

3. Analyse von Korpusdaten

3.1 Beschreibungsebenen

Um das sprachliche Material in einem Korpus klassifizieren, analysieren und interpretieren zu können, ist es sinnvoll, sich vorab über die verschiedenen Beschreibungsebenen klar zu werden. Zu unterscheiden sind dabei nicht nur die verschiedenen korpuslinguistischen Kategorien der Textwörter, Tokens und Types, wichtig ist auch die Abgrenzung der korpuslinguistischen Begriffe Textwort, Wortform-Type und Lemma-Type von den linguistischen Beschreibungskategorien Wort, Wortform und Lexem. Betrachten wir dazu erst einmal ein Beispiel (vgl. Linke *et al.* 2004). Wie viele Wörter enthält der folgende Satz:

(5) Wenn hinter Fliegen eine Fliege fliegt, fliegt eine Fliege Fliegen hinterher.

Hier gibt es mehrere richtige Antworten: elf, sieben oder sechs. Versteht man Wörter als Einheiten der Schriftsprache (orthografische Wörter), so muss die Antwort elf lauten, denn der Satz enthält elf Einheiten, die durch Leerzeichen oder Satzzeichen voneinander getrennt sind:

(6) 1. Wenn, 2. hinter, 3. Fliegen, 4. eine, 5. Fliege, 6. fliegt, 7. fliegt, 8. eine, 9. Fliege, 10. Fliegen, 11. hinterher

In einem Korpus werden solche, mithilfe der Schreibung definierten Worteinheiten als **Textwörter**, Wortform-Tokens oder laufende Wortformen bezeichnet. Definiert man Wörter hingegen als syntaktische Wörter oder Wortformen, d.h. als formal voneinander unterscheidbare Bestandteile des Satzes, so kommt man nur noch auf sieben Wörter:

(7) 1. wenn, 2. hinter, 3. Fliegen, 4. eine, 5. Fliege, 6. fliegt, 7. hinterher

Die Formen *Fliegen*, *eine*, *Fliege* und *fliegt* kommen im Satz doppelt vor. Sie werden bei dieser Betrachtung zusammengefasst und zählen jeweils als eine Wortform, da sie sich in ihrer Form nicht unterscheiden. In einem Korpus werden Worteinheiten, die über Unterschiede in ihrer Form definiert werden, als Wortform-Types bezeichnet. Schließlich kann man sich fragen, ob es sinnvoll ist, die Flexionsformen *Fliege* und *Fliegen* als zwei getrennte Wörter zu behandeln. Der einzige Unterschied zwischen den beiden Wortfor-

men besteht darin, dass *Fliegen* eine zusätzliche grammatische Markierung für den Plural trägt, nämlich das Flexionssuffix *-n*, während die Wortform *Fliege* keine zusätzliche Markierung für den Singular hat. In beiden Fällen handelt es sich aber um Vertreter derselben Tierart. Anders wäre das etwa im Fall von (8).

(8) Wenn hinter Fliegen eine Biene fliegt, fliegt eine Biene Fliegen hinterher.

Abstrahiert man also von grammatischen Markierungen, so gehören sowohl die Wortform *Fliege* als auch die Wortform *Fliegen* zum selben Wort, nämlich dem Lexem oder morphologischen Wort *FLIEGE*. Zählt man nun also die Lexeme im Beispielsatz, so kommt man auf sechs Wörter: die Konjunktion *WENN*, die Adverbien *HIINTER* und *HINTERHER*, das Nomen *FLIEGE*, den Artikel *EINE* und das Verb *FLIEGEN*. In einem Korpus werden Worteinheiten, die sich nur in ihren Flexionsmerkmalen unterscheiden, unter dem Begriff Lemma oder Lemma-Type zusammengefasst. Um im Folgenden auch optisch zu verdeutlichen, wenn explizit von Lemma-Types bzw. Lexemen die Rede ist, werde ich diese im Text in Kapitälchen angeben (*FLIEGE*).

