch gefundene Json Objekt beinhaltet folglich Hits mit  $\alpha$ , mit  $\beta$  und Hits, die beide Ausdrücke gefunden haben. Wenn Letzteres der Fall ist, lässt sich daraus schließen, dass  $\alpha$  und  $\beta$  in einem Satz zusammen aufgetreten sind. Letztlich bedeutet dies für das Ranking, dass jene Hits relevanter sind.

Die Ergebnisse werden an dieser Stelle jedoch noch nicht nach ihrer Qualität bzw. Relevanz sortiert. Die Hits sind derzeit nach ihrer Reihenfolge im Korpus sortiert. Wichtig für das Ranking ist es, eindeutig festzulegen, was einen guten Hit ausmacht. Ein Treffer, der nur  $\alpha$  gefunden hat, ist nicht besser als ein Treffer, der  $\alpha$  und  $\beta$  im gleichen Satz findet. Um die Ergebnisse sortieren zu können, müssen sie dementsprechend einen Rank zugeordnet bekommen. Demnach muss jedem Hit ein Rank übermittelt werden, der mit der Hilfe von Parametern und Features berechnet wird. Dazu wird jedes Feature zuerst implementiert und berechnet, danach werden alle Werte in einer Formel zusammengeführt.

Das Ranking kann nur für mehrere Suchbegriffe angewandt werden. Wenn nach einem Begriff gesucht wird, dann sollen weiterhin einfach nur die Token ausgegeben werden, die auf den Suchstring passen. Wie auch in 4.4.1 ersichtlich werden wird, ist es ohne eine Distanz zwischen zwei gefundenen Ergebnissen nicht möglich, die Formel komplett zu berechnen.

### 4.4 Die Features für das Ranking

Jede Art von Machine Learning und jeder Algorithmus beruht auf Features, die für die Berechnung des Scores unumgänglich sind. Auch für das Ranking einer Suchmaschine müssen diese Charakteristika festgelegt werden.

#### 4.4.1 Formel zur Berechnung des Rankings

$$R_S = \frac{Q}{d} \qquad Q = \sum_{m_i=1}^n \frac{|m|}{|q_m|}$$

Um einen Rank für einen Hit bzw. Satz  $R_S$  zu berechnen, wird der Quotient zwischen der Qualität Q und der Distanz d aller Matches berechnet. Es wird die Anzahl eines Matches m durch die Anzahl der Vorkommen des Matches im Query-String geteilt.

Zur Veranschaulichung:

Wenn der Query-String "ich denke" ist und die Matches aus einem Hit die Token:{ich, ich, denken} beinhalten, dann bedeutet das für die Formel:

$$Q = \frac{|ich|}{|q_{\rm ich}|} + \frac{|denken|}{|q_{\rm denken}|} = \frac{2}{1} + \frac{1}{1} = 3$$

Jedoch gibt es ein paar Belohnungen und Bestrafungen für das Ranking, die in 4.4.4 erklärt werden.

Ein Hit wird den größtmöglichen Wert für einen Rank annehmen, wenn seine Matches nebeneinander zu finden sind, die Matches genau der Suchanfrage entsprechen und es in dem Absatz, in dem der Hit vorkommt, keine anderen Hits gibt, die seinen Rank negativ

 $<sup>^{1}\</sup>alpha$  und  $\beta$  repräsentieren die Suchanfragen.

beeinflussen könnten. Genaueres hierzu in 4.4.5.

Gleichzeitig bekommt ein Hit den schlechtmöglichsten Rank zugeteilt, wenn es nicht alle Token aus der Query gefunden hat, die gefundenen Token zusätzlich weit auseinanderstehen und gleichzeitig nur Abwandlungen, d. h. Vollformen vom Query-Token sind. Dies wird in Kapitel 4.4.3 detailierter präsentiert.

#### 4.4.2 Distanz zwischen Suchtreffern

Als einfachstes, aber gleichzeitig entscheidendes Merkmal gilt die Distanz zwischen zwei Matches  $m_1$ ,  $m_2$  in einem Hit. Daraus folgt logisch, dass  $m_1$  und  $m_2$  höher im Gesamtranking abschneiden, je näher sie zueinander stehen. Die Distanz d nimmt den kleinstmöglichen Wert an, wenn  $m_1$  und  $m_2$  direkt nebeneinanderstehen. Je kleiner d, desto größer ist der Wert des Ranks. In 4.4.4 wird zu sehen sein, dass der Wert der Distanz verkleinert wird, wenn die Matches sehr gut in Bezug zur Suchanfrage waren. Somit kann d noch kleiner werden und  $R_S$  einen noch größeren Wert annehmen.

#### Berechnung der Distanz im Code

In Zeile 2 wird bis zum vorletzten Match im aktuellen Hit iteriert und anschließend die beiden Iteratoren konstruiert. Es wird nur bis zum vorletzten Match iteriert, damit in Zeile 4 der foll\_it die Position des letzten Matches im String annehmen kann. Die Distanz d wird für jeden Hit berechnet und infolgedessen wird die Entfernung zwischen  $m_1$  und  $m_2$ , ...  $m_{n-1}$  und  $m_n$  aufaddiert (vgl. Zeile 6).

```
if (r.matches.size()>=2){
    for(std::vector<Match>::size_type pos = 0; pos != r.
        matches.size()-1; pos++) { // iterate over all matches

    auto curr_it = r.matches[pos].first;
    auto foll_it = r.matches[pos+1].first;
    ...
    distance += tokenDistance(curr_it, foll_it);
    ...
```

Die Funktion **tokenDistance**(...) berechnet die Distanz zwischen zwei Iteratoren. Die Iteratoren repräsentieren zwei Matches. Nachdem beide dereferenziert wurden (6), werden aus dem Token-Match die Positionen geholt und mit einem  $static\_cast < double>$  auf einen double Value gecastet. Das ist nötig, damit im weiteren Verlauf auch rationale Zahlen berechnet werden können und  $size\_t$  nur unsigned integer Werte halten kann, die von  $size\_of$  zurückgegeben werden [10][122]. Wie in 3.2 und 3.4 gezeigt, hat jedes Match einen pos Attribut. Anschließend wird die Position des nachfolgenden Tokens vom aktuellen abgezogen und der Wert zurückgegeben.

```
double pos2 = static_cast < double > (token_2.pos());
distance = pos2 -pos1;
return distance;
}
```

#### 4.4.3 Qualität

Um die Qualität Q zu berechnen, wird eine  $unordered\_map<...>$  **query\\_map** erstellt (Zeile 4). Dafür muss der Query-String gesplittet werden und jedes Wort in einen std::wstring abgewandelt werden (Zeile 9). Das geschieht, weil die **find(...)** Funktion nur std::wstrings als Argument akzeptiert, um den Dictionary Entry für den wide String zu lokalisieren. Schlussendlich wird für jeden Entry die Anzahl der Vorkommen im Query-String gezählt (Zeile 11, 15). Dabei müssen auch wieder double-Werte verwendet werden, damit am Ende der Rank berechnet werden kann.

```
1 std::unordered map<const wf::dict::Entry *, double>
2 wf::RankedSearch::split(std::string query, int &query size) const
3 {
    std::unordered_map<const_wf::dict::Entry *, double>
       query entries;
    std::istringstream splitter(query);
5
    std::string word;
6
    while (splitter >> word) {
      query_size++;
8
      std::wstring wide str = wf::util::widen(word.c str());
      const wf::dict::Entry * token entry = dict ->find(wide str);
10
      if (token entry){
1\,1
        if (query_entries.find(token_entry) != query_entries.end())
12
          query entries [token entry] += (double) 1.0;
13
        } else {
14
          query_entries[token_entry] = (double) 1.0;
15
16
      }
17
18
    return query_entries;
19
20
```

Die Werte für den Nenner der Formel liegen nun vor. Für den Zähler wird ähnlich vorgegangen. Zuerst wird eine  $\mathbf{hit}$ \_ $\mathbf{map}$  erstellt, die eine Kopie der  $\mathbf{query}$ \_ $\mathbf{map}$  ist (Zeile 2). Jedoch werden nach jedem neuen Hit die Values in der  $\mathbf{hit}$ \_ $\mathbf{map}$  auf null zurück gesetzt. Dies sorgt dafür, dass für jeden Hit der  $R_S$  neu berechnet wird und nicht die Frequenzen der alten  $\mathbf{hit}$ \_ $\mathbf{map}$  miteinbezogen werden (Zeile 8).

```
std::unordered_map<const wf::dict::Entry *, double> hit_map;
hit_map.insert(query_map.begin(), query_map.end());

double ab_rank = 0.0;
for (auto& r : res) { // üfr jeden Hit in hits}

for (auto &it : hit_map){ // set each value to 0}
    it.second = 0.0;
}
```

10 . .

