

Empirisch ermittelte Muster in Rechtschreibfehlern für die Automatisierung qualitativer Rechtschreibdiagnostik

Jan Thomas Röhrig | jan.roehrig@uni-due.de | Universität Duisburg-Essen

Ziel

- Ziel ist es, einen sprachdidaktischen Beitrag zur Entwicklung eines Programms zu leisten, das Rechtschreibfehler in freien Texten qualitativ analysieren kann.
- Hauptmotivation dafür ist der hohe Arbeitsaufwand, den qualitative Fehleranalysen erfordern¹.

Annahmen

- Rechtschreibfehler in frei geschriebenen Texten von Kindern treten nicht willkürlich, sondern systematisch auf.
- Diese Systematizität lässt sich in den Rechtschreibfehlern in Form von prototypischen Mustern nachweisen, die für die automatische Erkennung und Analyse von Rechtschreibfehlern genutzt werden können.

Methode

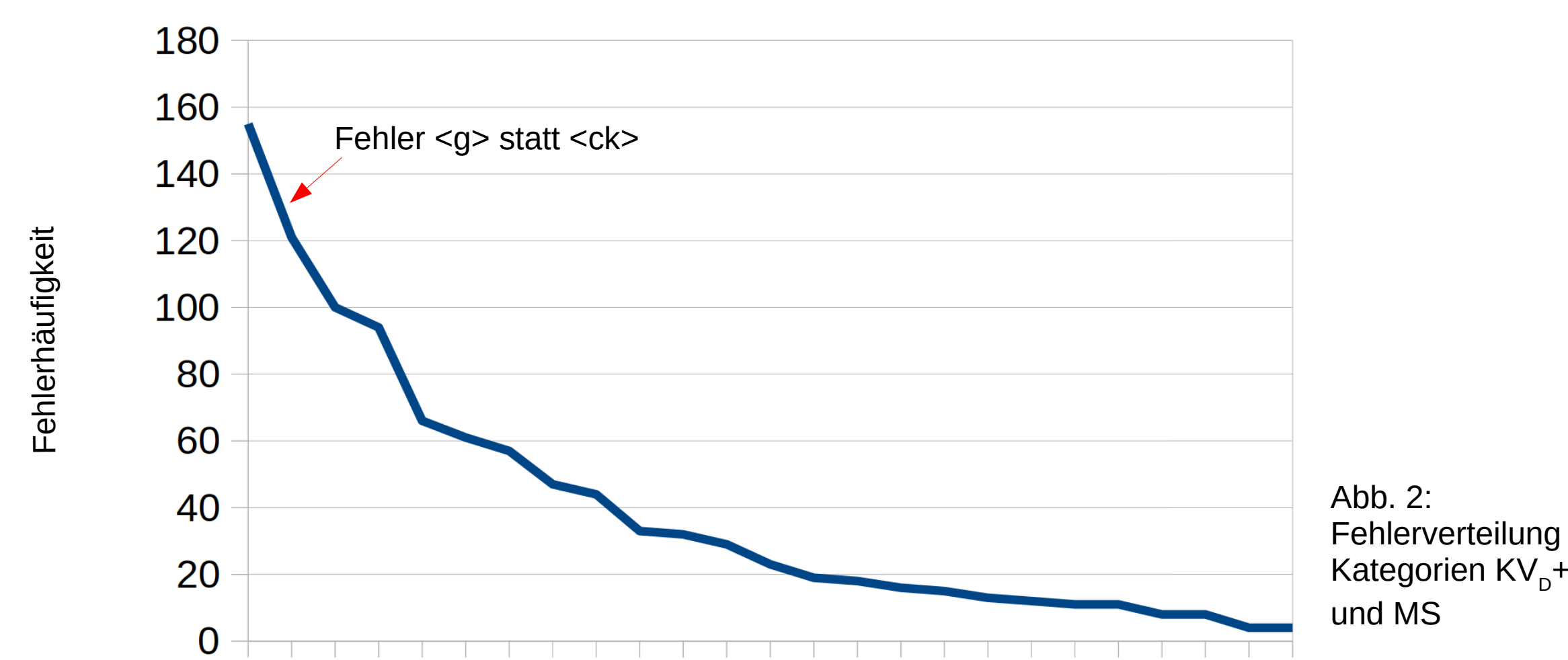
- Grundlage für die Untersuchung ist das LitKey-Korpus der Ruhr-Universität-Bochum, das etwa 1900 Kindertexte aus Grundschulen in NRW enthält².
- Alle Rechtschreibfehler des Korpus werden entlang der Fehlerkategorien der AFRA³ analysiert und kategorisiert.
- Die Annotationen werden auf Muster in den Fehlern hin verglichen.
- **Musterhaft sind Fehler, die**
 - (A) in mehreren Texten unterschiedlicher Kinder auftreten
 - (B) unter den gleichen Fehlergruppen der AFRA annotiert werden und
 - (C) signifikant häufiger auftreten als andere Fehler in der gleichen AFRA-Fehlerkategorie
- Ein Beispiel für ein solches Fehlermuster zeigt Abb. 1. Das Fehlermuster <g> statt <ck> nach Vokalgraphem und vor Flexionsendung tritt signifikant häufiger als andere Fehler der entsprechenden Kategorien im LitKey-Korpus auf.

