

Semi-automatische Differenzanalyse von komplexen Textvarianten

André Gießler `andre.giessler@informatik.uni-halle.de`

Marcus Pöckelmann `marcus.poeckelmann@informatik.uni-halle.de`

Jörg Ritter `joerg.ritter@informatik.uni-halle.de`

Einer der Schwerpunkte von Projekten der Editionsphilologie ist die Untersuchung alter Texte mit Mehrfachüberlieferungen sowie die Textgenese. Dabei stellt sich für die beteiligten Wissenschaftler die Aufgabe, die Verbindungen zwischen den einzelnen Textvarianten herauszuarbeiten und dabei Gemeinsamkeiten und Unterschiede zu erkennen. Oft sind große Textmengen der verschiedenen Varianten zu sichten, einander zuzuordnen und detailliert zu vergleichen, um anschließend als Edition präsentiert werden zu können. Bisher erfolgen die bei der Edition anfallenden, teils sehr gleichartigen und zeitaufwändigen Zwischenschritte in Handarbeit, belegen damit wertvolle Arbeitszeit und setzen einen Gesamtüberblick über das Textmaterial voraus. Die Informationstechnologie bietet heute Möglichkeiten, mit denen die Durchführung vieler dieser Schritte zumindest teilautomatisiert werden kann. Den Geisteswissenschaftlern können Werkzeuge zur Verfügung gestellt werden, die ihnen die Arbeit nicht nur wesentlich erleichtern, sondern auch die Fehleranfälligkeit reduzieren und neue Formen der Auswertung eröffnen.

Das hier vorgestellte, vom BMBF geförderte Gemeinschaftsprojekt von Geisteswissenschaftlern und Informatikern mit dem Ziel, Werkzeuge und Methoden zum Textvergleich und zur Erstellung kritischer und genetischer Editionen zu entwickeln. Diese Methoden sollen generisch und damit auf viele Textformen anwendbar sein. Dazu werden zwei Repräsentanten verschiedenartiger Textformen mit ihren Überlieferungen zur Grundlage genommen, für deren Anforderungen und Eigenheiten jeweils zugeschnittene Verfahren entwickelt und evaluiert werden. Dabei werden die Prozesse von der Erkennung und Lemmatisierung der Wörter, das Auffinden sich entsprechender Textstellen, die Herausarbeitung der Unterschiede und Gemeinsamkeiten, bis hin zur Darstellung in einer genetischen Edition abgedeckt. Im späteren Verlauf des Projektes findet die Verallgemeinerung der gewonnenen Erkenntnisse für möglichst viele Textformen statt.

Ein Teil des Projektes betrachtet einen handschriftlichen Lehrbuchtext aus der Zeit des Spätmittelalters, der in Frühneuhochdeutsch verfasst wurde. Unter der Leitung von Hans-Joachim Solms wird die „Wundarzney“ des Heinrich von Pfalzpaint aus dem 15. Jahrhundert in ihren zehn verfügbaren Überlieferungen untersucht. Ziel der Altgermanisten ist hier, die Varianten zu vergleichen und in einer kritischen Edition und einer Online-Edition mit Synopse und Variantenapparat darzustellen. Ausgangspunkt sind die Handschriften, die in einem ersten Arbeitsschritt von den Geisteswissenschaftlern diplomatisch transkribiert wurden. Da diese Texte in der Sprachstufe Frühneuhochdeutsch verfasst wur-

den, gibt es keine einheitliche Graphie, wodurch dieselben Wörter in verschiedenen Überlieferungen deutlich unterschiedlich geschrieben werden und mit jeder weiteren Überlieferung neue Schreibweisen entdeckt werden. Ein Beispiel ist das Wort Pfeil, welches in den Schreibweisen „pfeil“, „pffeil“, „pfejl“, „pffejl“, „pfeyl“, „pffeyl“ auftritt. Bevor ein Textvergleich stattfinden kann, ist somit eine philologische Aufbereitung nötig, bei der die Wörter erkannt und normalisiert werden. Die Normalisierung wird mit der Lemmatisierung jedes einzelnen Wortes erreicht. Für eine möglichst präzise Abbildung einer Handschrift auf eine andere werden die Wörter zusätzlich noch mit Part-Of-Speech-Tags und morphologischen Attributen wie Kasus, Numerus und Genus versehen. Für die Aufgabe der Lemmatisierung und Annotation existieren bereits verschiedene automatisierte Ansätze, die allerdings nicht auf Handschriften aus dem Frühneuhochdeutschen anwendbar sind, da sie nur eine sehr geringe Toleranz für abweichende Schreibweise (oder Schreibfehler) von Wörtern aufweisen. Die unstetige Graphie in den Handschriften der „Wundarzney“ führt bei ihnen zu geringen Trefferquoten in Bezug auf die Korrektheit der von ihnen vorgeschlagenen Annotationen.

