

# Verjüngungskur für sprachliche Altdaten

tekom-Jahrestagung 2017  
24. bis 26. Oktober, Stuttgart

Axel Theofilidis  
axel.theofilidis@iailc.de  
IAI Linguistic Content AG, Saarbrücken

# Kommt Ihnen das bekannt vor?

■■■■■■■■■■ RUNGSVENTIL, MANO ■■■■■■■■■■  
■■■■■■■■■■ RUNG- VENTIL  
■■■■■■■■■■ RUNGSVENTIL, UEBERS ■■■■■■■■■■  
■■■■■■■■■■ RUNGSVEN- TIL, UEBER ■■■■■■■■■■  
■■■■■■■■■■ RUNGSVENTIL  
■■■■■■■■■■ RUNGSVENT. DRUCKE ■■■■■■■■■■

■■■■■■■■■■ CHMESSER 450 MM MIT ZWISCHENRING 10 MM  
■■■■■■■■■■ CHMESSER 470 MM MIT ZWISCHEN-RING 10 MM  
■■■■■■■■■■ CHMESSER 470 MM MIT ZWISCHEN-RING 10 MM  
■■■■■■■■■■ CHMESSER 515 MM MIT ZWISCHEN-RING 10 MM  
■■■■■■■■■■ CHMESSER 515MM MIT ZWISCHEN- RING 25 MM  
■■■■■■■■■■ CHMESSER 520MM MIT ZWISCHEN- RING 5 MM

# Oder das?

■■■■■■■■■■	NGE HINTEN RECHTS
■■■■■■■■■■	NGE LI HI
■■■■■■■■■■	NGE LI VO
■■■■■■■■■■	NGE LINKS HINTEN
■■■■■■■■■■	NGE LINKS VORNE
■■■■■■■■■■	NGE RE VO
■■■■■■■■■■	NGE RECHTS HINTEN
■■■■■■■■■■	NGE RECHTS VORNE
■■■■■■■■■■	NGE VORNE RECHTS
■■■■■■■■■■	NGE LINKS VORNE

# Und was sagt Ihnen das?

■■■■■■■■■■ ENT.U.KRAFTST.VERBR.ANZ ■■■■■■■■■■  
IN.1FLGL.M.BEIDS.RUECKS ■■■■■■■■■■  
KLAP.LAG.LI.TLW.BEARB.LC ■■■■■■■■■■  
HYDR.ANL.P.P.LAG ■■■■■■■■■■  
■■■■■ PRITZP.EING.N.DATEN ■■■■■■■■■■  
■■■■■ PRITZP.EINGEST.N.DA ■■■■■■■■■■  
■■■■■■■■■■ M.VORM.HINT.LI.  
HYDR.ANL.STEUERU.LAENC ■■■■■■■■■■

**Eher wenig?!**

**Dann sind Sie hier richtig!**

- Die Beharrlichkeit sprachlicher Altdaten
- Was sprachliche Altdaten „auszeichnet“
- Ist die Bereinigung sprachlicher Altdaten machbar?
- Die stufenweise Bereinigung sprachlicher Altdaten
- Was die Technische Redaktion zum Gelingen beiträgt
- Fazit und Ausblick

# Die Beharrlichkeit sprachlicher Altdaten



Erfassung von Produktinformationen  
seit Einführung der ersten EDV-Systeme

- Rasante Zunahme des Datenvolumens ...
- ... bei beschleunigter Entwicklung der IT-Technologie
- ... und zunehmender Durchdringung der Betriebsabläufe  
(Arbeitspositionen, Kataloge, Vertrieb, Logistik, Shop, ...)

⇒ Notwendigkeit der **Datenstabilität**



## Notwendigkeit der **Datenstabilität**

- EDV-Systeme werden migriert,
- die Daten selbst bleiben aber unangetastet,
- althergebrachte Eingabeverfahren werden beibehalten!

⇒ Es entstehen massenhaft Daten, die den heutigen Anforderungen an die sprachliche Qualität nicht genügen ...

⇒ ... und der **Datenkompatibilität** zuwiderlaufen!

