

Migration und Konsolidierung

Dipl.-Ing. José Manuel de la Rosa Govantes

Vortrag am 26.04.2017 zur GSE Tagung Hamburg

Agenda



Über uns

Die semantics Kommunikationsmanagement GmbH

Ausgangspunkt und Ziel

Vorbemerkungen und Motivation

Ohne Fleiß kein Preis

Import und maschinelle Analyse

Automatisierung der Redaktionsarbeit

Maschinelle Konsolidierung

Ergebnis

Erwartungen und Ungeplantes

Status Quo und Ausblick

Anwendungen und Entwicklungen

Über uns

Die semantics Kommunikationsmanagement GmbH

semantics Kommunikationsmanagement GmbH



- Gründung 2000 als Spin-off der RWTH Aachen (Prof. Dr. Christian Stetter)
- Interdisziplinäres Team aus z.Zt. 38 Informatikern, Sprachwissenschaftlern und Psychologen
- Verbindung von Forschung und Praxis zur Textverständlichkeit und Texttechnologie
 - > Sprachschulungen, -richtlinien, -analysen, Workshops etc.
 - > Technische Lösungen für große Textbestände (Analyse, Management, Pflege)
 - > Empirische Studien
 - > Forschung in Linguistischer Informatik
- Ferner: Digitalisierungs- und Recherchesysteme für Bibliotheken, Archive und Museen (Marktführer D-A-CH)