Die bisher eingeführten Begriffe Lemma-Type, Wortform-Type und Wortform-Token haben jedoch einen Nachteil: sie beziehen sich auf Worteinheiten. Korpusanalysen sind jedoch nicht auf die Untersuchung von Worteinheiten beschränkt, sondern können auch auf Laut-, Satz-, Text- oder Bedeutungsebene durchgeführt werden. Aus diesem Grund unterscheidet man in der Korpuslinguistik üblicherweise unabhängig von der Sprachebene zwischen Types und Tokens. Bei einem **Token** handelt es sich ganz allgemein um das konkrete Vorkommen einer sprachlichen Einheit im Korpus. Das kann eine bestimmte Wortform, Lautäußerung oder Phrase sein. Ein **Type** ist hingegen die abstrakte sprachliche Einheit, die zusammengehörige Tokens wie Wortformen oder Lautvarianten zusammenfasst und dabei von konkreten Merkmalsausprägungen wie Flexions- oder Intonationsmerkmalen abstrahiert.

Verdeutlichen wir uns die Unterscheidung zwischen Types und Tokens an einem konkreten Beispiel, nämlich der Untersuchung von Wortbildungswandel im Mainzer Zeitungskorpus (Scherer 2005). Ziel der Untersuchung war es festzustellen, ob und wie sich die Möglichkeiten zur Bildung von Wörtern in den letzten vier Jahrhunderten verändert haben. Dazu wurden Nomen wie *Lehrer*, *Schüler* und *Lacher* untersucht, die mithilfe des Wortbildungssuffixes *-er* aus Verben (*lehren*, *lachen*), Nomen (*Schule*) und anderen Wortarten abgeleitet wurden. Der folgende Ausschnitt stammt aus dem ersten von insgesamt neun Teilkorpora.

Den 12. diß hat der Herr von Venesi zu Brüffel ein Pancket vnd Dantz gehalten / vnd seinem **Diener** oder **Gertner** befohlen / niemand ohn seinen willen in den Saal zulaffen / gleich hernach ist ein Spannifcher Hauptman vnd **Ritter** *Don Roderico Floris* genant / aber nur schlecht / wie ein **Diener** bekleidet für den Saal kommen / vnd sich mit gewalt eindringen wollen / dem **Gertner** zwey Maultaschen etliche Gemechtlöß / vnd ein Stich in Arm geben / darauff der **Diener** im Zorn seinen Dolchen außgezogen / vnd den Hauptman alßbald erfrochen / welchen man hernach in des Spinnola Hauß geführt / vnd den **Gertner** einziehen lassen / vnd die Freud eingestellt worden / als nun den dritten Tag / der Hauptman von allen *Officirn* vnd Herrn von Hoff ins Closter zu den **Augustinern** zur begrebnis begleitet / ist entzwischen der **Theter** zum Galgen geführt / gehenckt vnd ihme die rechte Hand abgehawen / vnd an Galgen genagelt worden / das gedünckt jederman ein frembder Sententz sein / weil der **Theter** seinem Befehlig nachkommen / vnd sich auch der Nothwehr gebrauchen müffen.

Abbildung 3: Mainzer Zeitungskorpus: Teilkorpus 1609 *Aviso* (Ausschnitt)

Analysieren wir den Ausschnitt zunächst im Hinblick auf die Bildung von Personenbezeichnungen mit dem Suffix *-er*. Der Ausschnitt enthält insgesamt zehn Tokens (vgl. 9a), die im Text fett markiert sind. Diese zehn Tokens verteilen sich auf insgesamt fünf Types (vgl. 9b).

- (9) a. Diener, Gertner, Ritter, Diener, Gertner, Diener, Gertner, Augustinern, Theter, Theter
 b. AUGUSTINER (1 Token), DIENER (3 Tokens), GÄRTNER (3 Tokens), RITTER (1 Token), TÄTER (2 Tokens)

Untersuchen wir den Ausschnitt hingegen im Hinblick auf eine andere Fragestellung, etwa die Verwendung von Präpositionen, so kommen wir zu einem anderen Ergebnis. Der Ausschnitt enthält in diesem Fall 16 Tokens (vgl. 10a), die im Text unterstrichen sind. Sie verteilen sich auf sieben Types (vgl. 10b).