In der gleichen for-Schleife wie zuvor wird nun der Dictionary Entry für das aktuelle Match cur\_match\_entry gesucht und damit fillHitMap(...) aufgerufen (Zeile 6). Im letzten Durchlauf der Schleife wird das gleiche auch für das letzte Match im Vektor gemacht (Zeile 11). Danach wird das Match zum aktuellen Hit hinzugefügt. In Kapitel 4.6 wird dies genauer erläutert.

```
if (r.matches.size()>=2){
    for(std::vector<Match>::size type pos = 0; pos != r.matches.
       size()-1; pos++) {
      auto curr it = r.matches[pos].first;
      auto foll it = r.matches[pos+1].first;
      const wf::dict::Entry* cur match entry = (*curr it).entry();
         // Dict-Entry von Match
      fillHitMap (query map, hit map, cur match entry);
      distance += tokenDistance(curr_it, foll_it);
      hit.push back(*curr it); // Match ühinzufgen
      if (pos = r.matches.size()-2){} //Im Letzten durchlauf
        const wf::dict::Entry* foll match entry = (*foll it).entry
10
           ();
        fillHitMap(query_map, hit_map, foll_match_entry);
11
        hit.push back(*foll it); // add last match
12
      }
13
14
```

Anschließend wird die gesamte Qualität Q berechnet und es werden noch einige Parameter festgelegt, die für spätere Belohnungen und Bestrafungen der Scores wichtig sind (Siehe 4.4.4).

```
bool perfect match = true;
          bool all found = true;
2
          double quality = 0.0;
3
          for (auto it : hit map) {
            if ((it.second < query map[it.first]) && (it.second >
                0.0)){
              perfect match = false;
               quality += it.second;
            else if (it.second == 0.0)
              perfect match = false;
              all found = false;
               quality -= query map[it.first] * (double)2; //
11
                  punishment
            } else {
12
               quality += it.second;
            }
14
```

#### Lemmata und inverse Lemmata

Der Funktion wird die **query\_map**, **hit\_map** und der Dictionary Entry des aktuellen Matches übergeben. Nun wird über die **query\_map** iteriert. Zuerst kommt es zu einer

Überprüfung, ob es sich beim Match um ein Partikelverb handelt. Falls das der Fall ist, wird die Qualität des Hits noch einmal verbessert.

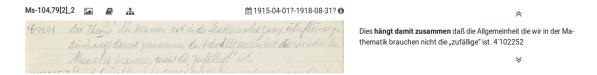


Abbildung 4.1: Suchergebnis zu Partikelverb

Wie in Abbildung 4.1 zu sehen ist, werden auch Partikelverben gefunden. Diese sollen eine Sonderrolle im Ranking einnehmen, indem sie sollen vor den anderen Ergebnissen stehen. Damit das Ranking nicht einen schlechten Wert berechnet, weil die Teile des Partikelverbs weit auseinanderstehen, wird das Ranking positiv angepasst. Das hätte zur Folge, dass das Ranking für Partikelverben sehr schlecht ausfallen würde, daher die Berichtigung (Zeile 6 ff.).

Jetzt kommt es zu einer Fallunterscheidung. Der Query Entry kann nun entweder exakt dem übergebenen Dictionary Entry entsprechen, ein Lemma, ein inverses Lemma sein oder gar nicht dem Match entsprechen. Je nachdem wird der Wert in der hit \_map an der entsprechenden Stelle um den Wert 1.0 oder 0.9 erhöht (Zeile 10, 13, 16). Die Keys der hit \_map bestehen aus Query Entries. Es hätte Komplikationen zur Folge, wenn man die Map mit den Dictionary Entries der gefundenen Treffer füllen würde. Die Konsequenz wäre, dass jede gefundene Vollform eines Wortes einen Eintrag in der hit \_map bekommt. Dadurch werden keine repräsentativen Werte gebildet und die Qualität würde falsch berechnet werden. Beim Berechnen der Qualität wird von folgender Bruchregel Gebrauch gemacht: [6][73-97]

Wenn für die Anzahl an gefundenen Matches m für das Wort denken gilt:

$$|m| = \sum_{m_i=1}^{n} 1$$

Dann gilt:<sup>2</sup>

$$\frac{|m|}{|q_m|} = \sum_{m_i=1}^n \frac{1}{q_n}$$

```
void wf::RankedSearch::fillHitMap(std::unordered map<const wf::
     dict::Entry *, double > &query map,
                       std::unordered_map<const wf::dict::Entry *,
2
                          double > &hit map,
                       const wf::dict::Entry * &match entry) const {
3
      for (auto query pair : query map) {
4
        auto query_entry = query_pair.first;
        size t vPartFound = match entry->gc.str().find(L"+V+#");
6
        if (vPartFound != std::string::npos){
          hit_map[query_entry] += (double) 3.0;
8
        if (query entry == match entry) {
10
          hit map [query entry] += (double) 1.0 / query map [
11
             query entry];
          break;
12
        } else if ((*query_entry).isLemmaOf(match_entry)) {
13
```

 $<sup>^{2}</sup>q_{m}$  ist die Anzahl an Vorkommen von m im Query-String.

#### 4.4.4 Belohnung und Bestrafung von Parametern

Es werden des Weiteren Paramerter festgelegt, die die Formel weiter beeinflussen sollen. Zum jetztigen Zeitpunkt ist die Gewichtung der Parameter festgelegt, soll aber in Zukunft vom User festgelegt werden können (Siehe 6). In der gleichen Schleife wird dann die Qualität des gesamten Hits ausgerechnet (Zeile 7). Ein **perfect\_match** ist dann gegeben, wenn der Suchtreffer genau der Suchanfrage entspricht oder optional noch mehr Token in dem Hit gefunden wurden. **All\_found** soll signalisieren, dass zwar alle Wörter aus der Suche gefunden wurden, aber sie nicht exakt dem Suchtreffer entsprechen. Als Beispiel hierfür wären Vollformen oder Flexionen zu nennen, die ein Query Token teilweise erfüllen. Wenn ein Query Token nicht gefunden wurde, d.h. wenn der **it** gleich 0.0 ist, dann wird die Qualität der Suchanfrage verschlechtert (Siehe 11). Für jedes fehlendes Token wird eine Bestrafung vollzogen.

```
bool perfect match = true;
    bool all found = true;
2
    double quality = 0.0;
3
    for (auto it : hit map) {
      if ((it.second < query map[it.first]) && (it.second > 0.0))
        perfect match = false;
        quality += it.second;
      else if (it.second == 0.0)
        perfect match = false;
        all found = false;
10
        quality -= query_map[it.first] * (double)2; // punishment
11
      }else {
12
        quality += it.second;
13
      }
14
```

r.s\_rank wird mit der in 4.4.1 beschrieben Formel berechnet. Zusätzlich werden hier Belohnungen berücksichtigt, indem der Wert der Distanz kleiner gemacht wird. Außerdem bekommt der Hit seinen Rank zugewiesen (Zeile 3)

```
if (perfect_match){
    r.s_rank = quality / (distance / (double)5.0);
    hit.SetSatzRank(r.s_rank);
4    } else if (all_found) {
    r.s_rank = quality / (distance/ (double)2.0);
    hit.SetSatzRank(r.s_rank);
    } else {
```

#### 4.4.5 Ranking über Absätze

Das Absatz-Ranking addiert für einen Absatz alle Satz-Ranks  $R_S$  zusammen und teilt den Wert durch die Anzahl der Sätze c, die mindestens einen Match hatten.

$$R_A = \frac{\sum_{i=1}^n R_S}{c}$$

Um ein Ranking über Absätze zu erstellen, muss eine unordered\_map erstellt werden, die die Satz-Ranks pro Absatz aufsummiert und gleichzeitig mitzählt, wie oft der einzelne Absatz gefunden wurde. Demnach hat absatz\_rank hat die Form <std::string, std::pair <double,double». Um die richtige Absatznummer zu wählen, kann nicht auf das Json Attribut "abnr" zurückgegriffen werden, denn für jedes Dokument werden die Absätze erneut gezählt. Das hätte zur Folge, dass sich Dokumente und Absätze überschneiden. Daher muss auf "n" (3.2) zurückgegriffen werden. Um sowohl das Dokument als auch den Absatz aus "n" zu extrahieren, muss am Unterstrich gesplittet werden (Zeile 4). In der Funktion getDocumentAbsatz(...) wird genau das umgesetzt.

```
std::string
wf::RankedSearch::getDocumentAbsatz(std::string n) const{
std::smatch match;
std::regex rgx("(.*)_");
if (std::regex_search(n, match, rgx)){
return match[1];
} else {
return n;
} }
```

Anschließend wird der Dokument-Absatz in den akutellen Hit eingetragen.

```
auto& hit = request.AppendHit()
...
setDocAbsatz(docAbsatz);
```

Schlussendlich muss, nachdem der Score für jeden Satz berechnet worden ist, die unordered\_map befüllt werden. Es wird überprüft, ob der Dokument-Absatz schon in der Map vorkommt (Zeile 6), sollte dies nicht der Fall ist, wird ein neues Paar aus s\_rank und dem Wert 1.0 gemacht und in die Map eingfügt (Zeile 7). Im anderen Fall wird einfach s\_rank und die Anzahl erhöht (Zeile 11 ff.).

Nachdem über alle Dokumente iteriert ist, ist die *unordered\_map* vollständig und kann an die Sortierung weitergegeben werden.