# Was sprachliche Altdaten „auszeichnet“



## Eingeschränkte Eingabemöglichkeiten

in älteren EDV-Systemen:

- limitierte Eingabefelder (30, 40, 80 Zeichen)
- limitierter Zeichensatz (typischerweise nur Großbuchstaben)

⇒ niedrige Hemmschwelle für  
**sprachliche Nachlässigkeiten** jeglicher Art

⇒ ... zumal es oftmals an der Technologie mangelt,  
redaktionelle Vorgaben durchzusetzen!

# Was sprachliche Altdaten „auszeichnet“

- durchgängige Großschreibung (Versalschreibung)
- keine Umlaute, kein  $\beta$

EUERGENAET HEISS  
HLAEUCHE FUER KUEHL  
JERTRAEGER ANGESCHWESST UN  
FUER ABREISSICHERUNG

# Was sprachliche Altdaten „auszeichnet“

- willkürliche Setzung (und Nicht-Setzung) von Leerzeichen
- Sperrsatz

STUFE, DRUCKW  
NGSSATZ, MIT SCHA  
RTEIL, ZWISCHEN RAD  
JUNGSTEIL (KONSOLE)|  
VOR-UND RUECK  
A C H T U N G ! V O R

# Was sprachliche Altdaten „auszeichnet“

- unsystematische Eingabe spezieller Datentypen  
(Maßangaben, Schraubenbezeichnungen, Teilenummern usw.)

[REDACTED]	KUNG	12V	140 AH
[REDACTED]	BATTERIEN	2X74AH	
[REDACTED]	3350 +	3700MM	[REDACTED]
[REDACTED]	3350 MM		[REDACTED]
[REDACTED]		DDEN 527 X 665 MM	[REDACTED]
[REDACTED]		DDEN	527X665 MM
[REDACTED]	NDE	M 42X1.5MM	VORN [REDACTED]
[REDACTED]	NDE	M42X1.5 MM	HINTE [REDACTED]

# Was sprachliche Altdaten „auszeichnet“

- massenhaft Abkürzungen  
(auch Spontanabkürzungen, Binnenabkürzungen)

■■■■■■■■■■ ENT.U.KRAFTST.VERBR.ANZ ■■■■■■  
IN.1FLGL.M.BEIDS.RUECKS ■■■■■■  
KLAP.LAG.LI.TLW.BEARB.LC ■■■■■■  
HYDR.ANL.P.P.LAG ■■■■■■  
■■■■■ PRITZP.EING.N.DATEN ■■■■■■  
■■■■■ PRITZP.EINGEST.N.DA ■■■■■■  
■■■■■■■■■■ M.VORM.HINT.LI.  
HYDR.ANL.STEUERU.LAENG ■■■■■■



# Was sprachliche Altdaten „auszeichnet“

- alte Rechtschreibung und orthografische Varianten

[REDACTED]	SEN NUMER	RIERT
[REDACTED]	BER TACHOGRAF	ELEK [REDACTED]
[REDACTED]	BER TACHOGRAPH	ELE [REDACTED]
[REDACTED]	KONTROLLAMPE RUN	[REDACTED]
[REDACTED]	KONTROLLANZEIGE	
[REDACTED]	UM TACHOGRAF	
[REDACTED]	UM TACHOGRAPH	
[REDACTED]	DER KRAFTSTOFFFILTER	



# Was sprachliche Altdaten „auszeichnet“

- vielfältige Falschschreibungen (insbesondere falsche Zusammen- oder Getrennschreibungen und Bindestrichsetzungen)

U RELAIS UND BATTERIE  
ER LINKS UND RECHTSOHNE AUS  
REINFACHTE AUS-FUEHRUNG  
ANSAUGUNG HIN-TEN  
UNG AUTOMATSCH GEREGLT  
AGGREGAT 8 -POLIG  
NGSSATZ 4- POLIG  
NG HIN TEN LINKS  
M EINGEBAU TIN  
ITLICH LINKSUND RECHTS 590X50  
NGSTRAEGER UND QUERT RAEGER  
ELEKTIRSCHEN LEITUNGSSAE  
WAND 180MMVERLAENGERT

# Ist die Bereinigung sprachlicher Altdaten machbar?



# Ist die Bereinigung sprachlicher Altdaten machbar?

Eine manuelle Bereinigung ist  
**ökonomisch und intellektuell kaum leistbar!**

- Redakteure müssten Spontanabkürzungen auflösen,
- sprachliche Fehler ausmerzen,
- terminologische Inkonsistenzen ausräumen,
- jeweils konsistent über den Gesamtdatenbestand!