- (10) a. von, zu, ohn, in, für, mit, in, im, in, von, von, ins, zu, zur, zum, an
 b. AN (1 Token), IN (5 Tokens), MIT (1 Token), OHNE (1 Token), VON (3 Tokens), VOR (= für) (1 Token), ZU (4 Tokens)

Aufgabe 10: Wie viele Textwörter enthält der oben stehende Ausschnitt aus dem Mainzer Zeitungskorpus? Bitte ermitteln Sie die Zahl der Types (Lemma-Types) und Tokens für die im Ausschnitt enthaltenen Nomen. Zählen Sie Eigennamen zu den Nomen.

Wie gesagt befasst sich aber nicht jede korpuslinguistische Analyse mit einer Untersuchung auf Wortebene. Vielmehr gibt es Fragestellungen, die sprachliche Einheiten betreffen, die größer oder kleiner sind als das Wort. So erfolgt die bereits erwähnte Untersuchung von Elter (2005) zur Kasusverwendung bei *wegen* auf der

syntaktischen Ebene der Phrase. Hier werden die Begriffe Type und Token auf die im Korpus enthaltenen *wegen*-Phrasen angewendet. Um Tokens handelt es sich also etwa bei den Phrasen *wegen des Mondscheinfrisierens* oder *wegen dem starken Wind* (vgl. Kapitel 1.2). Diese Tokens verteilen sich in Elters Studie auf zwei zugrunde liegende grammatische Muster, *wegen* + Genitiv und *wegen* + Dativ. Diese beiden Muster stellen die Types dar.

Untersucht man hingegen wie Dittmar/Bressem (2005) die Verbstellung in Nebensätzen mit *weil*, so bezieht sich die Zahl der Tokens auf die Anzahl der *weil*-Nebensätze und die der Types auf die beiden Verbstellungsvarianten: finites Verb an letzter Stelle wie in (11a) bzw. finites Verb an zweiter Stelle wie in (11b).

- (11) a. weil ich das immer so **mache**
 b. weil das **mache** ich immer so

Festzuhalten sind demnach zwei Dinge: Zum einen kann sich die Grundgesamtheit, auf die sich die Begriffe Type und Token beziehen, je nach Fragestellung verändern. Dahingegen ist die Anzahl der Textwörter in einem Korpus unabhängig von der untersuchten Fragestellung. Zum anderen haben die Begriffe Type und Token nicht zwangsläufig einen Bezug zur Wortebene. Vielmehr können die Begriffe auf sprachliche Einheiten unterschiedlicher Ebenen wie Wort, Satz, Text bzw. deren Bestandteile verweisen.

3.2 Methoden

Korpora können sowohl qualitativ als auch quantitativ ausgewertet werden. Der wesentliche Unterschied zwischen qualitativen und quantitativen Korpusanalysen besteht nicht darin, welche Fragestellungen untersucht werden, sondern wie diese untersucht werden. Nehmen wir das Beispiel Fremdwörter. In zwei unterschiedlichen Studien analysieren Schanke (2001) und O'Halloran (2002) den Einfluss englischer und französischer Entlehnungen, so genannter Anglizismen und Gallizismen, im Deutschen. Beide arbeiten mit einem selbst zusammengestellten Korpus aus Zeitungen bzw. Zeitschriften. Schankes Korpus enthält sämtliche Ausgaben des *Handelsblatts* aus dem März 2000, O'Hallorans Korpus umfasst ein Teilkorpus zur Modesprache mit mehreren Jahrgängen der Frauenzeitschrift *Brigitte* sowie ein Teilkorpus zur Standardsprache, das mehrere Jahrgänge des Nachrichtenmagazins *Stern* und der *Berliner Illustrierten Zeitung* enthält.

Abgesehen von der Größe der Korpora unterscheiden sich die beiden Studien auch in ihrer Methode. Während Schanke in seiner Korpusanalyse einen qualitativen Ansatz wählt, untersucht O'Halloran ihr Korpus unter quantitativen Gesichtspunkten.

Schanke's Ziel ist es, die gefundenen Fremdwörter im Hinblick auf ihre Wortart zu klassifizieren und sie bestimmten Themenbereichen wie Computerbranche, Börse oder Bankwesen zuzuordnen. Bei Schanke's Untersuchung geht es also darum, in einem Korpus die Existenz bestimmter sprachlicher Erscheinungen, nämlich Anglizismen, festzustellen, die einzelnen Wörter herauszusuchen und sie nach bestimmten Kriterien, konkret nach Wortfeldern, zu klassifizieren. Schanke's Vorgehen entspricht dem einer **qualitativen Korpusanalyse**. Qualitative Korpusanalysen legen ihren Schwerpunkt auf die Ermittlung, die Klassifizierung, die Einordnung und Interpretation von bestimmten Phänomenen.