### 4.5 Sortierung

Bei der Sortierung tritt die Schwierigkeit auf, dass die klassische sort(...) Methode von C++ zwar ein drittes Argument zulässt mit dem die Art der Sortierung festgelegt werden kann, aber es wird kein Json::Value akzeptiert. Um noch spezifischer zu sein, sort(...) nimmt nur sogenannte RandomAccessIteratoren [1], aber beim Iterieren über ein Json::Value werden Json::ValueIterator verwendet. Daher muss  $.json\_["hits"]$  in einen std::vector übertragen werden (Zeile 6). Dabei wird gleichzeitig der  $ab\_rank$  für jeden Hit berechnet und gesetzt. Die Summe über alle  $R_S$  und c ist schon in 4.4.5 berechnet worden. Jetzt kann der Quotient berechnet werden und dem Key im Hit zugewiesen werden (Zeile 5).

```
std::vector<Json::Value> copy;
for (auto i : request.json_["hits"]) {
    auto key = i["DocAbsatz"].asString();
    double ab_rank = absatz_ranks[key].first / absatz_ranks[key].
        second; //durch Anzahl der Hits von dem Absatz teilen.
    i["ab_rank"] = ab_rank;
    copy.push_back(i);
}
```

Die Kopie wird mit einem sogenannten *Compare* sortiert, mit dem festgelegt wird, nach welchem Value die Objekte angeordnet werden sollen. Es werden dabei zwei Objekte auf ihren Absatz-Rank verglichen. [2]. Der Aufruf wird von std::sort(...) ausgeführt.

```
bool
sortByAbsatz(const Json::Value& a, const Json::Value& b)
{// For sorting Hits
return a["ab_rank"].asDouble() > b["ab_rank"].asDouble();
}
```

Bevor das sortierte Array wieder an das .json\_["hits"] Array übergeben werden kann, werden noch die einzelnen Hit Nummern "no" gesetzt und die Größe des Arrays auf die "max" Variable reduziert. Da es sich wieder um Json::Values handelt, muss der Counter mit static\_cast auf einen Json unsigned integer gecastet werden.

```
Json::Value sorted_hits;
size_t c = 0;
for (auto it : copy){
   if (c <= max) { // max wird von search.js auf 200 gesetzt
```

### 4.6 Zusammensetzung des Programms

Da bestimmte Parameter sollen später auch vom User festgelegt werden können (Siehe 6), beschreibt folgender Ablauf nur den aktuellen Stand des Programms.

Zuerst wird überprüft, ob der Query-String mindestens zwei Token enthält bzw., ob der Query-String dem Pattern in Zeile 6 entspricht. Anschließend werden die | Operatoren eingefügt und der neue Query-String gesetzt (Zeile 7).

```
1 bool
2 wf:: Searcher:: qualifies for ranked (Message& request)
3 { //Checks if user query is fitted for ranked search, if so,
     changes query to perform rankedSearch
    std::string query = request.GetQueryStr();
    request.oldquery = query;
5
    if (std::regex match(query, std::regex("^[\\äöüßw]+(\\s[\\äöüßw
       ]+)+$"))){
      query = std :: regex replace (query, std :: regex("\\s"), " | ");
          //changes Query to ranked searching
      request.setRankedQueryStr(query);
8
      return true;
    } else return false;
10
11
```

Etwas weiter in der Searcher.cpp wird nun die gerankte Suche durchgeführt und die RankedSearch Klasse dafür initialisiert. Da buildAnswer(...) und ranked(...) jeweils die Referenz [10][97] des request-Message Obekts verändert, kann die Funktion einfach aufgerufen werden.

```
if (ranked_search)
{
    RankedSearch ranked(request, res, max, dict_);
    else buildAnswer(request, res, max);
}
```

In RankedSearch.cpp angekommen, wird zuerst die query\_map und hit\_map erstellt (Siehe 4.4.3). Anschließend werden über die einzelnen Hits iteriert, die in res gespeichert sind (Zeile 1). Daraufhin wird für diesen Hit der Dokument-Absatz erstellt (Siehe 4.4.5). In Zeile 4 ff. wird der neue Hit mit seinen Attributen ins Array hinzugefügt.

```
for (auto& r : res) { // üfr jeden Hit in hits
```

Im Anschluss wird über die Matches im aktuellen Hit iteriert, Distanz und Qualität werden berechnet, Belohnungen und Bestrafungen miteinbezogen und schlussendlich wird der Satz-Rank festgelegt, wie in 4.4.2, 4.4.3 und 4.4.4 beschrieben. Danach wird die unordered\_map absatz\_ranks gefüllt. Das richtige Setzen der Anzahl der Hits bildet den Abschluss des Rankingverfahrens.

```
request . SetNumberOfHits (request . json_ ["hits"] . size ());
```

Nachdem über alle gefundenen Hits iteriert wurde, wird die Sortierung vorgenommen, wie in 4.5 erklärt. Der sortierte Vektor wird am Ende wieder dem Hit-Array übergeben und hat damit seine finale Form.

### 4.7 Neues Json Objekt

Da das Ranking die Berechnung neuer Parameter erfordert, wird das Json Objekt verändern. Es werden sowohl der Rank für jeden Satz  $R_S$ , als auch der Rank für den Absatz  $R_A$  berechnet. Hinzu kommt, dass jeder Hit einen Absatznamen benötigt. Die normale Absatznummer "abnr" kann dazu nicht verwendet werden, da die Absatznummern für jedes Dokument neu gezählt werden. Daher wird das Json Objekt um drei Key Value Paare erweitert:

- "DocAbsatz" enthält den String, der aus "n" gesplittet wird. Er gibt an, in welchem Dokumentenabsatz sich der Hit befindet.
- "s\_rank" enthält einen double Wert, da jsoncpp Probleme mit Floats hat. Jsoncpp schient Floats nicht akzeptieren zu wollen, daher die double Alternative. Der Key gibt Auskunft über den Satz Rank.
- "ab\_rank" ist ähnlich zu "s\_rank" und kann sogar den gleichen Wert annehmen. Sobald es aber mehrere Hits in einem Absatz gibt, ist das nicht mehr der Fall und ab rank beschreibt den gesamten Absatz.

# 5 Ergebnisse des Rankingverfahrens

Die Ergebnisse des Rankingverfahrens sollen bestätigen, dass alle Zielsetzungen erreicht wurden. Die Suchanfragen wurden mit Hilfe von  $wf\_client$  und  $wf\_server$  abgefragt. Die Ergebnisse werden mit alten Suchergebnissen verglichen und sollen die Unterschiede zum neuen Ranking deutlich machen. Die alten Ergebnisse werden von der WiTTFind Website wittfind.cis.uni-muenchen.de bezogen. Auch wenn die Website alle Dokumente durchsucht und in den folgenden Beispielen nur fünf Dokumente berücksichtigt werden, wird durch die Bespiele dennoch deutlich, was sich verändert hat. Die einzelnen Abschnitte geben die Suchanfragen an. Zur Verdeutlichung, dass das Absatzranking auf mehreren Dateien funktioniert, wird der Client Befehl um vier weitere Dateien erweitert:

```
./bin/wf_client —query "SUCHANFRAGE" —f ../data/Ts-213\_NORM-tagged.xml —f ../data/Ts-211\_OA\_NORM-tagged.xml —f ../data/Ts-212\_OA\_NORM-tagged.xml —f ../data/Ms-142\_OA\_NORM-tagged.xml — max 20
```

Die Files müssen auch dem Ordner data hinzugefügt werden. Außerdem müssen die Namen in der files. txt Datei des Servers ergänzt werden.

Zur übersichtlicheren Darstellung des resultierenden Json Objekts wird ein jq Befehl genutzt:

Die Ansicht ist nur zu Zwecken der besseren und intuitiveren Verständlichkeit gewählt.

### 5.1 Beispiele

#### 5.1.1 Gedanken Welt

Wie in 1 schon erwähnt, werden zu diesen Suchbegriffen keine Ergebnisse gefunden. Jetzt werden in nur fünf Dokumenten über 514 Treffer gefunden. Selbstverständlich haben nicht alle dieser Treffer einen guten Rank, aber wie man dem unteren Json Objekt entnehmen kann, besitzt der erste Hit direkt beide Treffer innerhalb eines Satzes. Die Begriffe sind auch nur drei Positionen voneinander entfernt. Daher steht der Hit an erster Stelle. Beim zweiten Treffer ist es ähnlich. Bei den Hits vier bis neun sind die Hits auf den Absatz Ts-212,1062[1] verteilt. Im Anhang werden auf Seite 47 die gekürzten Hits aufgelistet.

Es werden jetzt also Hits gefunden und nach ihrem Rank geordnet, die nicht unmittelbar nebeneinanderstehen und nach ihrem Rank geordnet. Die Suche wurde damit von keinem Treffer auf mehrere verbessert.

#### 5.1.2 man weiß

Die Query Token wurden mit der alten Suche relativ oft hintereinander gefunden und sind im natürlichen Sprachgebrauch in der Reihenfolge verbreitet. Die neuen, zusätzlichen Hits sollen also nach den alten Hits angezeigt werden. Die Suchergebnisse zeigen zuerst vier perfekte Treffer, in denen beide Wörter hintereinander vorkommen. An sechster und siebter Stelle stehen wieder zwei Hits aus dem gleichen Absatz. In diesem Absatz wurden viele Wörter gefunden, die mit der Suchanfrage übereinstimmen. Daher fällt der berechnete Rank hoch aus. Hinzu kommt ein perfektes Match, bei dem beide Suchwörter wieder direkt aufeinanderfolgen. In Zeile 14 ist erkennbar, dass ein Token welches nicht dem Query Token entspricht zwar immer nohe gefunden wird, jedoch der s\_rank-Wert niedriger ist.

```
"DocAbsatz": "Ts-213,55r[1]",
2
     "ab rank": 5.58,
3
     "s_rank": 10,
     "matches": [
       {"token": ["ßwei", "man"],
          "pos": [3,4]
7
8
9
     "DocAbsatz": "Ts-213,55 r [1]",
10
     "ab rank": 5.58,
11
     "s rank": 1.16,
12
     "matches":
13
       \{"token": ["man", "wissen", "man"],
14
         "pos": [2, 4, 7]
15
16
17
```

Das Absatzranking funktioniert also auch wie erwünscht. Zwei gute Hits innerhalb eines Absatzes bilden zusammen wieder einen guten Score für das Ranking. Ein positiver Nebeneffekt des Rankings besteht darin, dass die Reihenfolge in der die Suchbegriffe eingegeben werden keine Rolle spielt.

Bei genauerem Hinsehen finden sich dennoch Probleme. Betrachtet man die Hits auf weiter hinterliegenden Positionen, ist beispielsweise erkennbar, dass ein Hit a eigentlich einen sehr hohen  $s\_rank$  zugewiesen bekommt. Jedoch kommt in einem anderen Satz b im gleichen Absatz nur ein Wort der Suchanfrage vor. Daraus folgt, dass das schlechte Ranking von Satz b das perfekte Ranking von Satz a zu sehr beeinflusst. Die Bestrafungen und Belohnungen müssen folglich noch genauer angepasst werden. Ferner gibt es Hits mit vielen

Matches, die aber nicht hoch genug eingestuft werden. Auf Position 97 ist folgendes Objekt zu finden:

```
{
1
    "DocAbsatz": "Ts-212,1812[1]",
2
    "ab rank": 0.85,
3
    "s rank": 0.85,
4
    "matches": [
5
      {"token": ["man", "wissen", "man", "wisse", "man", "man", "man"],
6
         "pos": [7,8,15,18,28,36,47]
       }]
8
 }
```

Es ist der einzige Hit im Absatz Ts-212,1812[1], jedoch wurde das Wort "weiß" nur in einer anderen Vollform gefunden. Der Hit wird dafür aber zu drastisch bestraft. Folglich muss noch einiges an "Finetuning" betrieben werden, um das Ranking zu optimieren. Das gesamte Ergebnis ist auf Seite 49 zu finden.

#### 5.1.3 Haus Stadt

Bei der Suche nach "Haus Stadt" wurde bisher nur ein Ergebnis gefunden.

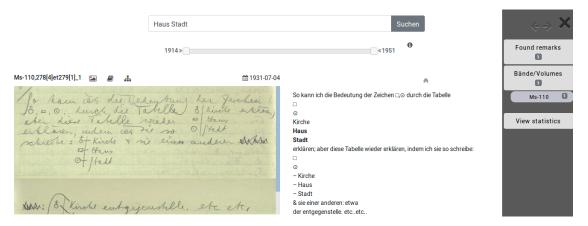


Abbildung 5.1: "Haus Stadt" Suche ohne Rankingalgorithmus

Nun werden Ergebnisse gefunden, in denen die Query Token im Kontext miteinander vorkommen. Der Rank wird schlechter, je weiter die Positionen der Matches auseinander liegen. Der erste Absatz, der in zwei Sätzen ein Match hat, ist *Ms-142,12[2]et13[1]*. Beim Betrachten der Treffer steht fest, dass der Absatz wohl zu Recht auf die ersten paar Treffer folgt.

```
1
     "DocAbsatz": "Ms-142,12[2] et13[1]",
2
     "ab_rank": 0.44625,
3
     "s rank": 0.58,
4
     "matches": [
5
       {"token": ["äHusern", "Stadt", "Stadt"],
6
         "pos": [4,11,14]
       }]
8
  },{
9
     "DocAbsatz": "Ms-142,12[2] et13[1]",
10
     "ab rank": 0.44625,
11
```

Auf Seite 50 ist wieder ein größerer Teil der Suchanfrage zu finden.

#### 5.1.4 Gott Religion Teufel

Auch längere Suchanfragen stellen kein Problem bei der Suche dar. Nur die Dauer der Suche verlängert sich, weil jeder Begriff einzeln gesucht werden muss. Die Menge der sinnvollen Ergebnisse nimmt mit zunehmender Hitnummer stark ab, weil keine Absätze mehr gefunden werden, in denen alle drei bzw. möglichst viele der Matches vorkommen. Das Beispiel soll verdeutlichen, dass auch Begriffe, die wahrscheinlich im Kontext voneinander Vorkommen, auch richtig durch die Rankingformel bewertet werden. Der untere Code-Abschnitt zeigt den am besten bewerteten Absatz. Es wurden alle drei Wörter in einem Satz gefunden und noch ein zusätzliches Wort in einem späteren Absatz.

```
{
1
   "DocAbsatz": "Ts-211,125[7]"
2
   3
   "matches": [
5
    {"token": ["Religion", "öGtter", "Teufeln"],
6
      "pos": [9,11,15]
7
    } |
8
9
   "DocAbsatz": "Ts-211,125[7]"
10
   11
   12
   "matches": [
13
    {"token": ["Teufeln"],
14
      "pos": [11]
16
17
```

Der Treffer wird dann auf folgende Art und Weise auf WiTTFind zu sehen sein.



Abbildung 5.2: Ansicht des Ergebnisses zur Suchanfrage "Gott Religion Teufel"

#### 5.1.5 Realität Vorstellung Ding selbstverständlich

Auch die Eingabe sehr langer Suchanfragen funktioniert. In diesem Beispiel soll auf die Eingliederung auf der Website aufmerksam gemacht werden. Wie die Abbildung 5.3 zeigt,

funktioniert die grafische Darstellung einwandfrei, jedoch kommt es zu Komplikationen mit dem *wfa* Server, der die Reihenfolge der Dokumente verwirft. Damit das Ranking auch die gewünschte Reihenfolge hat, muss an dieser Stelle noch etwas unternommen werden.

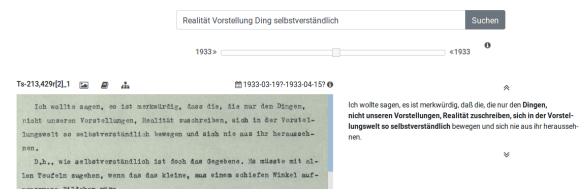


Abbildung 5.3: Ansicht des Ergebnisses zur Suchanfrage "Realität Vorstellung Ding selbstverständlich"

Der gezeigte Hit konnte hier nur durch Eingrenzen der Suche auf das Dokument Ts-213 an die erste Position gehoben werden. Aber  $wf\_client$  zeigt die Ergebnisse in der richtigen Reihenfolge an. Die Sortierung funktioniert, wird aber verworfen bevor sie auf der Website angezeigt werden kann.

```
"DocAbsatz" : "Ts-213,429 r [2]",

"DocAbsatz" : "Ts-213,429 r [2]",

"DocAbsatz" : "Ts-212,1190[1]",

"DocAbsatz" : "Ts-212,1190[1]",

"DocAbsatz" : "Ts-211,763[2] et764[1]",

"DocAbsatz" : "Ts-211,763[2] et764[1]",

"DocAbsatz" : "Ts-211,763[2] et764[1]",
```

Das obere Ranking zeigt die tatsächliche Reihenfolge der Treffer.

#### 5.2 Fazit

Das Grundgerüst für das Rankingverfahren steht und ist einsatzbereit, dennoch bietet es noch Verbesserungsmöglichkeiten. Somit wurden die formuliterten Zielsetzungen aus 4.2 erfolgreicht erreicht. Die Nutzer der Suchmaschine können jetzt Mehrwortanfragen eingeben und sich sicher sein, dass auf jeden Fall Resultate zurück kommen und, dass die Resultate nach ihrer Relevanz sortiert sind. Der Rankingalgorhithums wurde implementiert und beachtet Features wie Lemmata, Distanzen und die Qualität der Suchanfrage. Im Ergebnis funktioniert das Rankingverfahren.

## 6 Ausblick

#### 6.1 Verbesserungsmöglichkeiten

In zukünftigen Arbeiten kann das Rankingverfahren und die Suche allgemein weiter verfeinert werden.

#### 6.1.1 Einzelne Matches

Derzeit wird nur ein fester Value berechnet, wenn ein Satz nur ein Match hat, d. h. die Suchanfrage kann gar nicht zu 100% erfüllt sein. Hier kann noch eine Erkennung der Lemmata eingefügt werden. Dies führt wiederum zum Problem, dass keine Distanz ermittelt werden kann, denn dazu werden mindestens zwei Matches benötigt.

#### 6.1.2 Frequenzlisten

Bei bestimmte Suchanfragen, wie "Gedanken über die Welt", ist es eindeutig, dass der Suchende nach den Begriffen "Gedanken" und "Welt" sucht, als nach den beiden anderen Suchbegriffen. Derzeit wäre es jedoch so, dass alle Wörter gleich gewichtet sind. Zwar wurde mithilfe der Qualität schon der Anfang für eine Gewichtung gemacht, aber eine bessere Methode wäre es, die Informativität der Suchbegriffe zu ermitteln und diese Werte mit in die Formel einzubeziehen. Dabei kann auch von TF-IDF Verwendung gemacht werden. TF-IDF berechnet, wie relevant ein Wort in einem Dokument bzw. Korpus ist. Wörter, die nur in einer kleinen Menge an Dokumenten vorkommen, haben einen höheren TF-IDF Wert als beispielsweise Stoppwörter [7][134]. Jedoch muss so eine Datenbank zuerst aufgebaut werden und kann nicht bei jeder Suchanfrage neu berechnet werden. Wenn sie erstellt wird, kann der TF-IDF Faktor eine wichtige Rolle im Ranking einnehmen.

#### 6.1.3 Eingliederung in die Website

Die Ergebnisse können derzeit noch nicht auf der Website angezeigt werden, da der wfa-Server die sortierten Ergebnisse wieder verwirft und nur bestimmte Parameter zurückliefert. Wie dem unteren Code zu entnehmen ist, übergibt wfa nicht das gesamte Json Objekt, das es von wf zurückbekommt.

```
json response: Dict = \{\}
       for item in json wf response ["hits"]:
            matches = item ["matches"]
3
            if matches:
                date = item["date"] if "date" in item else ""
5
                date norm = item ["date norm"] if "date norm" in item
6
                current triple = {
                     "start": [matches [0] ["pagenr"], matches [0] ["
                         linenr"],
                                 matches [0] ["tokennr"]],
                     "end": [\text{matches}[-1][\text{"pagenr"}], \text{matches}[-1][\text{"}]
10
                         linenr"],
                               \operatorname{matches}[-1]["tokennr"]],
11
```

```
"name": item["n"],
"date": date,
"date_norm": date_norm}
document_name = item["path"]
number = item[anchored_nr]

return json_response
```

Daher muss die Sortierung veraussichtlich nach search. js ausgelagert werden und der wfa Server muss die Ranks der einzelnen Hits zurückliefern. In search. js findet auch schon die Sortierung nach Datum statt, daher wäre es sinnvoll, auch dort die Ergebnisse nach ihrem Ranking zu ordnen. Außerdem soll dem User der Suchmaschine ermöglicht werden einzustellen, ob die Suchergebnisse gerankt ausgegeben werden sollen, oder auf die normale, alte Weise. Dazu soll auf der Website ein Knopf implementiert werden, der einen Boolean auf True setzt und diesen an die wf Suche weitergibt. In 4.6 ist nachzulesen, dass wf schon einen solchen Boolean besitzt, der derzeit jedoch noch anders festgelegt wird.

#### 6.1.4 Index aufbauen für Suche

Ein Problem bei längeren Suchanfragen stellt die Dauer der Suchanfrage dar, da jeder Begriff gefunden werden muss. Um diesen Suchprozess zu verschnellern, könnte ein invertierter Index aufgebaut werden, der jedes Wort in einem Dictionary zu seinen Postings mapped [5][7]. Dadurch hätte man zu jedem Suchbegriff schon eine Liste von Dokumenten oder auch Absätzen, bevor überhaupt gesucht wird. Es müssten nicht aufs Neue alle Dokumente durchsucht werden, sondern nur noch eine limitierte Anzahl an Dokumenten. Bei mehreren Suchbegriffen könnte dann auf die Keys der einzelnen Wörter oder Teilwörter zugegriffen werden und aus den erhaltenen Dokumenten eine Vereinigung gebildet werden. Indexierung ist bei den meisten Suchmaschinen schon Standard und könnte WiTTFind auch weiter helfen, obgleich der Nachlass von Wittgenstein nicht mit der Suche im Internet vergleichbar ist.

#### 6.1.5 Rankingstatistik

Es ist sinnvoll eine Statistik über die Verteilung der Ranks in Bezug zur Suchanfrage anzulegen. Dabei sollte man über längere Zeit Suchanfragen auf der Website tracken, um einen großen Datensatz zu haben, den man auswerten kann. Dadurch könnte man die Gewichte für die Features am besten ermitteln. Auch könnte man die durchschnittliche Länge einer Suchanfrage herausfinden und besser erkennen, wie WiTTFind genutzt wird.

#### 6.2 Schluss

Dank Sphinx wurde eine grafisch ansprechende Website erstellt, die neuen Bachelor-/Masterstudierende einen schnelleren Einstieg in das Gesamtprojekt WAST ermöglicht. Zusätzlich ist es einfach die Website instand zu halten, da nur eine README Datei neu eingelesen werden oder neu hinzugefügt werden muss. Wie in der Arbeit aufgezeigt wird, kann man das Ranking noch weiter verfeinern, um die Suchmaschine und Ergebnisse zu verbessern, damit WiTTFind noch effektiver, schneller, handlicher und nützlicher wird. Einige der genannten Aspekte würden auch die Suche der Token so stark beschleunigen, dass man noch viel mehr Features einbauen könnte, um das Ranking weiter zu optimieren. Das Ziel, jedem Hit einen Rank zuzuordnen und nach diesem Wert zu sortieren, wurde erreicht. Außerdem kann jetzt nach Begriffen gesucht werden, die nicht zwangsweise hintereinanderstehen müssen. Die Scores für das Ranking wurden anhand eigener Features berechnet, die für das Ranking als relevant eingestuft worden sind. Es ist einfach, das Ranking nach Belieben anzupassen und

neue Gewichtungen für Features auszuprobieren. Im Ergebnis besitzt WiTTFind jetzt ein funktionierendes Mehrwort-Rankingverfahren.

# 7 Anhang

#### **7.1** Code

find(...) aus Search.hpp

```
1 template < class It >
3 wf::find(const Automaton& a, It& b, It e, const Skip& skip)
    using Pair = std::pair < It, const State *>;
    {\tt using Stack} = {\tt std} :: {\tt stack} {<\! Pair \! >};
    using Set = std::set<std::pair<const State*, const Token*>>;
    Stack stack;
    Set visited;
    stack.emplace(b, a.initial()); // Fill stack
                                                     Inital() liefert
       state aus Automaton
    while (not stack.empty()) {
11
      auto top = stack.top(); // take first element from Stack
12
      auto p = std :: make\_pair(top.second, \&(*(top.first))); //Pair
13
         aus state-Iterator
      stack.pop(); // Removes top element from Stack
14
      if (visited.count(p) == 0) { // if Element not in visited (
15
         set)
        visited.insert(p); // add to visited
16
        if (not top.second->final()) { // if atomaton not final //
17
           ßHEIT?
          // top.second is autoamaton. Automaton has State.==>
              calls delta from state.cpp
          top.second->delta(top.first, e, skip, stack); // do
19
              something with iterator ...
          //stack wird auch erweitert //
        } else { // if Word is found
          b = top.first; // b wird zu b, dass geliefert wurde.
          return true; // found something
24
25
    return false;
28
```

#### 7.1.1 Searcher.cpp

```
void
void
searcher::search(Message& request, bool ranked_search) const
{
```

```
using future = std::future < Search Document::Message >;
5
    std::vector<future> futures;
    auto query = request.GetQuery(); // query string
    auto builder = request.GetBuilder(dict_); // getbuilder abstact
8
        query builder.// Baut Suchautomat in Message.cpp
    assert (query);
9
    assert (builder);
10
11
    request.ForEachPath([&](const std::string& path) {
12
      auto automaton = query->GetAutomaton(*builder); // query
13
          token to SimpleQuery.hpp.
      auto skipper = request.GetSkipper(); //Skipper for
14
          Satzzeichen
      auto search = SearchDocument (
                                             // search bekommt die
15
        std::move(automaton), std::move(skipper), files ->Get(path)
16
            );
17
        futures.push back(std::async(std::launch::async, std::move(
18
            search)));
    });
19
    SearchDocument:: Message res; // vektor aus Hit Objekten
20
    for (auto& future : futures) { // for each hit
21
      auto tmp = future.get(); //tmp = hits
22
      std::copy(begin(tmp), end(tmp), std::back inserter(res));
23
^{24}
    auto max = request.GetMax();
25
    \max = \max ? \max : std :: numeric limits < size t > :: max();
26
    if (ranked search)
27
28
      //buildRankedAnswer(request, res, max);
29
      RankedSearch ranked(request, res, max, dict);
       else buildAnswer(request, res, max);
31
32 }
```

```
2 wf:: Searcher:: build Answer (Message& request,
                                   const SearchDocument:: Message& res,
3
                                   size t max) const
4
5
       size_t = 0;
6
       size t i = 0;
       \mathtt{std}:: \mathtt{unordered\_map} {<} \mathtt{std}:: \mathtt{string} \ , \ \mathtt{size\_t} {>} \ \mathtt{counts} \ ;
       for (const auto\& r : res) \{ // "ufr jeden Hit" in hits]
         n += r.matches.size(); // n + Anzahl der Matches
          for (const auto& m : r.matches) { // over all Matches
11
               if (\max \le \text{counts}[r.path]++) // Abbruchkriterium
12
                 break;
13
            // Hits werden erstellt.
14
            auto& hit = request.AppendHit()
15
                              . SetPath (r.path)
16
```

```
. SetN(r.n)
17
                            . SetF (r.f)
18
                            . SetAbnr(std::stoul(r.abnr))
19
                            . SetSatznr(std::stoul(r.satznr))
                            . Set Counter (i++);
21
           // insert other
22
                  for (const auto& p: r.other) {
23
                    hit.SetKeyVal(p.first, p.second);
24
           // insert Matches
26
           for (auto i = m. first; i != m. second; ++i)
27
              hit.push back(*i);
28
         }
29
       }
       request . SetNumberOfHits(n);
^{32}
```