Im Projekt wurde das Werkzeug *Lemmano* entwickelt, das einen semiautomatischen Ansatz verfolgt. Es erlaubt die manuelle Annotierung jedes einzelnen Wortes, in dem es zu dem Wort ähnliche Wortformen ableitet, diese mit zugehörigen Annotationen in Lexika sucht und dem Benutzer so Vorschläge für das aktuelle Wort unterbreitet. Ähnlich heißt hier, dass sich die neue Wortform mittels von Altgermanisten erarbeiteten Ersetzungsregeln, die auf Äquivalenzen bestimmter Buchstabenfolgen basieren, aus dem gegebenen Wort ableiten lässt. Im Gegensatz zu automatischen Ansätzen liegt die Entscheidung für das passende Lemma und die passenden morphologischen Daten beim Anwender.

Die Benutzeroberfläche ist intuitiv verständlich gestaltet und auf die Massenverarbeitung ausgerichtet. *Lemmano* hat sich bei der Annotation der Handschriften der „Wundarzney“ durch Germanisten als sehr große Arbeitserleichterung erwiesen. Da es ein webbasiertes Werkzeug ist, können mehrere Nutzer gleichzeitig annotieren und profitieren von den gelernten Eingaben der anderen Nutzer. Bild 1 zeigt den Dialog für eine Wortform.

Nach der Lemmatisierung der Handschriften lassen sich diese nun mit weiteren digitalen Werkzeugen, die im Rahmen des Projektes realisiert werden, detailliert vergleichen. Der nächste Schritt auf dem Weg zu einer Edition mit synoptischer Darstellung ist die Alignierung der Varianten zueinander. Dabei werden sich entsprechende Textstellen identifiziert, gegenübergestellt und mit Auflistung der Unterschiede in einem Variantenapparat präsentiert.

Das zweite Teilprojekt widmet sich einem neuphilologischen Text in Fremdsprache. Unter der Leitung von Thomas Bremer wird die im späten 18. Jahrhundert in Französisch verfasste „Histoire philosophique et politique des établissements et du commerce des Européens dans les deux indes“ von Abbé Guillaume Thomas François Raynal untersucht. Sie gilt auf Grund ihrer Thematik, die Auseinandersetzung mit der europäischen Kolonialpolitik dieser Zeit, als bedeutendes Werk der Aufklärung. Nach dem Verbot der Erstauflage von 1770 erschienen zwei weitere Auflagen sowie eine postume Textfassung. Im Rahmen des Projektes soll eine genetische Edition der Lateinamerika-Bände entstehen, die insbesondere durch ihre Interaktivität die Evolution dieses Werkes nachvollziehbar macht. Dabei liegt das Hauptaugenmerk auf einer abschnittsweisen