Im Gegensatz dazu steht O'Halloran's Arbeit. O'Halloran untersucht die Verbreitung von englischen und französischen Fremdwörtern innerhalb der letzten einhundert Jahre. Dabei stellt sie fest, dass der Anteil an Fremdwort-Types im Gesamtkorpus steigt, und zwar von 0,6% im Jahr 1902 auf 2,0% im Jahr 1997. Darüber hinaus beobachtet sie, dass der Anteil von Fremdwort-Tokens im Teilkorpus zur Modesprache zu jedem Zeitpunkt den Fremdwortanteil im Teilkorpus zur Standardsprache übersteigt. Im Jahr 1997 liegt der Fremdwortanteil in der Standardsprache z.B. bei 4,0%, in der Modesprache hingegen bei 14%. O'Halloran geht es in ihrer Untersuchung also darum, die **Frequenz** von bestimmten Phänomenen zu ermitteln und miteinander zu vergleichen, um daraus Rückschlüsse über die untersuchte Fragestellung ziehen zu können. Das Bestimmen von Häufigkeiten im Korpus und die sich daraus ergebende Möglichkeit, Ergebnisse unmittelbar miteinander zu vergleichen, ist das Kennzeichen **quantitativer Korpusuntersuchungen**.

An quantitativen Kennzahlen wird standardmäßig die Korpusgröße ermittelt, die üblicherweise in Textwörtern gemessen wird (vgl. Kapitel 4.8). Sie bildet die wichtigste Bezugsgröße für alle quantitativen Auswertungen. Ist die Größe eines Korpus nicht bekannt, sind quantitative Analysen nur dann sinnvoll, wenn die Ergebnisse für mehrere ähnlich geartete Phänomene innerhalb des Korpus verglichen werden können.

Von besonderem Interesse ist für den Forscher die Anzahl der Types und Tokens des untersuchten Phänomens, da diese beiden Zahlen Auskunft darüber geben, wie oft ein Phänomen insgesamt im Korpus belegt ist (Tokens) und auf wie viele unterschiedliche

Ausprägungen des Phänomens (Types) sich die Tokens verteilen. So fanden sich im obigen Ausschnitt aus dem Mainzer Zeitungskorpus insgesamt fünf Types und zehn Tokens für die untersuchten nominalen *-er*-Derivate (vgl. Abbildung 3).

Wichtig ist, die Zahl der Types und Tokens jeweils im Verhältnis zur Korpusgröße zu sehen. Zehn Tokens in einem kleinen Korpus können relativ gesehen eine höhere Frequenz darstellen als hundert Tokens in einem großen Korpus. Liegt die Zahl der Types und Tokens vor, kann man daraus das Verhältnis von Types zu Tokens berechnen. Dieses **Type-Token-Verhältnis** gibt Auskunft darüber, wie viele Tokens durchschnittlich auf einen Type entfallen. Liegt die Anzahl der Tokens je Type sehr hoch, handelt es sich bei den meisten Tokens vermutlich um häufig verwendete Ausdrücke, die eine gewisse Formelhaftigkeit aufweisen. Die Anzahl der spontanen, neuen Formen, die dem untersuchten Muster folgen, ist dann gering. Umgekehrt ist ein niedriges Verhältnis von Types zu Token ein Indiz dafür, dass viele Types nur selten vorkommen. Kommt ein Type nur ein einziges Mal im Korpus vor, spricht man von einem **Hapax Legomenon**. Enthält ein Korpus viele Hapax Legomena und andere seltene Types, ist die Wahrscheinlichkeit hoch, dass das untersuchte sprachliche Muster von den Sprechern bzw. Schreibern produktiv eingesetzt wird und dass nach seinem Vorbild neue Bildungen vorgenommen werden. Allgemein kann der Anteil der Hapax Legomena an der Zahl der Tokens dazu verwendet werden, die **Produktivität** eines sprachlichen Musters zu bestimmen. Je höher der Anteil der Einmalbelege, desto höher ist die Wahrscheinlichkeit, dass das Muster Neubildungen hervorbringt.