Ausgangspunkt und Ziel

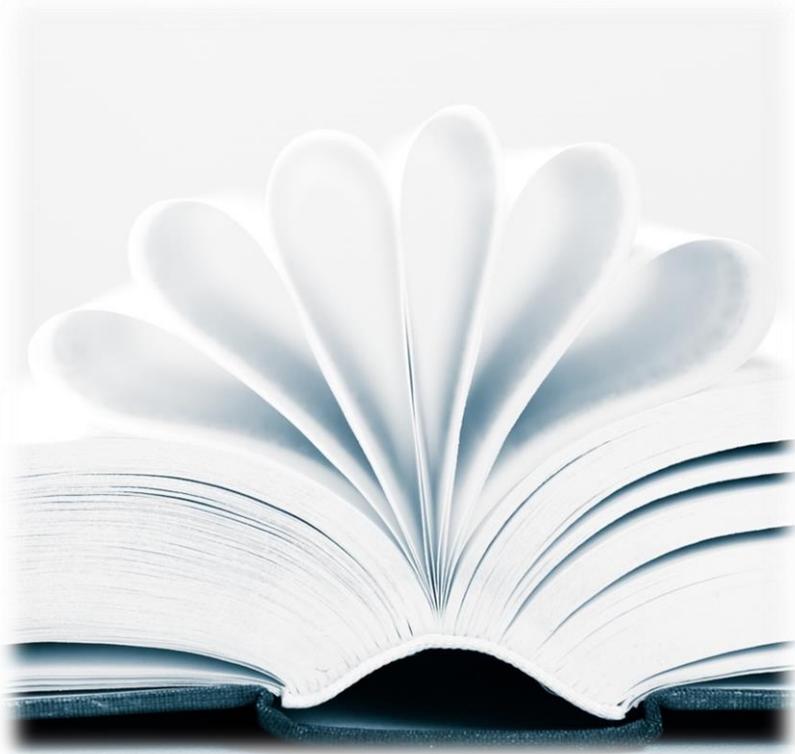
Vorbemerkungen und Motivation

Kontext – Zeitalter der Digitalisierung



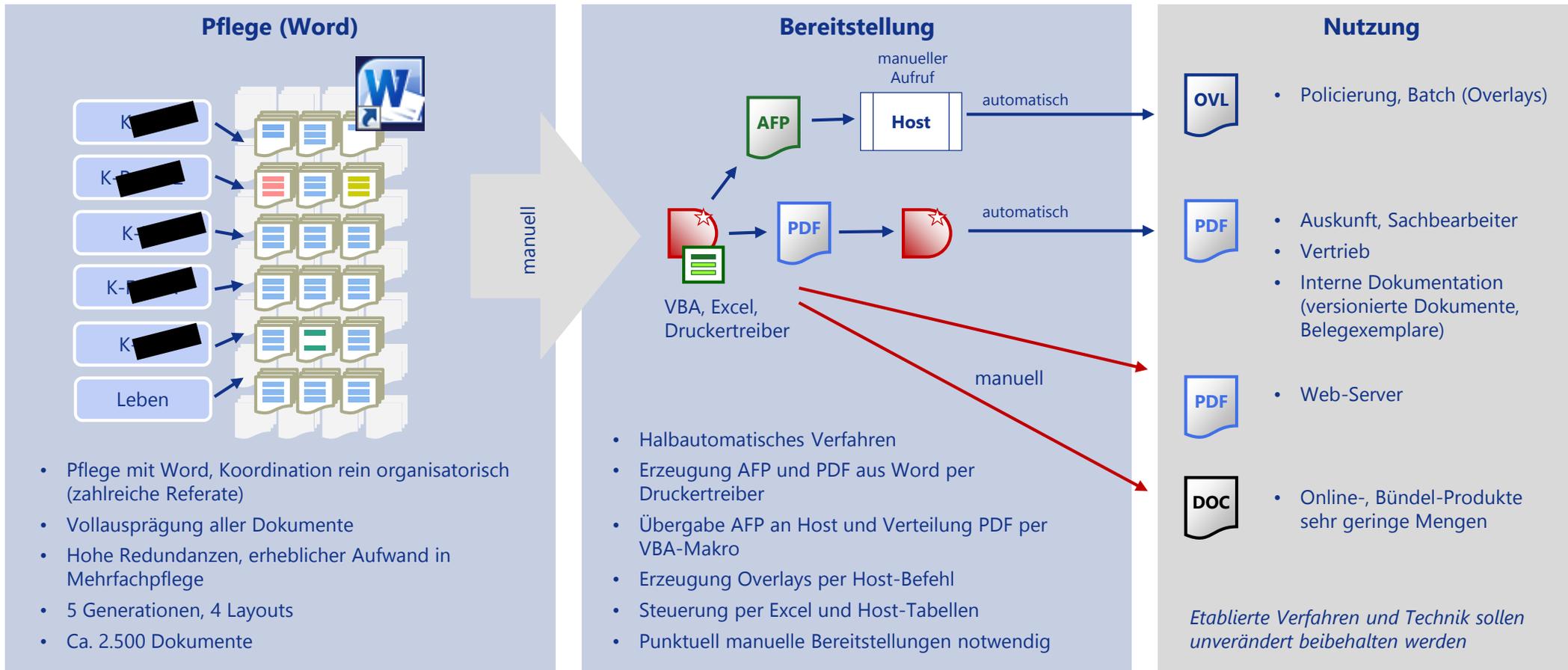
- Zeitalter der *Digitalisierung*, BigData, künstlichen Intelligenz
auch Chatbots, Fake News, Toxic Comments, ...
- Ansprüche der Verbraucher *steigen*
Bzgl. Kommunikation: Inhalt, Reaktionszeit, Kanal
- Kundenansprache optimieren: Bindeglied *Sprache*
Die richtigen Infos, zur richtigen Zeit, im richtigen Umfang und richtigen Kanal – und dabei den richtigen Ton treffen
- Kommunikation muss *entworfen* und *geplant* werden
Brief, Email, Portale, Apps, Kurznachrichten, Social Media... sind unterschiedlich zu bedienen; Aufwand steigt
- Potential *Gestaltung der Sprache*:
 - ⇒ *Nachfrage- und Beschwerdeaufkommen reduzieren*
 - ⇒ *Empfehlungen, Kundenbindung und -zufriedenheit erhöhen*