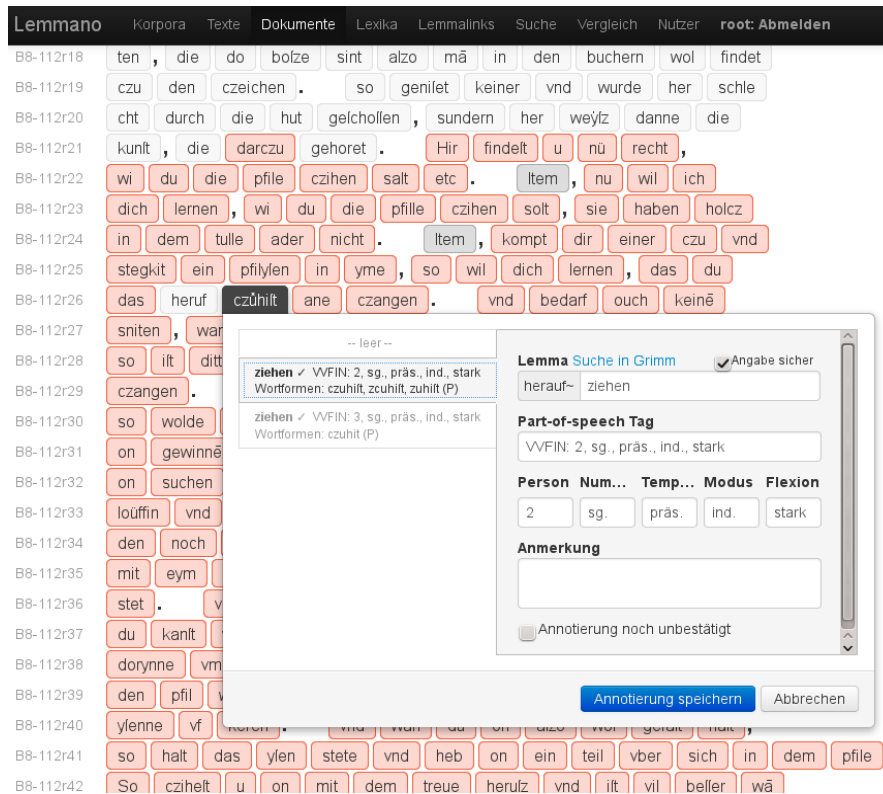


Abbildung 1: Annotationsdialog für eine Wortform in Lemmano

Gegenüberstellung der vier Varianten als Fließtext mit einer übersichtlichen Form des Apparats.

Ausgangspunkt hier sind digitale Faksimiles, die als Scan der Erstausgaben angefertigt wurden. Diese wurden mittels existierender Software zur Texterkennung in eine maschinenlesbare Form gebracht, anschließend von den Romanisten von fehlerhaft erkannten Stellen bereinigt sowie mit speziellen Markierungen versehen, die beispielsweise Überschriften oder Seitennummern kenntlich machen. Erleichtert wird die Suche nach Fehlern dabei durch ein softwaregestütztes Verfahren, das auffällige Kombinationen von Symbolen anzeigt. So liefert beispielsweise die Suche nach Kombinationen aus Buchstaben und Zahlen ohne trennendes Leerzeichen eine Reihe von fehlerhaft kodierten Jahreszahlen, wie „15o6“ oder „158o“. Aus den so entstehenden Textdokumenten wird automatisch eine TEI-konforme XML-Darstellung generiert, die als Grundlage für die folgenden Arbeitsschritte dient. Für einen der beiden algorithmischen Schwerpunkte aus Sicht der Informatik, die Alignierung der Absätze, werden derzeit verschiedene automatische Verfahren geprüft. Für den zweiten Schwerpunkt, die Bestimmung und Visualisierung der Unterschiede auf Absatzebene, wurde bereits ein erster Ansatz implementiert. Dieser vergleicht eine beliebige Anzahl von Textpassagen untereinander. Auf Basis der Levenshtein-Distanz werden die Differenzen zwischen den Varianten ermittelt und daraus eine synoptische Dar-