#### 7.1.2 reply json

```
{
1
       // Die antwort des servers enthaelt alle variablen der
2
          Anfrage.
       "builderType" : "SimpleBuilder",
3
        // hits ist ein array in denen die treffer stehen.
5
        // jeder hit repraesentiert eine gefundene textstelle
6
       "hits":
7
8
           // jeder hit ist ein json objekt, dass alle informationen
9
                ueber
           // den treffer enthaelt.
10
1.1
         // absatznummer der textstelle
12
         "abnr": 54,
13
         // f idenifikation der textstelle
         "f" : "Ts-213, iii -r",
15
         // matches ist ein array, das informationen ueber die
16
             gematchten token enthaelt
         "matches" :
17
18
       {
19
           // grammatical code (lexikon information des gematchten
^{20}
               token)
           "gc": "+V+refl+tr",
21
           // inflectional code (lexikon information des gematchten
22
               token)
           "ic": ":2 \,\mathrm{mGi}:3 \,\mathrm{eGi}",
23
           // eindeutige interne id des token
24
           "id": 10383,
25
           // offset position des token im satz
26
           "pos" : 3,
27
28
           // treetagger tag
```

```
"tag" : "VVFIN",
29
            // der gematchte token
30
            "token" : "denkt"
31
       }
33
         // n identifikation der textstelle
34
         "n" : "Ts-213, iii -r [8]_1",
35
         // nummer im array (unnoetig?)
36
         "no" : 0,
         // pfad zum dokument, in dem die textstelle steht.
38
         "path": "/home/flo/devel/wast/wf/build/../data/Ts-213_NORM
39
             -tagged.xml",
         // satznummer der textstelle
40
         "\,satznr\,"\ :\ 128
41
     },
42
43
         "abnr" : 73,
44
         "f": "Ts-213, iv-r",
4.5
         "matches":
46
47
       {
            "gc": "+V+intr+refl+tr",
49
            "ic": 1 \text{ mGc} : 1 \text{ mGi} : 3 \text{ mGc} : 3 \text{ mGi} : OI",
50
            "id": 7938,
51
            "pos" : 14,
52
            "tag" : "VVFIN",
53
            "token" : "denken"
54
       }
55
56
         "n": "Ts-213, iv-r[8] 1",
57
         "no": 1,
58
         "path": "/home/flo/devel/wast/wf/build/../data/Ts-213 NORM
59
             -tagged.xml",
         "\,satznr\,"\ :\ 179
60
     },
61
62
63
       "\max" : 25,
       // insgesamte anzahl der treffer die gefunden worden waehren,
65
            wenn \max = 0
       "numberOfHits" : 275,
66
          ausgabedatei (name der datei, in die diese antwort
67
           geschrieben wurde (wird nur von wf server, wf client und
           wf display benoetigt)
       "outputFile" : "/dev/stdout",
68
       "paths":
69
70
     ".../ data/Ts-213 NORM-tagged.xml"
71
       ],
72
       "query" : "denken",
       "queryType" : "UserQuery",
74
       // antwortzeit des servers (wird von wf_client hinzugefuegt)
7.5
       "responseTimeMs" : 402,
76
```

#### 7.1.3 RankedSearch.cpp/.hpp

```
1 #ifndef wf_RankedSearch_hpp__
2 #define wf_RankedSearch_hpp__
5 #include "SearchDocument.hpp"
6 #include "Message.hpp"
7 #include <regex>
10 namespace wf {
11 namespace dict {
12 class Dictionary;
14 struct Entry;
15 class Message;
 class RankedSearch
18 {
 public:
19
    RankedSearch (Message& request, SearchDocument:: Message& res,
       size t max, const dict::Dictionary * dict );
    void buildRankedAnswer (Message& request,
^{21}
                              SearchDocument:: Message& res,
                              size_t max) const;
23
24
25
  private:
26
    double tokenDistance(std::vector<Token>::const_iterator&
       curr_it,
                       std::vector<Token>::const iterator &foll it)
28
                           const;
    std::unordered_map<const_dict::Entry *, double> split(std::
29
       string query,
                                                              int &
30
                                                                  query_size
                                                                  const;
    void fillHitMap(std::unordered map<const dict::Entry *, double>
31
        &query_map,
                     std::unordered map<const dict::Entry *, double>
32
                         &hit_map,
                     const dict::Entry * &match entry) const;
33
    double fillSingleHitMap(std::unordered_map<const dict::Entry *,
34
        double > &query_map,
                     int query_size,
35
                     const dict::Entry * &match entry) const;
36
```

```
1 #include "RankedSearch.hpp"
2 #include "Message.hpp"
3 #include "dict/Dictionary.hpp"
4 #include "util/string.hpp"
5 #include <algorithm>
6 #include <iostream>
8 bool
9 sortByAbsatz (const Json :: Value& a, const Json :: Value& b)
10 {// For sorting Hits
      return a ["ab rank"]. asDouble() > b ["ab rank"]. asDouble();
11
12
13
  wf::RankedSearch::RankedSearch(Message& request,
                                    SearchDocument:: Message& res,
16
                                    size t max,
17
                                    const dict::Dictionary* dict)
18
    : dict (dict)
19
    buildRankedAnswer(request, res, max);
^{21}
22 };
23
25 wf::RankedSearch::buildRankedAnswer(Message& request,
                                    SearchDocument:: Message& res,
26
                                    size t max) const
27
28
    std::unordered_map<std::string, std::pair<double,double>>
29
       absatz ranks; // um Absatz Rank zu äzhlen
    // map each query entry to its frequenzy
    int query size = 0;
31
    if (res.size() == 0){
32
      return; // Funktion verlassen!
33
34
35
    std::unordered_map<const wf::dict::Entry *, double> query_map =
36
        split (request.oldquery, query_size); // Get Query Entrys
    // maps each query entry to its hit frequenzy
37
    std::unordered_map<const wf::dict::Entry *, double> hit_map;
38
    hit map.insert(query map.begin(), query map.end());
39
40
```

```
double ab rank = 0.0;
41
    for (auto& r : res) { // üfr jeden Hit in hits
42
43
      for (auto &it: hit map) { // set each value to 0
        it.second = 0.0;
4.5
46
        std::string docAbsatz = getDocumentAbsatz(r.n);
47
        auto& hit = request.AppendHit()
48
                        . SetPath (r. path)
49
                        . SetN(r.n)
50
                        . SetF ( r . f )
51
                        . SetAbnr(std::stoul(r.abnr))
52
                        . SetAbsatzRank (ab rank)
53
                        . SetSatznr (std :: stoul (r.satznr))
54
                        . SetDocAbsatz (docAbsatz);
55
        double distance 0.0;
57
        if (r.matches.size()>=2){
58
          for(std::vector < Match > :: size\_type pos = 0; pos != r.
59
              matches.size()-1; pos++) { // Matches iterieren}
             auto curr it = r.matches[pos].first;
             auto foll it = r.matches[pos+1].first;
61
             const wf::dict::Entry* cur_match_entry = (*curr_it).
62
                entry(); // Dict-Entry von Match
             fillHitMap(query_map, hit_map, cur_match_entry);
63
             distance += tokenDistance(curr_it, foll it);
             hit.push_back(*curr_it); // add match
66
67
             if (pos = r.matches.size()-2)\{ //Im Letzten durchlauf
68
                auch den nachfolgenden Iter adden
               const wf::dict::Entry* foll match entry = (*foll it).
69
                  entry();
               fillHitMap (query map, hit map, foll match entry);
70
               hit.push_back(*foll_it); // Letztes Match zu Matches
71
                  adden
             }
72
          }
74
          // calculate match quality
7.5
          bool perfect_match = true;
76
          bool all found = true;
77
          double quality = 0.0;
78
          for (auto it : hit map) {
             if ((it.second < query_map[it.first]) \&\& (it.second >
                0.0)){
               perfect_match = false;
81
               quality += it.second;
82
             else if (it.second == 0.0)
83
               perfect match = false;
               all found = false;
85
               quality -= query_map[it.first] * (double)2; //
86
                  punishment
```

```
}else {
87
                quality += it.second;
88
89
           }
91
92
           if (perfect match) {
93
           r.s rank = quality / (distance / (double) 5.0);
94
              hit.SetSatzRank(r.s_rank);
           } else if (all_found) {
96
              r.s\_rank = quality / (distance / (double) 2.0);
97
              hit.SetSatzRank(r.s_rank);
98
           } else {
99
              r.s rank = quality / distance;
100
              hit.SetSatzRank(r.s rank);
101
102
           fill unord map(absatz ranks, docAbsatz, r.s rank); //
103
               increase absatz rank
104
         } else {
105
           r.s. rank = (double) 1 / (query. size *2.0); //durch. die
106
               Query Size teilen.
107
           hit.SetSatzRank(r.s_rank);
108
           fill unord map(absatz ranks, docAbsatz, r.s rank); //
109
               count Absatz Rank mit
110
           for (auto m : r.matches) {
111
              for (auto i = m. first; i != m. second; ++i)
112
                hit.push back(*i);
113
114
           }
115
         }
117
           for (const auto& p: r.other) {
118
              hit.SetKeyVal(p.first, p.second);
119
           }
120
122
     request.SetNumberOfHits(request.json_["hits"].size());
123
124
125
     // Json array needs to be coppied,
126
        otherwise not sortable because sort needs
127
        RandomAccessIterator not Json::ValueIterator)
     std::vector<Json::Value> copy;
128
     for (auto i : request.json_["hits"]) {
129
       auto key = i["DocAbsatz"].asString();
130
       double ab_rank = absatz_ranks[key].first / absatz_ranks[key].
131
          second; //durch Anzahl der Hits von dem Absatz teilen.
       i ["ab_rank"] = ab_rank;
132
       copy.push_back(i);
133
134
```

```
std::sort(copy.begin(), copy.end(), sortByAbsatz); // Sort Hits
135
         by Absatz
136
     //sortierten Vector wieder in Json::Value verwandeln
137
     Json::Value sorted hits;
138
     size t c = 0;
139
     for (auto it : copy) {
140
       if (c \le max) { // max wird von search.js auf 200 gesetzt
141
         it ["no"] = static cast < Json :: UInt > (c); // Hit-Nummer
142
             anpassen
         sorted_hits.append(it); // Reduce Hits to max size
143
144
       c++;
145
146
     // request bekommt sortierten Vector
147
     request.json ["hits"] = sorted hits;
149
150
151
153 double
wf::RankedSearch::tokenDistance(std::vector<Token>::
      const iterator &curr it,
                                  \mathrm{std}::\mathrm{vector} < \mathrm{Token} > ::\mathrm{const} iterator &
                                      foll it) const
    // Calculate distance between the two matches
     double distance = 0.0;
157
     auto token 1 = *curr it;
158
     auto token_2 = *foll_it;
159
     double pos1 = static cast < double > (token 1.pos()); // wegen
160
        size t
     double pos2 = static cast < double > (token 2.pos());
161
     distance = pos2 - pos1;
162
     return distance;
163
164
165
166
167 std::unordered map<const wf::dict::Entry *, double>
  wf::RankedSearch::split(std::string query, int &query_size) const
169
     std::unordered\_map{<}const\_wf::dict::Entry~*,~double{>}
170
        query entries;
     std::istringstream splitter(query);
171
     std::string word;
172
     while (splitter >> word) {
173
       query size++;
       std::wstring wide str = wf::util::widen(word.c str());
175
       //std::wcout << wide_str << std::endl;
176
       const wf::dict::Entry * token_entry = dict_->find(wide str);
177
```

```
if (token entry) {
178
         if (query entries.find(token entry) != query entries.end())
179
            query entries [token entry] += (double) 1.0;
180
         } else {
181
           query_entries[token_entry] = (double) 1.0;
182
           //std::wcout << "Query: "<< token entry->str<< std::endl;
183
         }
184
       }
186
     return query_entries;
187
188
189
190 /
191 void
192 wf::RankedSearch::fillHitMap(std::unordered map<const wf::dict::
      Entry *, double > &query map,
                         std::unordered map<const wf::dict::Entry *,
193
                            double > &hit_map,
                         const wf::dict::Entry * &match entry) const
194
195
196
       for (auto query pair : query map) {
197
         auto query_entry = query_pair.first;
198
         size t vPartFound = match\_entry -> gc.str().find(L"+V+#");
199
         if (vPartFound != std::string::npos) {
200
           hit map[query entry] += (double) 3.0;
201
202
         if (query_entry == match_entry){
203
           hit map | query entry | += (double) 1.0 / query map |
204
               query entry ;
           //std::wcout << "Literal: " << hit map[query entry] <<
205
               std::endl;
           break;
206
         } else if ((*query entry).isLemmaOf(match entry)) {
207
           hit map [query entry] += (double) 0.9 / query map [
               query_entry |;
           //std::wcout << "OTHER LEMMA!" << hit map[query entry] <<
209
               std::endl;
           break;
210
         } else if ((*match entry).isInverseLemmaOf(query entry)){
211
           //std::wcout << "inverse" << std::endl;
212
           hit_map[query_entry] += (double)0.9 / query_map[
               query entry];
           break;
214
215
216
           //std::wcout << "no match" << std::endl;
           continue;
218
219
       }
220
```

```
221
222
223 void
224 wf::RankedSearch::fill unord map(std::unordered map<std::string,
      std::pair<double, double>> & m,
                         std::string docAbsatz,
225
                         double s_rank) const
226
227
     auto it = m. find ( docAbsatz );
228
     if (it == m.end()) 
229
       auto\ value\_pair\ =\ std::make\_pair\ (s\_rank\ ,1.0\ )\ ;\ //\ 1.0\ Anzahl
230
           des Absatzes
       m. insert (std::make_pair(docAbsatz, value_pair));
231
232
     else {
233
     m[docAbsatz].first += s_rank;
234
     m[docAbsatz].second += 1.0; // öErhhe anzahl der äAbstze
235
236
237
238
  std::string
  wf::RankedSearch::getDocumentAbsatz(std::string siglum) const{
     std::smatch match;
241
     std::regex rgx("(.*)_");
242
     if (std::regex search(siglum, match, rgx)){
243
       return match [1];
     } else {
^{245}
       return siglum;
246
247
248
```