Bei  $\pm 1$  Minute pro Kurztext von 50 bis 80 Zeichen

→  $\sim 200$  Personentage für 100.000 Datensätze

**... optimistisch geschätzt!**

# Ist die Bereinigung sprachlicher Altdaten machbar?



## Zielsetzung

- signifikant reduzierte Aufwände
- bei konstanter Korrekturleistung
- durch Einsatz von weitestgehend automatischen Verfahren

## Was bedarf es hierfür mindestens?

- Einfache **Programmroutinen**, die spezielle Zeichenkettenersetzungen ausführen, aber auch ...
- **linguistisch intelligente Sprachverarbeitungssoftware** zur automatischen Etablierung der korrekten Rechtschreibung und der terminologischen Konsistenz

# Ist die Bereinigung sprachlicher Altdaten machbar?



## Anforderungen an eine Sprachverarbeitungssoftware

- umfassende **Abdeckung des relevanten Vokabulars**
- Funktion zur **Umwandlung von Versalschreibung** in normale Groß- und Kleinschreibung
- automatische **Korrektur systematischer Falschschreibungen** (z.B. ausgeschriebene Umlaute, alte Rechtschreibung)
- gewichtete Vorschläge für die **Korrektur von Tippfehlern**
- **Auflösung von Mehrdeutigkeiten** basierend auf einer grammatischen Analyse
- Idealerweise **Abgleich mit Unternehmensterminologie**

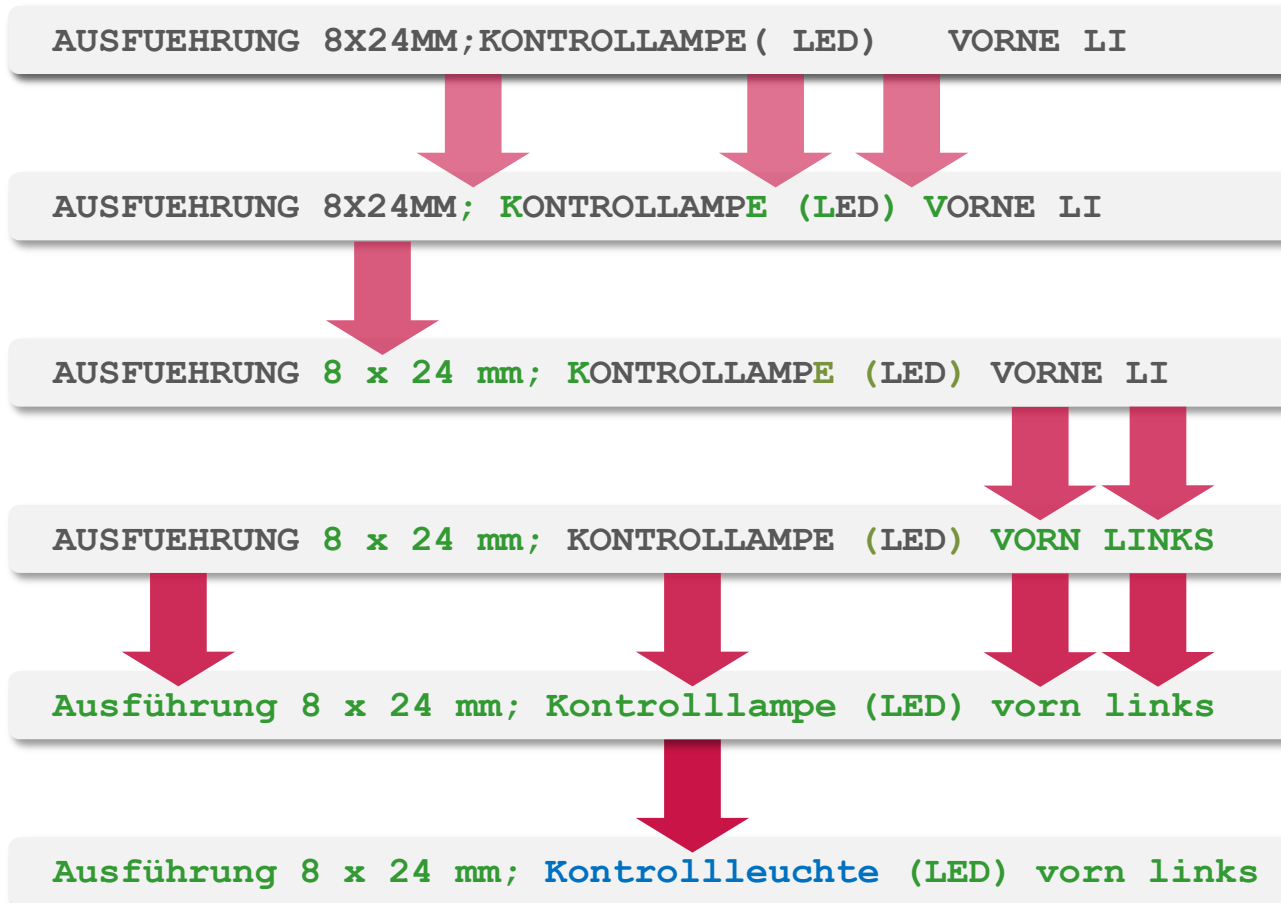
# Ist die Bereinigung sprachlicher Altdaten machbar?