Abb. 1

Anna verste~~g~~gt sich hinter dem Baum.

Fehlerbeschreibung: *<g> statt <ck> nach kurzem Vokalgraphem und vor konsonantischer Flexionsendung

- (A) kommt auch in anderen Texten im Korpus vor,
- (B) wird immer in den gleichen Fehlerkategorien annotiert (konkret: KV_D+ und MS),
- (C) und tritt signifikant häufiger als andere Fehler in KV_D+ und MS auf (siehe Abb. 2)

Abb. 2: Fehlerverteilung Kategorien KV_D+ und MS

Ablauf der automatischen Fehleranalyse

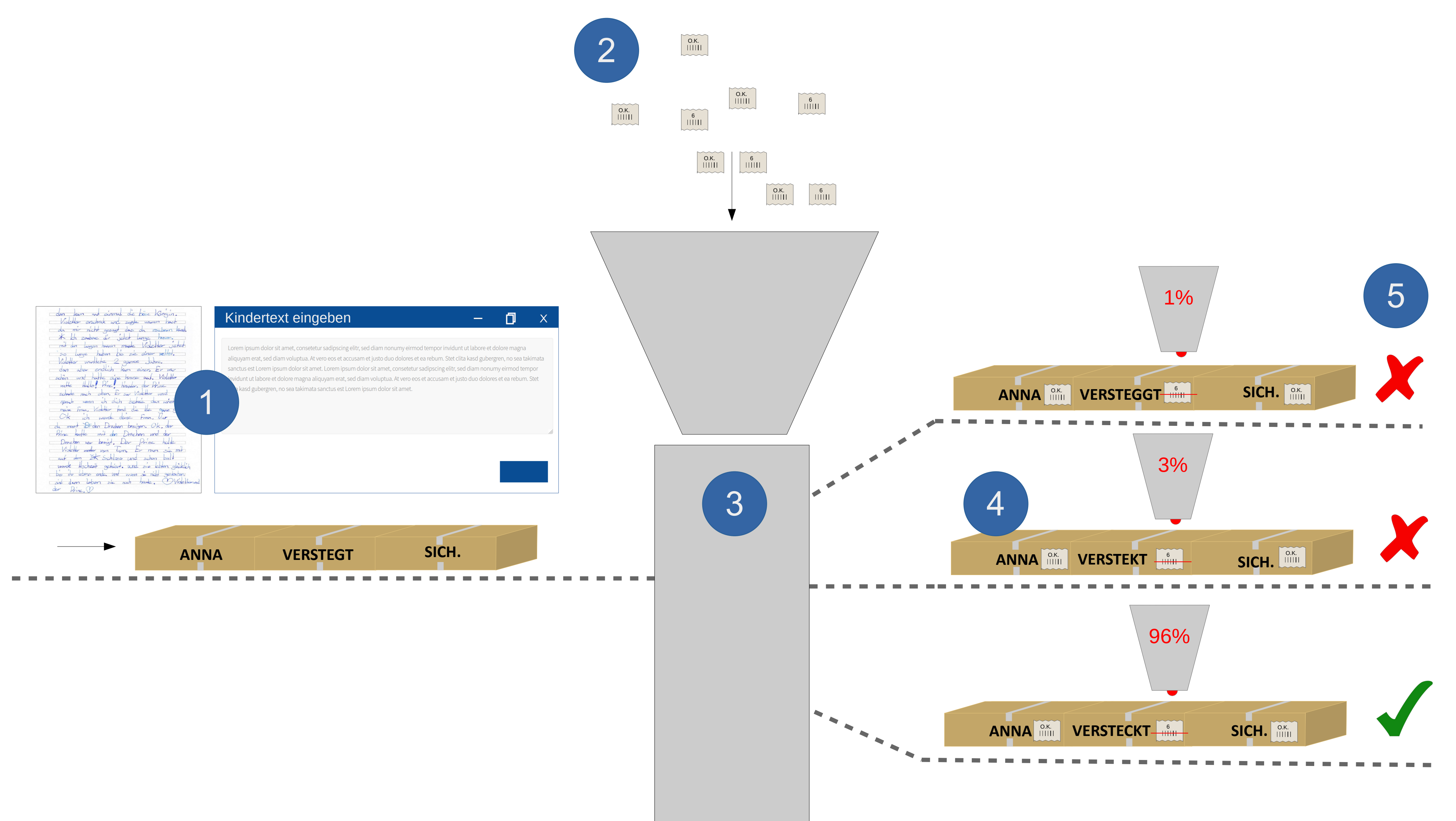


Abb. 3

Ergebnisse

- Die erwartete Systematizität in den Rechtschreibfehlern zeigt sich im Korpus.
- Die ermittelten Fehlermuster lassen sich in drei Gruppen unterscheiden:
 - (1) **morphemspezifisch**, d.h. Fehler treten wort- oder morphemspezifisch auf,
 - (2) **umgebungsspezifisch**, d.h. Fehler treten in einer bestimmten graphematischen Umgebung auf, und
 - (3) **unspezifisch**, d.h. Fehler treten unabhängig von einer bestimmten graphematisch- oder morphologischen Struktur auf.
- Die Fehlermuster machen es möglich, von einer fehlerhaften Schreibung durch „Suchen und Ersetzen“ automatisiert auf die orthografisch korrekte Version zu schließen (Abb. 3).

Ausblick

- Im Vordergrund steht die (weitere) technische Umsetzung des Verfahrens und die anschließende umfassende Evaluation, die dadurch erst möglich wird.
- Eine Erweiterung des Verfahrens um syntaktische Auswertungen ist wünschenswert, damit mehr Fehler in der automatischen Analyse erfasst werden (z.B. Groß- und Kleinschreibung).
- Die empirisch ermittelten Fehlermuster können auch für eine Weiter- und Neuentwicklung von qualitativen Fehleranalyseverfahren unabhängig von einer Automatisierung genutzt werden.