### 7.2 Suchergebnisse

Suchergebnisse zu "Gedanken Welt"

```
{{
1
   "DocAbsatz": "Ts-212,1028[1]",
2
   3
   4
   "matches": [
5
    {"token": ["Welt", "Gedanken"],
6
      "pos": [ 7, 10 ]
7
     }]
8
 }, {
9
   "DocAbsatz": "Ts-211,399[2]",
10
   11
   12
   "matches": [
13
    {"token": ["Welt", "Gedanken"],
14
      "pos": [ 7, 10]
15
16
17
 }, {
```

```
"DocAbsatz": "Ts-213,370 \, \text{r} \, [2] \, \text{et} \, 369 \, \text{v} \, [2]",
18
     "ab_rank": 2,
19
     "s rank": 2,
20
     "matches": [
^{21}
       {"token": ["Welt", "Gedanken"],
22
          "pos": [ 8, 13 ]
23
        }]
24
25
     "DocAbsatz": "Ts-212,1062[1]",
26
     "ab rank": 0.9570512820512821,
27
     "s rank": 0,
28
     "matches":
29
       {"token": ["Gedanken", "Gedanken"],
30
          "pos": [ 9, 28 ]
31
        } |
32
33
     "DocAbsatz": "Ts-212,1062[1]",
34
     "ab rank": 0.9570512820512821,
35
     "s rank": 0.25,
36
     "matches": [
37
       {"token": ["Gedanken"],
          "pos": [ 4 ]
39
        }
40
   },{
41
     "DocAbsatz": "Ts-212,1062[1]",
42
     "ab\_rank":\ 0.9570512820512821\,,
43
     "s rank": 0.25,
44
     "matches": [
45
       { "token": [ "Gedanken" ],
46
          "pos": [53]
47
48
49
     "DocAbsatz": "Ts-212,1062[1]",
50
     "ab rank": 0.9570512820512821,
51
     "s\_rank": -0.007692307692307699,
52
     "matches": [
53
       { "token": ["Gedanken", "Gedanke"],
54
          "pos": [17,30]
        }]
56
   },{
57
     "DocAbsatz": "Ts-212,1062[1]",
58
     "ab rank": 0.9570512820512821,
59
     "s rank": 5
60
     "matches": [
61
       {\text{"token": [ "Welt", "Gedanken"]}}
62
          "pos": [ 12,14 ]
63
        } |
64
65
     "DocAbsatz": "Ts-212,1062[1]",
66
     "ab_rank": 0.9570512820512821,
     "s rank": 0.25,
68
     "matches": [
69
       { "token": ["Gedanken"],
70
```

```
71 "pos": [5]
72 }]
```

#### Suchergebnisse zu "man weiß"

```
{
1
     "DocAbsatz": "Ts-213,149v[2]",
2
     "ab_rank": 10,
3
     "s_rank": 10,
4
     "matches": [
5
       \{ \text{"token": [ "ßwei","man"]}, 
6
          "pos": [1, 2]
       }
8
   },{
9
     "DocAbsatz": "Ms-142,105[2] et106[1]",
10
     "ab rank": 10,
11
     "s\_rank": 10,
12
     "matches": [
13
       {"token": ["man", "ßwei"],
14
          "pos": [0, 1]
1.5
       } ]
16
   },{
17
     "DocAbsatz": "Ts-211,656[4]et657[1]",
18
     "ab rank": 10,
19
     "s_rank": 10,
^{20}
     "matches": [
21
       {"token": ["ßwei", "man"],
22
          "pos": [8, 9]
23
        }]
^{24}
   },{
25
     "DocAbsatz": "Ts-212,1894[1] et1895[1]",
26
     "ab_rank": 10,
27
     "s_rank": 10,
28
     "matches":
^{29}
       {"token": ["ßwei", "man"],
30
          "pos": [8,9]
31
       }
^{32}
33
34
     "DocAbsatz": "Ts-211,498[2]",
35
     "ab rank": 5.58,
^{36}
     "s_rank": 1.16,
^{37}
     "matches": [
38
         "token": ["man", "wissen", "man"],
39
          "pos": [2,4,7]
40
       } ]
41
42
     "DocAbsatz": "Ts-213,55r[1]",
43
     "ab_rank": 5.58,
44
     "s\_rank": 10,
45
     "matches": [
46
       {"token": ["ßwei", "man"],
47
```

```
"pos": [3,4]
48
       }
49
50
     "DocAbsatz": "Ts-213,55r[1]",
51
     "ab rank": 5.58,
52
     "s rank": 1.16,
53
     "matches": [
54
       {"token": ["man", "wissen", "man"],
55
         "pos": [2,4,7]
       }
57
58
     "DocAbsatz": "Ts-211,498[2]",
59
     "ab rank": 5.58,
60
     "s rank": 10,
61
     "matches": [
62
       {"token": ["ßwei", "man"],
         "pos": [3,4]
64
       }]
65
66
     "DocAbsatz": "Ts-212,187[2]",
67
     "ab rank": 5.58,
     "s rank": 1.16,
69
     "matches": [
70
       {"token": ["man", "wissen", "man"],
71
         "pos": [2,4,7]
72
  },{
74
     "DocAbsatz": "Ts-212,187[2]",
75
     "ab rank": 5.58,
76
     "s rank": 10,
77
     "matches": [
78
       {"token": ["ßwei", "man"],
         "pos": [3,4]
80
       }]
81
82
```