Kontext – Sprache gestalten

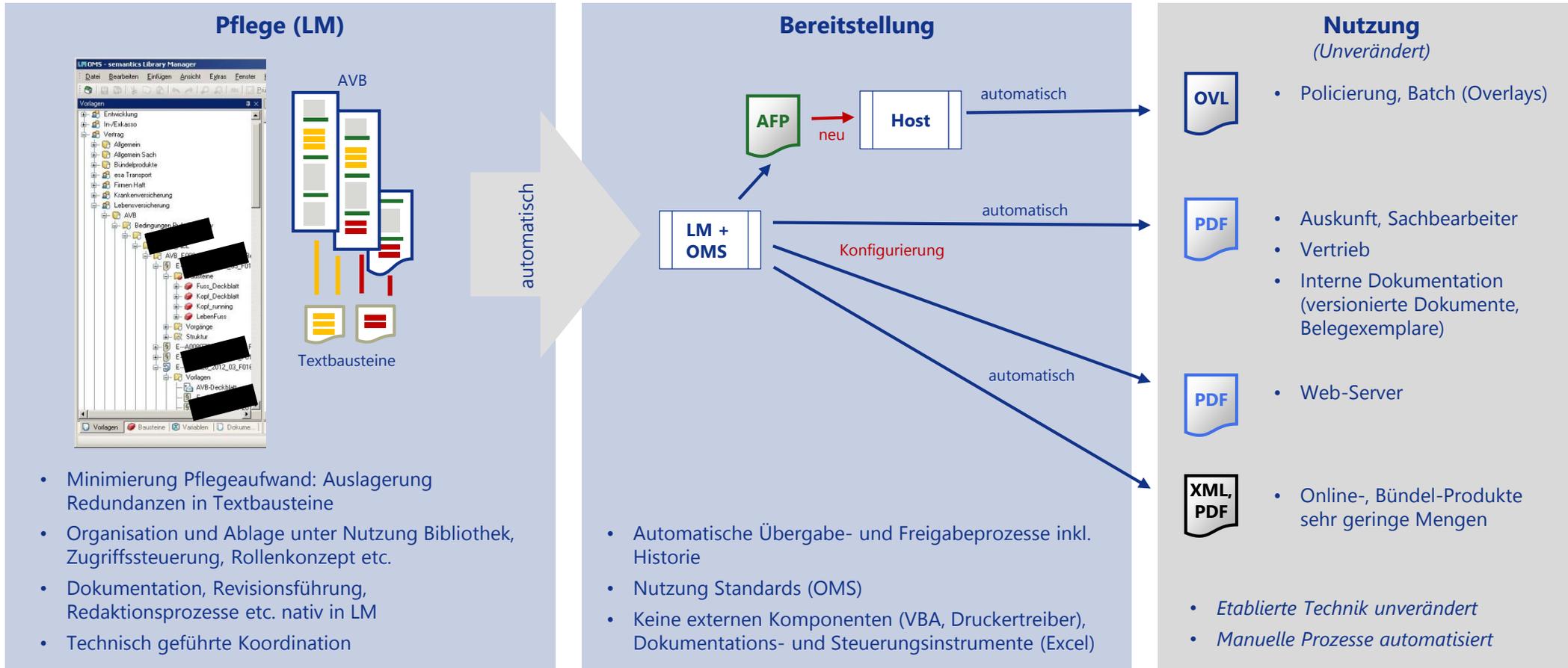


- Heute kein (sprach-)wissenschaftlicher Vortrag
- Unsere Technik adressiert das *systemische Problem*: Spannungsfeld Textpflege
 - > Sehr große Textbestände, lange Historie
 - > Große Organisationen, viele beteiligte Personen
 - > Verschiedene Kompetenzen, Rollen und Interessen
 - > Gemeinsame Verantwortung für den Textbestand
- Anwendung dieser Technik im vorgestellten Projekt
- 2 abgeschlossene Großprojekte (D, CH); 8 laufende Projekte
- Migration ohne Konsolidierung möglich und umgekehrt

Ausgangssituation – Redundante AVB



Ziel: Konsolidierte Redundanzen

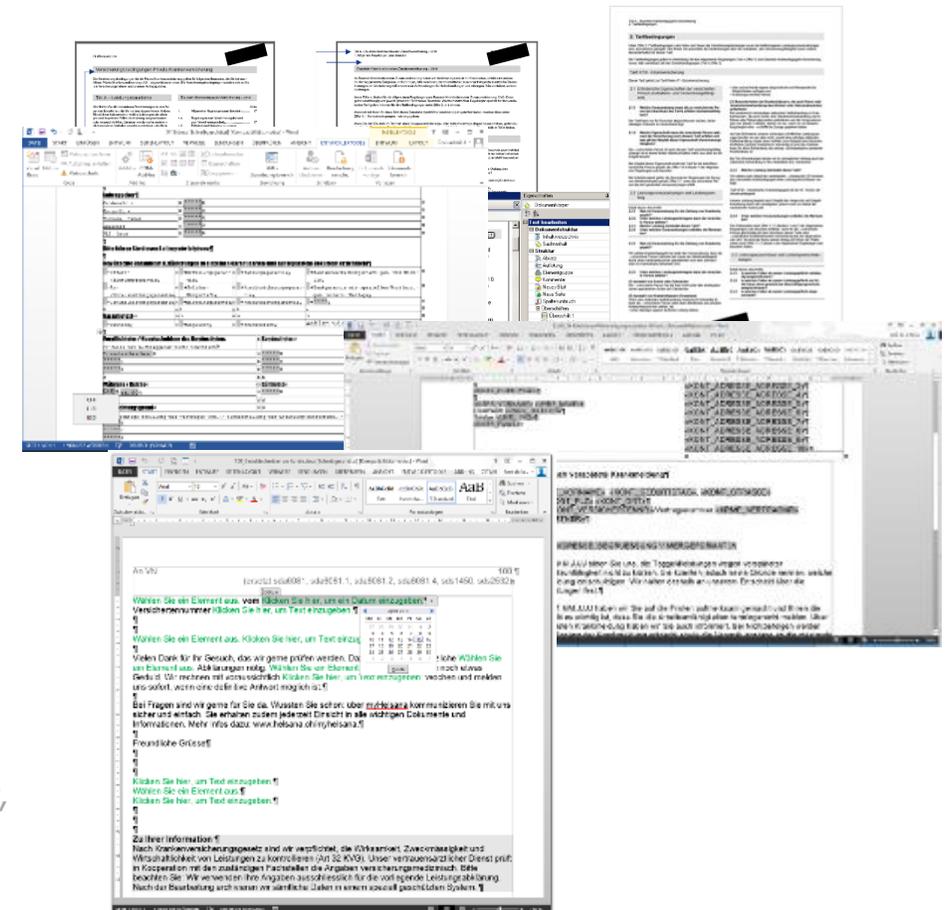


Ohne Fleiß kein Preis

Import und maschinelle Analyse

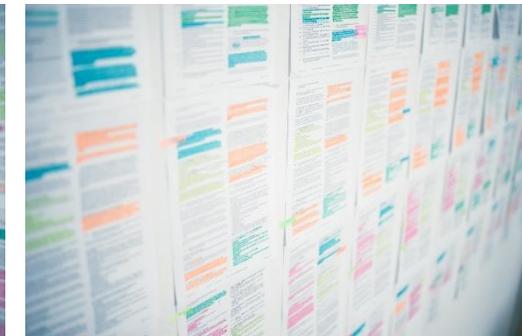
Import – Standard Word, daher „ganz einfach“ ...

- Import Word-Format (andere ebenfalls möglich)
hier: Jahrzehnte gewachsene Variationen...
- Material moderat anspruchsvoll
ein-/zweispaltiges Layout, Formatierungen, Inhaltsverzeichnisse, ...
- Metadaten aus Dateinamen, Textbereichen
auch zur Zuweisung Layout, Dokumenttypen
- Herausforderungen durch „manuelles Layouting“
lokale Formatierung, manuelle Umbrüche und Silbentrennung, Tabellen durch Tabs, manuelle Zähler / Inhaltsverzeichnisse, ... (Lösung z.B. durch Ersetzungstabelle und Baustein-Templates)
- Vielfalt durch Variablen und Eingabeelemente...
Textmarken, Formularfelder, Feldfunktionen, Inhaltssteuerelemente, gelbe Farbe. Oder: Sprachliche Anweisungen. (Einrichtung des Imports je Textbestand)