1 Der handgeschriebene Kindertext muss zuerst für die automatisierte Auswertung digitalisiert werden. Dafür wird er inkl. Rechtschreibfehler in die Programmoberfläche eingegeben.

3 Der eingegebene Kindertext wird nach den definierten Fehlermustern durchsucht und Fehlerstellen werden durch die korrekten Schreibungen aus den definierten Fehlermustern ersetzt (Suchen und Ersetzen). Jede Ersetzung wird mit dem jeweiligen Fehlermuster markiert.

5 Die Rechtschreibfehler werden mit der jeweils wahrscheinlichsten Korrektur korrigiert. Alternativ kann der Nutzer selbst die passende Korrektur auswählen. Die Markierung der Fehlermuster in der korrigierten Version erlaubt Rückschlüsse auf die Fehlerkategorie und macht damit die automatische qualitative Fehleranalyse neben der Korrektur möglich.

2 Grundlage für die Fehlererkennung sind die empirisch ermittelten Fehlermuster. Jedes Fehlermuster besteht aus

- einer Fehlerkategorie für die Analyse,
- der Häufigkeit des Fehlermusters,
- den Graphemen der Originalschreibung inkl. graphematischer Umgebung und
- den Graphemen der intendierten korrekten Schreibung.

4 Durch das Suchen und Ersetzen der Fehlermuster ergeben sich mehrere mögliche Korrekturen. Jede Korrektur ist durch die Markierung mit denjenigen Fehlermustern verknüpft, die zur Korrektur geführt haben. Aus den verschiedenen Möglichkeiten wird die wahrscheinlichste Korrektur ausgewählt. Die Wahrscheinlichkeit ergibt sich aus der Häufigkeit des Fehlermusters und kann um Wahrscheinlichkeitswerte aus N-Gramm-Sprachmodellen ergänzt werden.

Vorschau:

