

Semi-automatische Differenzanalyse von komplexen Textvarianten

André Medek (né Gießler) andre.medek@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 „Wundarznei“ 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 „Wundarznei“ führt bei ihnen zu geringen Trefferquoten in Bezug auf die Korrektheit der von ihnen vorgeschlagenen Annotationen.

Im Projekt wurde das Werkzeug *LAKomp* 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. *LAKomp* hat sich bei der Annotation der Handschriften der „Wundarznei“ 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 diese sich nun innerhalb des Werkzeugs *LAKomp* detailliert vergleichen. Der Textvergleich auf dem Weg zu einer kritischen Edition erfolgt zweistufig: Zuerst werden in einem Vergleich auf Makroebene sich entsprechende Textstellen mit Hilfe von Signaturen identifiziert und einander gegenübergestellt. Das Ergebnis dieses Schrittes ist eine Alignierung. Im nächsten Schritt werden die innerhalb einer Alignierung einander zugeordneten Textstellen der verschiedenen Überlieferungen auf Mikroebene verglichen und durch 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

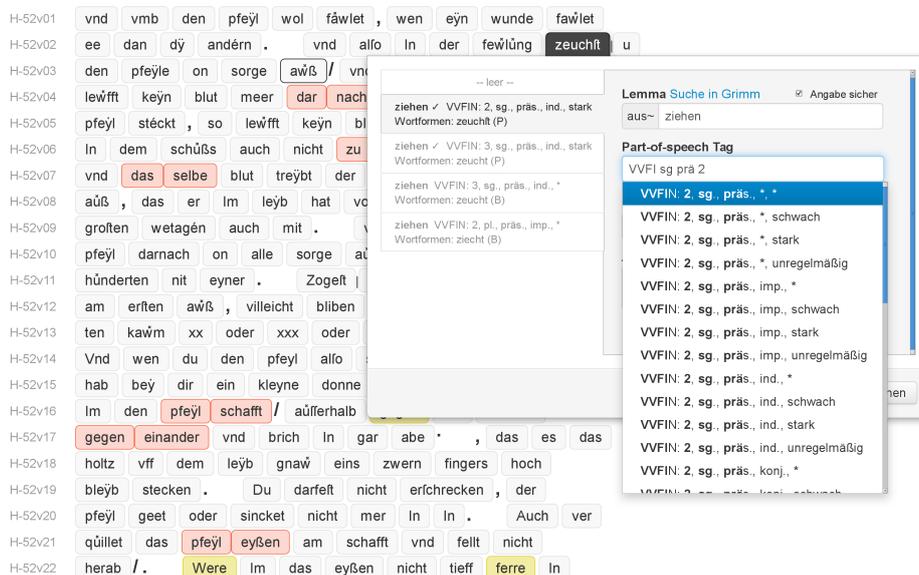


Abbildung 1: Annotationsdialog für eine Wortform in LAKomp

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

Ausgangspunkt hier sind digitale Faksimiles, die als Scan der Erstaussgaben 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 Darstellung in LaTeX erzeugt (Abbildung 2). Die Philologen können für die Visualisierung 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 gemein-

Les **ministres**¹ de cette **princesse**² prirent d'a-bord pour un **viñomairne**³ un homme qui vou-lait découvrir un monde. Ils le traitèrent long-temps avec cette hauteur infultante⁴ que les hommes communs, quand ils font en place, ont pour les hommes de génie. Colomb ne fut pas rebuté par les difficultés. Il **avait**⁵ comme tous ceux qui forment des projets ex-traordinaires⁶ et enthousiastes⁷ qui les roidit contre les jugemens de l'ignorance⁸; les dédains de l'orgueil, les **petiteffes**⁹ de l'avarice, les dédais de la **pareffe**¹⁰. Son **ame**¹¹ ferme, élevée, couragieuse, fa¹² prudence et son adresse¹³ le firent enfin triompher de tous ces obstacles¹⁴. On lui accorda trois petits **vaiffeaux**¹⁵ et quatre-vingt-dix¹⁶ hommes. Il **partit**¹⁷ le 3 Août 1492¹⁸ avec le titre d'Ami-rail¹⁹ et de Vice-Roi²⁰ des îles²¹; des terres qu'il découvrit²².

Les **ministres**¹ de cette **princesse**² prirent d'a-bord pour un **viñomairne**³ un homme qui voulait découvrir un monde. Ils le traitèrent long-temps avec cette hauteur infultante⁴ que les hommes **en place affectent** si souvent avec ceux qui n'ont que du génie. Colomb ne fut pas rebuté par les difficultés. Il **avait**⁵ comme tous ceux qui forment des projets ex-traordinaires⁶ et enthousiastes⁷ qui les roidit contre les jugemens de l'ignorance⁸; les dédains de l'orgueil, les **petiteffes**⁹ de l'avarice, les dédais de la **pareffe**¹⁰. Son **ame**¹¹ ferme, élevée, couragieuse, fa¹² prudence et son adresse¹³ le firent enfin triompher de tous les obstacles¹⁴. On lui accorda trois petits **vaiffeaux**¹⁵ et quatre-vingt-dix¹⁶ hommes. Il **partit**¹⁷ le 3 Août 1492¹⁸ avec le titre d'Ami-rail¹⁹ et de Vice-Roi²⁰ des îles²¹; des terres qu'il découvrit²².

Les **ministres**¹ de cette **princesse**² prirent d'a-bord pour un **viñomairne**³ un homme qui vou-lait découvrir un monde. Ils le traitèrent long-temps avec cette hauteur infultante⁴ que les hommes en place affectent si souvent avec ceux qui n'ont que du génie. Colomb ne fut pas rebuté par les difficultés. Il **avait**⁵ comme tous ceux qui forment des projets ex-traordinaires⁶ et enthousiastes⁷ qui les roidit contre les jugemens de l'ignorance⁸; les dédains de l'orgueil, les **petiteffes**⁹ de l'avarice, les dédais de la **pareffe**¹⁰. Son **ame**¹¹ ferme, élevée, couragieuse, fa¹² prudence et son adresse¹³ le firent enfin triompher de tous les obstacles¹⁴. On lui accorda trois petits **vaiffeaux**¹⁵ et quatre-vingt-dix¹⁶ hommes. **Sur cette foible escadre, dont l'armement ne coûtait pas cent mille francs, il mit à la voile**¹⁷ le 3 Août 1492¹⁸ avec le titre d'Amiral¹⁹ et de Vice-Roi²⁰ des îles et²¹ des terres qu'il découvrit²², et arriva aux Canaries où il s'é-tait proposé de relâcher²³.

Les **ministres**¹ de cette **princesse**² prirent d'a-bord pour un **viñomairne**³ un homme qui vou-lait découvrir un monde. Ils le traitèrent long-temps avec cette hauteur infultante⁴ que les hommes **en place affectent si souvent** avec ceux qui n'ont que du génie. Colomb ne fut pas rebuté par les ⁸ difficultés. Il **avait**⁵ comme tous ceux qui forment des projets extraordinaires⁶ et enthousiastes⁷ qui les roidit contre les jugemens de l'igno-rance⁸; les dédains de l'orgueil, les **petiteffes**⁹ de l'avarice, les dédais de la **pareffe**¹⁰. Son **ame**¹¹ ferme, élevée, couragieuse, fa¹² prudence et son adresse¹³ le firent enfin triompher de tous les obstacles¹⁴. On lui accorda trois petits navires¹⁵ et quatre-vingt-dix¹⁶ hommes. Sur cette **faible escadre, dont l'ar-mement ne coûtait pas cent mille francs, il mit à la voile**¹⁷ le 3 Août 1492¹⁸ avec le ¹⁹ titre d'Amiral¹⁹ et de Vice-Roi²⁰ des îles et²¹ des terres ²² qu'il découvrit²², et arriva aux Canaries ²³ où il s'était proposé de relâcher²³.

H70	ministres	H70	princesse	H70	viñomairne	H70	voulait	H70	traitèrent long-temps	H70	infultante	H70	communs, quand ils font en place, ont pour les hommes de												
1	H74	ministres	2	H74	princesse	3	H74	voulait	5	H74	infultante	7	H74	en place affectent si souvent avec ceux qui n'ont que du											
H80	ministres	H80	princesse	H80	viñomairne	H80	voulait	H80	traitèrent long-temps	H80	infultante	H80	en place affectent si souvent avec ceux qui n'ont que du												
H20	ministres	H20	princesse	H20	viñomairne	H20	voulait	H20	traitèrent long-temps	H20	infultante	H20	en place affectent si souvent avec ceux qui n'ont que du												
H70	H70	avait	H70	extraordinaires,	H70	enthousiastes	H70	l'ignorance,	H70	petiteffes	H70	pareffe.	H70	ame	H70	élevée, couragieuse, fa									
8	H74	9	H74	avait,	10	H74	extraordinaires,	11	H74	enthousiastes	12	H74	l'ignorance,	13	H74	petiteffes	14	H74	pareffe.	15	H74	ame	16	H74	élevée, couragieuse, fa
H80	H80	avait,	H80	extraordinaires,	H80	enthousiastes	H80	l'ignorance,	H80	petiteffes	H80	pareffe.	H80	ame	H80	élevée, couragieuse, fa									
H20	H20	avait,	H20	extraordinaires,	H20	enthousiastes	H20	l'ignorance,	H20	petiteffes	H20	pareffe.	H20	ame	H20	élevée, couragieuse, sa									
H70	son adresse	H70	ces obstacles.	H70	vaiffeaux,	H70	quatre-vingt-dix	H70	il partit																
H74	son adresse,	18	H74	les obstacles.	H74	vaiffeaux	H74	quatre-vingt-dix	21	H74	il partit														
H80	son adresse,	H80	les obstacles.	H80	navires	H80	quatre-vingt-dix	H80	Sur cette foible escadre, dont l'armement ne coûtait pas cent mille francs, il mit à la voile																
H70	Août 1492,	H70	d'Amiral	H70	Vice-Roi	H70	îles,	H70	découvrit.																
22	H74	Août 1492,	23	H74	vice-roi	24	H74	îles et	27	H74	découvrit.														
H80	Août 1492,	H80	d'Amiral	H80	vice-roi	H80	îles et	H80	découvrit, et arriva aux Canaries où il s'était proposé de relâcher.																
H20	Août 1492,	H20	d'Amiral	H20	vice-roi	H20	îles et	H20	découvrit, et arriva aux Canaries, où il s'était proposé de relâcher.																

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

samen, webbasierten Arbeitsumgebung eingebettet werden.

Prototypische Werkzeuge, unter anderem *LAKomp*, 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.

Weitere Informationen zum Projekt sind auf der Website <http://www.informatik.uni-halle.de/sada/> zu finden.