⇒ Mit einer entsprechenden Sprachverarbeitungssoftware kann die Bereinigung sprachlicher Altdaten in einem gestuften Verfahren angegangen werden.

# Die stufenweise Bereinigung sprachlicher Altdaten



# Die stufenweise Bereinigung sprachlicher Altdaten



- Leerzeichen normalisieren
- Datentypen standardisieren
- Zeichenketten ersetzen
- Rechtschreibung korrigieren
- Terminologie konsolidieren



## Normalisierung der **Leerzeichensetzung** durch einfache Ersetzungsrountinen

- Mehrfachleerzeichen auf einzelne Leerzeichen reduzieren
- Leerzeichen am Zeilenanfang oder -ende löschen
- Leerzeichen im Kontext von Satzzeichen setzen bzw. löschen (vor/nach Klammer, Doppelpunkt, Semikolon, Komma)
- ...

→ Berücksichtigung **besonderer Datenkonstrukte**:  
gegliederte Nummern (13 45.00 04), Dezimalzahlen (20,5MM),  
Ordinalzahlen (1.STUFE), Verhältnisangaben (30:50)!

Standardisierung der **Darstellung spezieller Datentypen** durch einfache (Muster-basierte) Ersetzungsroutinen

- Datumsangaben (24.9.2017)
- Maßangaben (30X50X200MM)
- Schraubenbezeichnungen (M42X1.5)
- ...

→ Setzung von **Leerzeichen bzw. Festabständen**

→ Darstellung von **Sonderzeichen**  
(z.B. **✕** versus × als Mal-Zeichen)

# Die stufenweise Bereinigung sprachlicher Altdaten



Nach Normalisierung der Leerzeichensetzung und Standardisierung der Darstellung spezieller Datentypen **erweisen sich erste Texte als Dubletten!**

⇒ Ein Anteil von **3% oder mehr Textdubletten** ist an dieser Stelle keine Seltenheit!

# Die stufenweise Bereinigung sprachlicher Altdaten

## Analyse des Zeichenketteninventars durch Datenauswertungsroutinen und Wortanalysen

- Welche Zeichenketten kommen im Datenbestand wie oft vor?

144557	,
45645	/
41006	<b>MIT</b>
33004	;
28938	<b>UND</b>
24761	<b>FUER</b>
24169	:
21800	<b>LINKS</b>
21735	<b>RECHTS</b>
19578	<b>AB</b>
17477	<b>BEI</b>
16734	<b>IN</b>
14779	(
14676	)
12972	<b>VON</b>
12610	<b>MM</b>
12099	<b>CODE</b>
10892	<b>EINGEBAUT</b>
10329	<b>TEILNUMMER</b>

# Die stufenweise Bereinigung sprachlicher Altdaten

## Analyse des Zeichenketteninventars durch Datenauswertungsroutinen und Wortanalysen

- Welche Zeichenketten kommen im Datenbestand wie oft vor?
- Welche weiteren speziellen Datentypen lassen sich ablesen?

2092	<b>1X</b>	
849	<b>24V</b>	
827	<b>12V</b>	
727	<b>4X2</b>	
535	<b>2X</b>	
429	<b>4X4</b>	
407	<b>V6</b>	
361	<b>6X2</b>	
302	<b>5T</b>	
272	<b>6X4</b>	
218	<b>8X4</b>	
188	<b>4X</b>	
180	<b>400L</b>	
144	<b>300L</b>	
121	<b>A31</b>	
120	<b>3X</b>	
103	<b>1X322</b>	
81	<b>28V</b>	
79	<b>1B</b>	

# Die stufenweise Bereinigung sprachlicher Altdaten

## Analyse des Zeichenketteninventars durch Datenauswertungsroutinen und Wortanalysen

- Welche Zeichenketten kommen im Datenbestand wie oft vor?
- Welche weiteren speziellen Datentypen lassen sich ablesen?
- Welche Abkürzungen sind im Datenbestand belegt?

1564	<b>FA.</b>
659	<b>BZW.</b>
385	<b>NR.</b>
311	<b>FUSSN.</b>
128	<b>BEF.</b>
113	<b>U.</b>
102	<b>ELEKTR.</b>
93	<b>MIN.</b>
68	<b>NO.7</b>
43	<b>MECH.</b>
39	<b>F.</b>
31	<b>ERN.</b>
27	<b>LI.</b>
24	<b>M.P.H.</b>
22	<b>L.</b>
21	<b>BZW.BIS</b>
20	<b>BZW.AB</b>
20	<b>DURCHM.</b>
18	<b>S.</b>

## Analyse des Zeichenketteninventars

durch Datenauswertungsroutinen und Wortanalysen

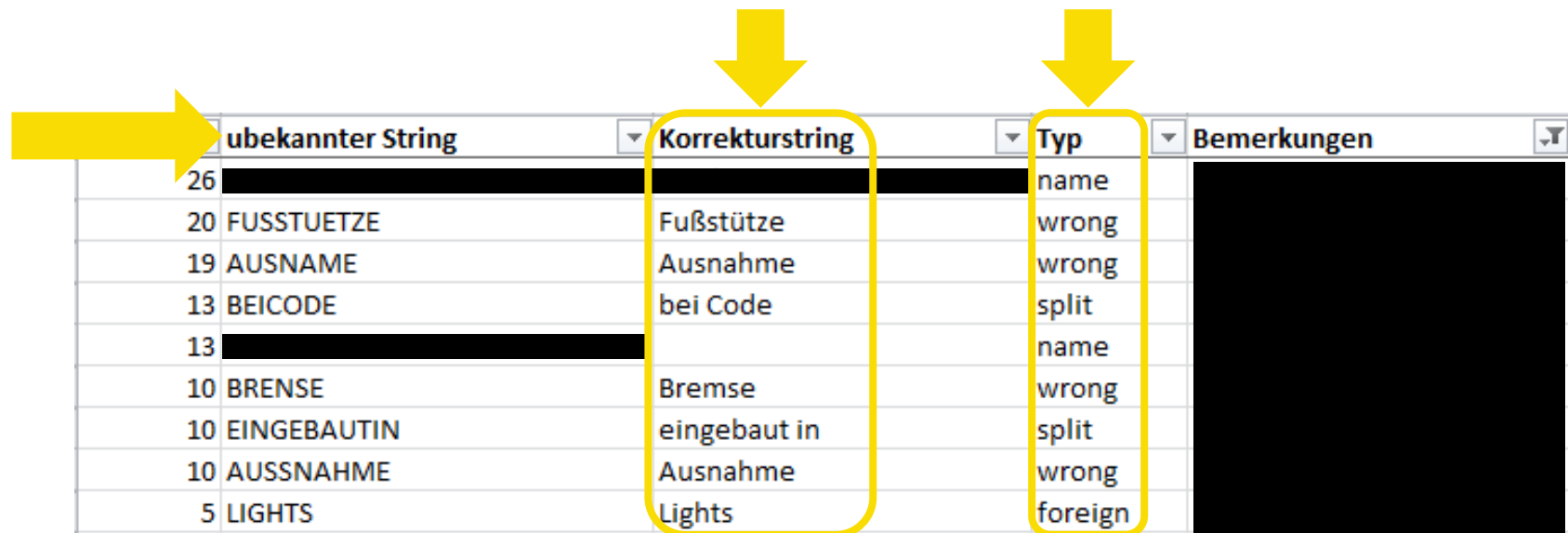
- Welche Zeichenketten kommen im Datenbestand wie oft vor?
- Welche weiteren speziellen Datentypen lassen sich ablesen?
- Welche Abkürzungen sind im Datenbestand belegt?
- Und welche unbekannten Wörter oder Falschschreibungen kommen eigentlich vor?

20	AUSSNAHME
20	VERSCHLUSSCHRAU
20	VORDER
19	FUSSTUETZE
19	AUSNAME
16	STUEZWINKEL
15	AUSNHAME
14	AUSRALIEN
13	AUSERDEM
12	OESTEREICH
11	UEBERSETUNG
10	VONMITTE
10	AUSFUERUNG
9	HAEUSE
9	EINGEBAUTIN
9	NUMMERIC
8	LOESCEN
7	ELECTRIC
7	INSIEHE

# Die stufenweise Bereinigung sprachlicher Altdaten

## Bereitstellung der Ergebnisse

der Zeichenkettenanalyse als einfaches Tabellenformat zur Bearbeitung seitens der Redaktion



The table displays the results of a string analysis. It has four columns: 'unbekannter String', 'Korrekturstring', 'Typ', and 'Bemerkungen'. The 'Bemerkungen' column is currently empty. A yellow arrow points to the first row, and two yellow arrows point to the 'Korrekturstring' and 'Typ' columns respectively.

	unbekannter String	Korrekturstring	Typ	Bemerkungen
26			name	
20	FUSSTUETZE	Fußstütze	wrong	
19	AUSNAME	Ausnahme	wrong	
13	BEICODE	bei Code	split	
13			name	
10	BRENSE	Bremse	wrong	
10	EINGEBAUTIN	eingebaut in	split	
10	AUSSNAHME	Ausnahme	wrong	
5	LIGHTS	Lights	foreign	



## Generierung von Ersetzungsregeln

aus den bearbeiteten Tabellen der Zeichenkettenanalyse

→ Ausführung der Ersetzungen über den Gesamtdatenbestand

- **Expansion von Abkürzungen**, ggf. auch unter Berücksichtigung vorhandener Abkürzungslisten, z.B.

LI	→ LINKS
HI	→ HINTEN
BEZ	→ BEZEICHNUNG
REP.SATZ	→ REPARATURSATZ

- **Auftrennung falscher Zusammenschreibungen, z.B.**

BEICODE           → BEI CODE  
EINGEBAUTIN      → EINGEBAUT IN

- **Auswahl bei Wortvarianten, z.B.**

WAAGERECHT      → WAAGRECHT  
VORNE             → VORN

- **Korrektur von unbekanntem Wörtern ohne eindeutigen Korrekturvorschlag, z.B.**

AUSNAME           → AUSNAHME;AUSNAHMEN  
BRENSE             → BRENNE;BREMSE

- **Bekanntmachung unbekannter „Wörter“**, die nicht – oder auf bestimmte Weise – korrigiert werden sollen (Akronyme, Codes, Firmennamen usw.), z.B.

DFT

Z97G67

TOTAL

→ ggf. **Aufnahme in ein Spezialwörterbuch** der Sprachverarbeitungskomponente

⇒ Der Datenbestand ist „reif“ für die linguistische Analyse und Rechtschreibkorrektur!

**Rechtschreibkorrektur** basierend auf einer linguistisch fundierten Sprachverarbeitungs-komponente

- **95%** und mehr der in Versalschreibung vorliegenden Texte können vollständig in **korrekte Groß- und Kleinschreibung** umgewandelt werden.
- Ausgeschriebene Umlaute und sonstige **Falschreibungen** sind **umfassend und konsistent korrigiert**.

## Rechtschreibkorrektur basierend auf einer linguistisch fundierten Sprachverarbeitungs-komponente

ELEKTRISCHE LEITUNG AN ABSCHIRMBLECH	
ELEKTRISCHE LEITUNG AN AMPEREMETER	
ELEK	Elektrische Leitung an Abschirmblech
ELEK	Elektrische Leitung an Amperemeter
ELEK	Elektrische Leitung an Anschlussstück
ELEK	Elektrische Leitung an Aufhängeöse
ELEK	Elektrische Leitung an Batterie Hauptschalter
ELEK	Elektrische Leitung an Druckschalter
ELEK	Elektrische Leitung an Gehäuse
ELEK	Elektrische Leitung an Geräteträger
ELEK	Elektrische Leitung an Glühlampe
	Elektrische Leitung an Halter Zwischengehäuse

Nur wenige Texte enthalten jetzt noch **einzelne Wörter, die nicht korrigiert wurden.**

→ Grund hierfür: Wortmehrdeutigkeiten, die im gegebenen Kontext nicht aufgelöst werden konnten, z.B.

MASSE Ladedose → Masse / Maße ?

⇒ **Spezielle Auszeichnung** solcher Fälle zur Unterstützung einer gezielten manuellen Nachbereinigung

⇒ Auch **datenspezifische Vereindeutigungsregeln** sind denkbar!

Kann die Rechtschreibkorrektur auch **über das Ziel hinauschießen?**

→ Fehlanalysen und somit auch Fehlkorrekturen sind vereinzelt möglich, z.B.

ANSCHLUSS AN TRAEGER PUMPE → Anschluss an träger Pumpe

⇒ Erfassung aller Wortmehrdeutigkeiten

⇒ Implementierung **datenspezifischer Vereindeutigungsregeln**

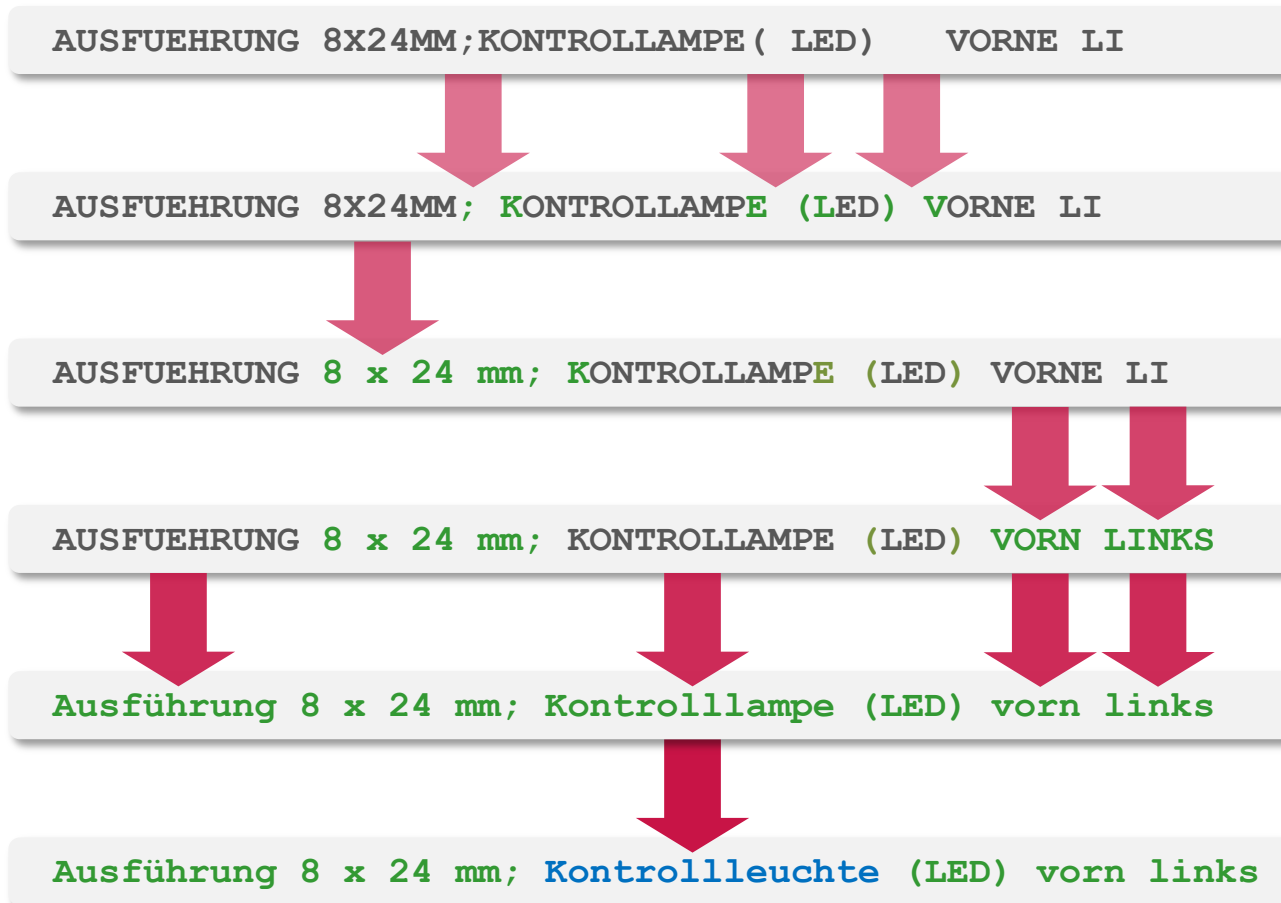
**Terminologiekonsolidierung** basierend auf einem Abgleich mit der Unternehmensterminologie

- Erkennung von verbotenen Termen und Termvarianten
- Ersetzung durch Vorzugsbenennungen, z.B.  
Kontrolllampe → Kontrollleuchte
- ... oder zumindest Auszeichnung als terminologische Problemfälle

⇒ Die stufenweise Bereinigung sprachlicher Altdaten auf der Grundlage automatischer Verfahren ist damit fürs Erste abgeschlossen!



# Die stufenweise Bereinigung sprachlicher Altdaten



- Leerzeichen normalisieren
- Datentypen standardisieren
- Zeichenketten ersetzen
- Rechtschreibung korrigieren
- Terminologie konsolidieren

# Was die Technische Redaktion zum Gelingen beiträgt



# Was die Technische Redaktion zum Gelingen beiträgt



Die **Technische Redaktion** unterstützt maßgeblich den Prozess der Bereinigung sprachlicher Altdaten.

- Abstimmung mit den redaktionellen Vorgaben
- Anpassung an die spezifischen Datengegebenheiten

⇒ Erhöhung der Bereinigungsquote

⇒ geringere Aufwände für manuelle Nachbereinigung

⇒ Vermeidung von Fehlkorrekturen

## Bereitstellung von Wissensquellen

- Redaktionsleitfaden  
→ Muster-Spezifikationen für spezielle Datentypen
- Abkürzungslisten  
→ Expansion der Langformen
- Terminologie  
→ Einsetzung der Vorzugsbenennungen
- Marketingbenennungen, Produktnamen, Zulieferfirmen ...  
→ Vermeidung von Fehlkorrekturen

**Iterative Bewertung von Zwischenergebnissen,**  
insbesondere der Analyse des Zeichenketteninventars

- Expansion von Spontan- und Binnenabkürzungen
- Festlegungen zu Tippfehlerkorrekturen und Wortvarianten
- Klassifikation unbekannter Wörter als Spezialvokabular
- Stichprobenprüfungen

## Manuelle Nachbereinigung

- Sichtung speziell ausgezeichnete Textstellen
  - nicht korrigierte Wörter (wegen Wortmehrdeutigkeiten)
  - terminologische Problemfälle

⇒ Auch die letzten verbleibenden Problemfälle sind intellektuell bewertet und korrigiert!

# Fazit und Ausblick



**Die über viele Jahre (Jahrzehnte) aufgeschobene  
Bereinigung sprachlicher Altdaten ist machbar!**



## Ein Bereinigungsgrad von 95% und mehr ist erreichbar

- mit einigen grundlegenden Ersetzungsroutinen
- ... und dem Einsatz einer hochwertigen Sprachverarbeitungskomponente.

Ein höherer Bereinigungsgrad wird durch die **Abstimmung des Bereinigungsprozess auf die unternehmensspezifischen Gegebenheiten** erzielt.

→ Berücksichtigung redaktioneller Vorgaben und Referenzdaten

→ Einbeziehung der Technischen Redaktion

## Variantenermittlung

auf der Grundlage der bereinigten Daten

- Benennungsvarianten
- Wortstellungsvarianten
- Formulierungsvarianten

FONDTUERE RECHTS  
FONDTUER, RECHTS  
FONDTUER,RECHTS  
RECHTS FONDTUER  
SEIDENTUER HINTEN RECHTS  
TUER HINTEN RECHTS  
TUER HINTEN,RECHTS



Fondtür, rechts

- Eliminierung unnötiger Textvarianten
- Reduktion der informationellen Redundanz im Gesamtdatenbestand
- ⇒ Fit für den effektiven Einsatz sprachtechnologischer Anwendungen wie die eines TM oder AM

⇒ **Fit für die Zukunft!**

**Vielen Dank für Ihre Aufmerksamkeit.**