#### Suchergebnisse zu "Haus Stadt"

```
{
1
   "DocAbsatz": "Ts-213,727r[2]et728r[1]",
2
   3
   "matches": [
5
    {"token": ["Haus", "Stadt"],
6
      "pos": [22,25]
7
    }
8
 },{
9
   "DocAbsatz": "Ts-211,630[2]et631[1]",
10
   11
   12
   "matches": [
13
    {"token": ["Haus", "Stadt"],
14
      "pos": [23,26]
15
```

```
}]
16
  },{
17
    "DocAbsatz": "Ts-212,1766[1] et1767[1]",
18
    19
    20
    "matches": [
21
       {"token": ["Haus", "Stadt"],
22
         "pos": [22,25]
23
       }]
24
  } ,{
^{25}
    "DocAbsatz": "Ts-213,185r[8] et 186r[1]",
26
    "ab_rank": 2.5,
27
    "s rank": 2.5,
28
    "matches": [
^{29}
       {"token": ["Haus", "Stadt"],
30
         "pos": [26,30]
31
       }]
32
  },{
33
    "DocAbsatz": "Ts-212,571[2]",
34
    "ab rank": 2.5,
35
    "s rank": 2.5,
36
    "matches": [
37
       {"token": ["Haus", "Stadt"],
38
         "pos": [24,28]
39
       }]
40
41
    "DocAbsatz": "Ts-211,297[3]",
42
    "ab_rank": 2.5,
43
    "s rank": 2.5,
44
    "matches": [
45
       {"token": ["Haus", "Stadt"],
46
         "pos": [26,30]
47
       }]
49
    "DocAbsatz": "Ms-142,12[2] et13[1]",
50
    "ab_rank": 0.44625,
51
    "s rank": 0.58,
52
    "matches": [
53
       {"token": ["äHusern", "Stadt", "Stadt"],
54
         "pos": [4,11,14]
55
       }]
56
  },{
57
    " DocAbsatz": "Ms-142,12[2]et13[1]",
58
    "ab rank": 0.44625,
59
    "s_rank": 0.3125,
60
    "matches":
61
       {"token": ["Stadt", "äHusern", "äHusern", "äHusern", "äHusern", "
62
          äHusern"|,
         "pos": [8,35,50,66,69,96]
63
       }]
64
  },{
65
    "DocAbsatz": "Ts-211,579[2] et580[1]",
66
     "ab_rank": 0.25,
67
```

```
"s_rank": 0.25,
68
     "matches": [
69
       {"token": ["Haus"],
         "pos": [6]
       }]
72
73
     "DocAbsatz": "Ts-211,580[2]et581[1]",
74
     "ab rank": 0.25,
75
     "s rank": 0.25,
     "matches": [
77
       {"token": ["Haus"],
78
         "pos": [7]
79
80
       } |
```

#### 7.3 READMEs aus Gitlab

#### 7.3.1 Sphinx

```
# Prerequisites:
2 * Sphinx: `sudo pip3 install -U Sphinx`
3 * recommonmark `sudo pip3 install -U recommonmark`
4 * Sphinx-Markdown-Tables: `pip install sphinx-markdown-tables`
5
6 # Create HTML-Page:
7
8 1. `cd PATH_TO_FOLDER/`
9 2. `make html` (Warnings can be ignored for now.)
10 3. `firefox build/html/index.html`
```

#### 7.3.2 wf

```
1 There are two simple tools for wittfind:
2
   * the wf tool searches for queries and gives out a list of
3
      matches.
   * the wf display tool uses the input of wf to display the
      matches in the original file.
   * the wf server tool starts a server listening on a unix domain
      socket.
   * the wf client tool queries a server that waits on a unix
      domain socket.
   * all tools know the -h [ --help] option.
9 All tools are built in the build/bin folder
10 (if you use the custom Makfile, they are copied directly into
     your current directory).
11
12
      $ ./wf —help
      $ ./wf -d dictionary -f input -q "query" -m max -o outfile
13
```

```
14
       $ ./wf -d dictionary -f input -Q query-file -m max -o outfile
       $ ./wf -L dicdictionary -f input -q "query"
15
16 max specifies the maximal number of hits shown. Default is 25, if
       0 all matches are shown.
17
18
       $ ./wf display —help
       $ ./wf_display [hits]
19
       $ ./wf -d data/witt_WAB_dela.txt -f data/witt_input_tagged.
20
          xml -q / d+/ | ./wf_display
       $ ./wf -v -d data/witt WAB dela.txt -f data/witt input tagged
21
          .xml -q /d+/ | ./wf display -v
22
23 creating wf display xml output for the web front end
24 (-t B \text{ marks hits with } < B > \text{hit} < /B >, -r \text{ hits puts the results in } <
      hits > ... < / hits > tags):
25
      $ ./wf display -t B -r hits pfad/zur/serverdatei.xml
26
27 starting the server:
       $ cd wf/build/
28
29
       $ ./bin/wf server — dictionary ../data/witt WAB dela XI.txt
         -- files ... / data / files.txt
30
31 using the client:
32
33
       $ ./bin/wf client —query '[VVFIN] & denken' -f ../data/Ts
          -213 NORM-tagged.xml \
34
      --max 10 --threads 2
35
36 ## Query syntax
  * path/to/graph.json(arg1, arg2, ..., argn) loads a subgraph
      from path/to/graph.json
38
     and replaces the arguments $1$, $2$, ... $n$ with arg1, arg2,
         \dots, argn.
39
      It is not possible to concatenate subgraph expressions with
         other expressions in a node.
   * token matches any token that is either 'token' or has 'token'
40
       as its lemma.
   * "token" or 'token' matches any token that is equal to token.
41
   * \<GC\> matches any token with the grammtical code 'GC'.
42
43
   * [TAG] matches any token with an annotation equal to 'TAG'
44
   * /regex/ matches any token that matches the regular expression
       'regex'.
45
     Note: if you want the regex to match the whole token you have
         to use '/^regex$/'.
   * /regex/i matches any token that matches the regular expression
46
        'regex' ignoring case.
   * [/regex/] matches all token whose tag matches the regular
47
       expression 'regex'.
   * You can prefix any query expression with '!' to prohibt the
48
       higlighting of this particular match.
     E.g. "/\w+/!<PUNCT> /\w+/" matches words, followed by
49
         punctuation and another word,
```

- 50 but the punctuation is not highlighted as a match (the two words are, though):
- 51 [[[foo]]] , [[[bar]]]
- 52 \* You use boolean operators '(', ')', '&', '|' or '~' in a node to form complex expressions:
- \* '/en\$/ & [N] ' matches token that end on 'en' and have the tag N.
- $^{54}$  \*  $^{\prime}/en\$/$  | [N]  $^{\prime}$  matches token that end on 'en' or have the tag N.
- \* '~ /en\$/' matches token that \_don't\_ end with 'en'.
- \* use bracets to form more complex expressions.
- \* the parser for complex expressions is not finished yet. You need to use explicit whitespace
- to seperate operators and expressions:  $'(/a/|/b/)\&^{\sim}/c/'$  is invalid.
- 59 Use: '(/a/|/b/) & ~/c/'.
- 60 \* There are some unitex special expressions: \<MAJ\>, \<MIN\>, \<MOT\>, \<PRE\> and \<NB\>
- wich all map to speical regex pattern.
- 62 \* It is possible to use '?' to make the previous expression optional.
- The query 'a b? c' matches the token 'a b c' or 'a c'.
- 64 \* It is possible to specify a range expression '{f, t}' to match the previous match f to t times.
- The query 'x  $\{3, 5\}$ ' matches at least 3 xs up to 5 xs.
- The Expression  $\{x\}$  is shorthand for  $\{x, x\}$ .
- $67\ *\ You\ can\ append\ '+'\ to\ an\ expression\ to\ make\ it\ match\ one\ or\ more\ times\,.$
- The Query a <MOT>+ b matches any expression 'a' followed by one or more words and a 'b'.
- 69 \* You can append '\*' to an expression to make it match zero or more times.

## Literaturverzeichnis

- [1] CPP-Reference. Algorithm std::sort. Verfügbar unter https://en.cppreference.com/w/cpp/algorithm/sort [Accessed: 01-Jun-2019].
- [2] CPP-Reference. C++ named requirements: Compare. Verfügbar unter https://en.cppreference.com/w/cpp/named\_req/Compare [Accessed: 01-Jun-2019].
- [3] Dr. Max Hadersbeck, Daniel Bruder. WAST Wittgenstein Advanced Search Tools Dokumentation. CIS, LMU, München, 2014.
- [4] Bharat Krishna. Ranking search results by reranking the results based on local interconnectivity. U.S. Patent No. 6.526.440., 25 Feb. 2003.
- [5] Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze.Introduction to information retrieval. Cambridge University Press. 2008.
- [6] Friedhelm Padberg, Sebastian Wartha. Didaktik der Bruchrechnung. Herder, Freiburg, 1978.
- [7] Juan Ramos. Using tf-idf to determine word relevance in document queries. In: Proceedings of the first instructional conference on machine learning. Vol. 242. 2003, [S. 133-142].
- [8] Prof. Dr. Klaus U. Schulz. Formale Sprachen und Automaten. CIS, LMU, München, 2011
- [9] Sphinx-Doc. Projects using Sphinx. verfügbar unter http://www.sphinx-doc.org/en/master/examples.html [Accessed: 01-Jun-2019].
- [10] Bjarne Stroustrup. The C++ programming language, Pearson Education India, 2000.
- [11] Wikipedia. Sphinx (Software). Verfügbar unter https://de.wikipedia.org/wiki/Sphinx\_(Software) [Accessed: 01-Jun-2019].

# Abbildungsverzeichnis

2.1	WAST Dokumentationswebsite	6
	wf Schaubild	
4.1	Suchergebnis zu Partikelverb	20
5.2	Suchergebnisse WiTTFind	30

# Inhalt des beigelegten USB-Sticks

Folgende Dateien beinhalten neuen Code:

- RankedSearch.cpp/.hpp
- Searcher.cpp/.hpp
- Message.cpp/.hpp
- $\bullet \ SearchDocument.cpp/.hpp$
- Server\_test.cpp
- CMakeLists.txt

Folgende Dateien beinhalten nur neue Kommentare:

- $\bullet$  Search.hpp
- Dictionary.cpp./.hpp

Sonstige Dateien:

• BA.pdf