stellung in LaTeX erzeugt (Abbildung 2). Die Philologen können für die Visuali-

1770	1774	1780	1820
Les ¹ ministres de cette ² princesse prirent d'a- bord pour un ³ visionnaire un homme qui vou- loit ⁴ découvrir un monde. Ils le traitèrent long-temps avec cette hauteur ⁵ infultante que les hommes communs, quand ils font en place, ont pour les hommes de ⁶ génie. Colomb ne fut pas rebuté par les difficultés. Il ⁷ avoit comme tous ceux qui forment des projets ex- traordinaires, ⁸ et enthousiasme ⁹ qui les roidit contre les jugemens de l'ignorance, ¹⁰ les dédains de l'orgueil, les ¹¹ petiteffes ¹² de l'ava- rice, les délais de la ¹³ pareffe. ¹⁴ Son ¹⁵ ame ferme, élevée, courageuse, fa ¹⁶ prudence et son adresse ¹⁷ le firent enfin triompher de tous ces ¹⁸ obstacles. On lui accorda trois petits vaiffeaux ¹⁹ et quatre-vingt-dix ²⁰ hommes. Il partit ²¹ le 3 Août 1492, ²² avec le titre d'Ami- ral ²³ et de Vice-Roi ²⁴ des îles, ²⁵ des terres qu'il découvrirait. ²⁶	Les ¹ ministres de cette ² princesse prirent d'a- bord pour un ³ visionnaire un homme qui vou- loit ⁴ découvrir un monde. Ils le traitèrent long-temps avec cette hauteur ⁵ infultante que les hommes en place affectent si souvent avec ceux qui n'ont que du ⁶ génie. Colomb ne fut pas rebuté par les difficultés. Il ⁷ avoit comme tous ceux qui forment des projets ex- traordinaires, ⁸ et enthousiasme ⁹ qui les roidit contre les jugemens de l'ignorance, ¹⁰ les dédains de l'orgueil, les ¹¹ petiteffes ¹² de l'ava- rice, les délais de la ¹³ pareffe. ¹⁴ Son ¹⁵ ame ferme, élevée, courageuse, fa ¹⁶ prudence et son adresse ¹⁷ le firent enfin triompher de tous les ¹⁸ obstacles. On lui accorda trois petits vaiffeaux ¹⁹ et quatre-vingt-dix ²⁰ hommes. Il partit ²¹ le 3 Août 1492, ²² avec le titre d'Ami- ral ²³ et de vice-roi ²⁴ des îles et ²⁵ des terres qu'il découvrirait. ²⁶	Les ¹ ministres de cette ² princesse prirent d'a- bord pour un ³ visionnaire un homme qui vou- loit ⁴ découvrir un monde. Ils le traitèrent long-temps avec cette hauteur ⁵ infultante que les hommes en place affectent si souvent avec ceux qui n'ont que du ⁶ génie. Colomb ne fut pas rebuté par les difficultés. Il ⁷ avoit comme tous ceux qui forment des projets ex- traordinaires, ⁸ et enthousiasme ⁹ qui les roidit contre les jugemens de l'ignorance, ¹⁰ les dédains de l'orgueil, les ¹¹ petiteffes ¹² de l'ava- rice, les délais de la ¹³ pareffe. ¹⁴ Son ¹⁵ ame ferme, élevée, courageuse, fa ¹⁶ prudence et son adresse ¹⁷ le firent enfin triompher de tous les ¹⁸ obstacles. On lui accorda trois petits vaiffeaux ¹⁹ et quatre-vingt-dix ²⁰ hommes. Sur cette foit ²¹ il étoit, dont l'armement ne cou- toit pas cent mille francs, il mit à la voile ²² le 3 Août 1492, ²³ avec le titre d'Amiral ²⁴ et de vice-roi ²⁵ des îles et ²⁶ des terres qu'il découvrirait, et arriva aux Canaries où il s'é- toit proposé de relâcher. ²⁷	Les ¹ ministres de cette ² princesse prirent d'a- bord pour un ³ visionnaire un homme qui vou- loit ⁴ découvrir un monde. Ils le traitèrent long-temps avec cette hauteur ⁵ infultante que les hommes en place affectent si souvent avec ceux qui n'ont que du ⁶ génie. Colomb ne fut pas rebuté par les ⁷ difficultés. Il ⁸ avoit comme tous ceux qui forment des projets ex- traordinaires, ⁹ et enthousiasme ¹⁰ qui les roidit contre les jugemens de l'igno- rance, ¹¹ les dédains de l'orgueil, les ¹² petit- esses ¹³ de l'avarice, les délais de la ¹⁴ pareffe. ¹⁵ Son ame ¹⁶ ferme, élevée, courageuse, fa ¹⁷ pru- dence et son adresse ¹⁸ le firent enfin triom- pher de tous les ¹⁹ obstacles. On lui accor- da trois petits navires ²⁰ et quatre-vingt-dix hommes. Sur cette faible escadre, dont l'ar- mement ne coûtait pas cent mille francs, il mit à la voile ²¹ le 3 août 1492, ²² avec le titre d'Amiral ²³ et de vice-roi ²⁴ des îles et des terres ²⁵ qu'il découvrirait, et arriva aux Canaries où il s'était proposé de relâcher. ²⁶

Abbildung 2: Automatisch generierte Synopse mit Variantenapparat für einen Absatz des franz. Textes

sierung der Textvarianten zwischen konfigurierbaren Darstellungsarten wählen. Die genannten Arbeitsschritte, von der Generierung der XML-Dateien bis hin zum Entwurf der elektronischen Edition, sollen perspektivisch in einer gemeinsamen, webbasierten Arbeitsumgebung eingebettet werden.

Prototypische Werkzeuge, unter anderem *Lemmano*, werden in den nächsten Monaten zur Demonstration als Webanwendungen öffentlich verfügbar gemacht.

Anmerkungen

Diese Arbeit wurde durch das Bundesministerium für Bildung und Forschung (BMBF) [Projektkürzel: 01UG1247 / human-325-010 / SaDA] im Rahmen des Projekts „Semi-automatische Differenzanalyse von komplexen Textvarianten“ unter Leitung von Prof. Dr. Thomas Bremer, Prof. Dr. Paul Molitor, Dr. Jörg Ritter und Prof. Dr. Hans-Joachim Solms gefördert. An dieser Stelle möchten wir auch unseren Projektmitarbeiterinnen Sylwia Kösser, Dr. Aletta Leipold und Susanne Schütz danken.