Zielt eine Korpusuntersuchung weniger auf die Wortebene als auf die Satzebene ab, so ist es sinnvoll, die Anzahl der Sätze in einem Korpus oder in einem Text zu ermitteln. Zudem können die durchschnittliche Satzlänge, die Zahl der Sätze mit einer bestimmten Länge sowie der Anteil von Sätzen mit einer bestimmten Zahl an Wörtern dazu dienen, syntaktische Charakteristika eines Korpus, eines Textes oder einer Varietät herauszuarbeiten.

Um zu gewährleisten, dass es sich bei den ermittelten Ergebnissen nicht um bloßen Zufall handelt, empfiehlt es sich, die Ergebnisse statistisch abzusichern. Dies geschieht mittels eines Signifikanztests, der sicherstellt, dass die Ergebnisse nicht allein dem Zufall geschuldet sind. Die meisten Signifikanztests gehören jedoch der höheren Mathematik an, sodass es sich empfiehlt, entsprechende Statistikprogramme zu benutzen.

3.3 Vergleichbarkeit von Daten

Beim Vergleich von Daten aus unterschiedlichen Korpora ist es wichtig, qualitative und quantitative Charakteristika der Korpora zu beachten. Zum einen sollte überlegt werden, ob sich die Korpora aufgrund ihrer Konzeption überhaupt vergleichen lassen und wenn ja, in welchem Rahmen. Auf den ersten Blick scheint es wenig sinnvoll zu sein, ein Korpus der Kindersprache mit einem historischen Korpus oder einem Korpus zur Fachsprache der Biologie zu vergleichen. Dennoch kann ein solcher Vergleich sinnvoll sein, wenn untersucht werden soll, ob Parallelen in der kindlichen und der historischen Sprachentwicklung bestehen oder ob Biologiebücher für die Schule den Entwicklungsstand der Kinder angemessen berücksichtigen, was die Bezeichnung von Pflanzen, Tieren und deren Teilen betrifft.

Als Vergleichsgrundlage dienen häufig Referenzkorpora, die als Standard verwendet werden, um Abweichungen zwischen Varietäten und Standardsprache festzustellen (vgl. Kapitel 2.9). Ein Referenzkorpus kann also dazu dienen, festzustellen, inwieweit sich das Bairische, die Fachsprache der Medizin oder das Mittelhochdeutsche von der Standardsprache unterscheiden. Daneben ist es aber auch wichtig, auf die quantitative Vergleichbarkeit zu achten. Tabelle 1 zeigt die Ergebnisse einer Suchabfrage in drei verschiedenen Korpora des IDS, dem Bonner Zeitungskorpus, dem LIMAS-Korpus und dem Mannheimer Korpus 1. Gesucht wurde nach den Wortformen *Buch*, *Hochhaus* und *Universität*.

Suchbegriff	Bonner Zeitungskorpus	LIMAS-Korpus	Mannheimer Korpus 1
<i>Buch</i>	313	166	229
<i>Hochhaus</i>	16	3	6
<i>Universität</i>	315	116	219

Tabelle 1: Ergebnisse der Stichwortsuche (absolut)

Wie man sieht, finden sich im Bonner Zeitungskorpus jeweils die meisten und im LIMAS-Korpus jeweils die wenigsten Belege für alle drei Suchbegriffe. Woran liegt das? Nun, zum einen könnte es an der Zusammensetzung der Korpora liegen: Während das Bonner Zeitungskorpus ausschließlich Zeitungstexte enthält, ist der Anteil an Zeitungstexten in den anderen beiden Korpora gering. Sie bestehen überwiegend aus Textsorten wie Belletristik, Gebrauchsliteratur und wissenschaftlichen Texten. Eine mögliche Folgerung ist also, dass sich alle drei Suchbegriffe überdurchschnittlich häufig in Zei-

tungstexten finden. Bevor man jedoch einen solchen Schluss zieht, sollte man einen Blick auf die Größe der Korpora werfen, die verglichen werden sollen.

Aufgabe 11: In einem Korpus A finden sich 80 Belege für das Wort *BLUMENTOPF*, in Korpus B 100 Belege für dasselbe Wort. Korpus A und Korpus B enthalten je eine Million Textwörter. Korpus C enthält ebenfalls 100 Belege für *BLUMENTOPF*, aber anderthalb Millionen Textwörter. In welchem der drei Korpora finden sich die meisten Belege für das Wort *BLUMENTOPF*?

Obwohl sich im Bonner Zeitungskorpus mehr Belege für die Wortform *Buch* finden als in den anderen beiden Korpora, bedeutet dies nicht unbedingt, dass *Buch* im Bonner Zeitungskorpus häufiger ist als im LIMAS-Korpus oder dem Mannheimer Korpus 1. Dies liegt daran, dass eine Aussage über die Häufigkeit eines Wortes immer im Verhältnis zur Größe des Korpus gesehen werden muss.

Ein direkter Vergleich von Korpusdaten ist aufgrund unterschiedlicher Korpusgröße im Normalfall nicht möglich. Insofern ist es beim Vergleich von Daten aus verschiedenen Korpora von größter Wichtigkeit, die jeweiligen Ergebnisse ins Verhältnis zur Korpusgröße zu setzen. Geht es darum, die Frequenz bestimmter Wörter anzugeben, so kann dies wie bei O'Halloran (2002) in Form von Prozentangaben geschehen, die sich auf die Zahl der Textwörter oder Lemma-Types im Korpus beziehen. Befasst sich die Untersuchung hingegen nicht mit Einheiten der Wortebene, so kommt nur eine **Normalisierung** infrage. Bei der Normalisierung werden die Ergebnisse auf eine bestimmte Anzahl von Textwörtern, etwa 10.000 oder eine Million, umgerechnet. Dabei sollte sich die Normalisierung an der typischen Textlänge im Korpus orientieren.

Vergleichen wir das Bonner Zeitungskorpus mit dem LIMAS-Korpus und dem Mannheimer Korpus 1, so ergibt sich folgendes Bild: Das Bonner Zeitungskorpus enthält über 3,6 Millionen Textwörter und ist damit fast dreimal so groß wie das LIMAS-Korpus mit rund 1,2 Millionen Textwörtern und etwa anderthalbmal so groß wie das Mannheimer Korpus 1, das rund 2,6 Millionen Textwörter beinhaltet. Es ist also nicht verwunderlich, dass sich im größten Korpus die meisten Belege finden und im kleinsten Korpus die wenigsten! Als Konsequenz aus diesen Größenunterschieden müssen sämtliche Ergebnisse auf eine genormte Korpusgröße umgerechnet werden (vgl. Tabelle 2). Sinnvoll erscheint in diesem Fall eine Normgröße von einer Million Textwörtern.

Suchbegriff	Bonner Zeitungskorpus	LIMAS-Korpus	Mannheimer Korpus 1
<i>Buch</i>	86	135	89
<i>Hochhaus</i>	4	2	2
<i>Universität</i>	87	94	85

Tabelle 2: Ergebnisse der Stichwortsuche (normalisiert je 1 Mio. Textwörter)

Die normalisierten Daten in Tabelle 2 zeigen im Vergleich zu Tabelle 1 ein ganz anderes Bild: Lediglich die Wortform *Hochhaus* kommt mit vier Belegen je Million Textwörter im Bonner Korpus häufiger vor als in den anderen beiden Korpora. Dahingegen finden sich die meisten Belege für *Buch* (135) und *Universität* (94) im LIMAS-Korpus. In den anderen beiden Korpora haben *Buch* und *Universität* hingegen fast dieselbe Frequenz.

Nach dem Vergleich der normalisierten Ergebnisse lautet die Frage also nicht mehr, warum die Wortformen *Buch*, *Hochhaus* und *Universität* im Bonner Zeitungskorpus am häufigsten sind, sondern vielmehr, warum *Buch* im LIMAS-Korpus deutlich öfter vorkommt als in den anderen Korpora, warum *Hochhaus* im Bonner Zeitungskorpus doppelt so oft belegt ist wie in den anderen Korpora und warum *Universität* in allen drei Korpora fast gleich häufig vorkommt.

Bei dem Vergleich von Daten aus verschiedenen Korpora sind demnach zwei Dinge wichtig: Zum einen sollte man sich fragen, ob es im Hinblick auf die zu untersuchende Fragestellung überhaupt sinnvoll ist, zwei gegebene Korpora miteinander zu vergleichen. Zum anderen ist es unerlässlich, bei einem Vergleich von korpusbasierten Häufigkeiten die Korpusgröße zu berücksichtigen, da bei unterschiedlich großen Korpora andernfalls die Ergebnisse des Vergleichs verfälscht werden.

Aufgabe 12: Unten finden Sie die Ergebnisse aus dem Erlanger Dürer-Korpus (Müller 1993) und dem Würzburger Korpus der Wissensliteratur (Brendel *et al.* 1997) zur Wortbildung in frühneuhochdeutschen Fachtexten.

In welchem der beiden Korpora finden sich die meisten Nominalisierungen mit den Suffixen *-er*, *-heit/-keit* und *-ung* (Types, Tokens)? Bitte berechnen Sie zudem das Type-Token-Verhältnis für die einzelnen Suffixe und vergleichen Sie die Ergebnisse miteinander.

Korpus	Textwörter	<i>-er</i>		<i>-heit/-keit</i>		<i>-ung</i>	
		Types	Tokens	Types	Tokens	Types	Tokens
Dürer-Korpus	440.000	93	700	76	326	193	2.443
Wissensliteratur	1.073.000	510	4.505	454	6.575	1.025	5.213

3.4 Stichwortsuche – die Suche nach Wörtern, Wortformen und Wortteilen

Die einfachste Möglichkeit an Informationen in einem Korpus zu kommen, ist die Suche nach einem bestimmten Wort, einer Wortform oder einem Wortteil wie *HAUS*, *liest* oder *un-*. Genau diese Möglichkeit der Stichwortsuche haben Günther (2002) und Hämmer (2001) genutzt, um die Verwendung des Wortes *stolz* bzw. des Wortteils *-park* zu analysieren.

Anlass für Günthers Untersuchung des Wortes *stolz* war die öffentliche Diskussion um den Satz *Ich bin stolz darauf, ein Deutscher zu sein*, den ein Politiker in einem Interview geäußert hatte. Günther wollte jenseits der gesellschaftlichen Debatten klären, in welchem Zusammenhang das Wort *stolz* verwendet wird. *Stolz sein*, so ergab Günthers Recherche in den Textkorpora des IDS, kann man nicht nur auf eine Leistung (vgl. 12a), eine berufliche oder private Tätigkeit (vgl. 12b), sondern auch auf eine bestimmte nationale oder geografische Herkunft (vgl. 12c).

- (12) a. stolz darauf, Abgeordneter geworden zu sein
 b. stolz darauf, ein Bauer/ein Zeitungsleser zu sein
 c. stolz darauf, ein Schweizer/ein Münchner zu sein

Wie Günther feststellt, bringt die Äußerung *stolz darauf, ein X zu sein* somit zwar ein gewisses Maß an Selbstbewusstsein zum Ausdruck, sie muss jedoch nicht zwangsläufig ein Zeichen von Überheblichkeit seitens des Sprechers sein.

Anders als Günther suchte Hämmer nicht nach vollständigen Wörtern, sondern lediglich nach einem Wortbestandteil. Gegenstand ihrer Analyse bildeten Komposita mit dem Zweitglied *-park*. Hämmer's Suche im Korpus des Projekts Deutscher Wortschatz in Leipzig ergab, dass die Komposita mit *-park* semantisch in zwei Gruppen zerfallen. Bei der größeren Gruppe handelt es sich um klassische Determinativkomposita, bei denen *-park* als Grundwort auftritt (vgl. 13a).

- (13) a. Schlosspark, Tierpark, Vergnügungspark
 b. Gerätepark, Unternehmenspark, Windpark

Ein *Schlosspark*, *Tierpark* oder *Vergnügungspark* ist eine bestimmte Art von Park, die durch das Erstglied näher bestimmt wird: ein Park am Schloss, ein Park mit Tieren, ein Park, den man zur Vergnügung besucht. Daneben fand Hämmer aber auch Beispiele wie in (13b), wo *-park* im Sinne von 'Ansammlung, Gesamtheit von X' interpretiert werden muss. Ein *Gerätepark* ist nicht ein Park