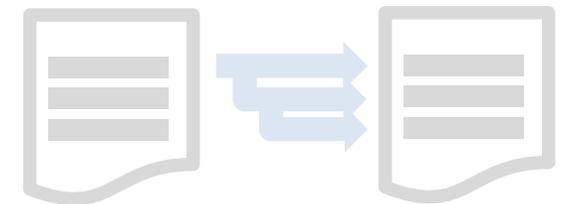
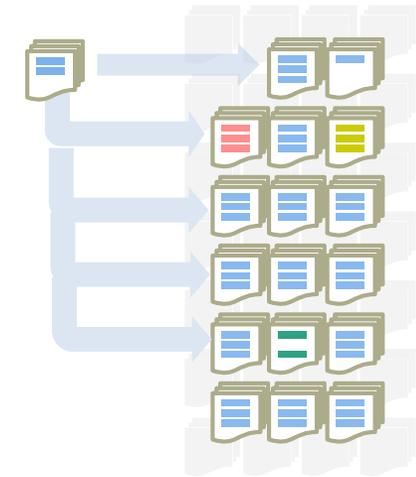
Analyse – Erfahrungen

- Erfahrung mit „kleinen“ Projekten: z.B. Redundanzanalyse schweizerische KV, 120 Dokumente, überschaubar (...?)
- Nutzung Vergleichstool; Ergebnisse konservieren schwierig
- Nach kurzer Zeit entsteht Raumbedarf – wie übergreifende Redundanzen visualisieren?
- Überblick gelang durch Nutzung aller Dimensionen...
- Interesse an weiteren Projekten sehr gering (bei der Belegschaft)



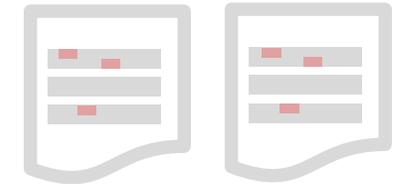
Analyse – Herleitung Zahlen

- Vergleich d Dokumente führt zu $D = \frac{d*(d-1)}{2}$ Operationen
- Vergleich jedes Absatzes zweier Dokumente $A = a_1 * a_2$
- Insgesamt $G = D * A$ Textpassagen zu vergleichen
- Im Beispiel
 - > 120 Dokumente, je 15-20 Absätze
 - > $D = 7.080, \bar{A} \approx 289 \Rightarrow G \approx 2.046.120$
 - > Optimierung durch intuitive Bewertung potentieller Ähnlichkeit
- Im vorgestellten Projekt
 - > 2.500 Dokumente, je ca. 200 Absätze
 - > $D = 3.123.750, \bar{A} \approx 40.000 \Rightarrow G \approx \mathbf{124.950.000.000}$
 - > Maschinelle Verarbeitung und Optimierung notwendig (*sonst bei 1.000/s ca. 4 Jahre*)

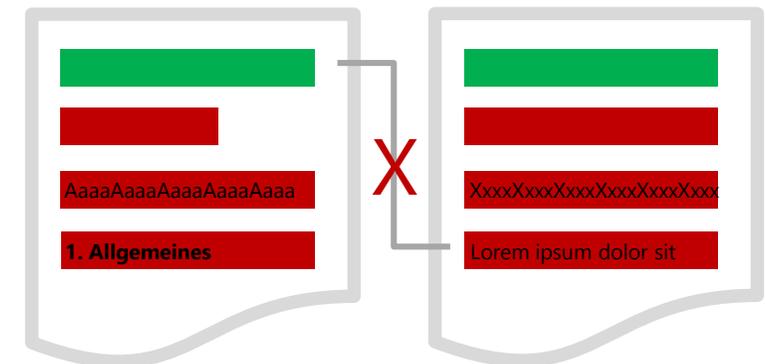


Analyse – Algorithmen

- Suche nach *identischen* und *ähnlichen* Textpassagen; kein Tokenizing o.ä. sinnvoll (*punktueller Abweichungen: Tarifbezeichner, Erstattungsquote, ...*)
- Auswahl: Levenshtein-Distanz, Entfernung zwischen Wörtern; auf Absätze anwenden (*35% Toleranz*)
- Bewährte Algorithmen vorhanden, aber höherer Ordnung (*gewisse Laufzeit, aufgrund Masse signifikant*)
- Optimierung bzw. Vorausschluss u.a. durch
 - > Identität (*Speichervergleich in „Nullzeit“ – keine weitere Prüfung notwendig*)
 - > Länge (*Zeichenanzahl Absatz*)
 - > Fingerprint (*Checksumme Absatz*)
 - > Absatzformat (*Überschriften vs. Text vs. Aufzählungen vs. ...*)
 - > Entfernung (*Position im Dokument*)

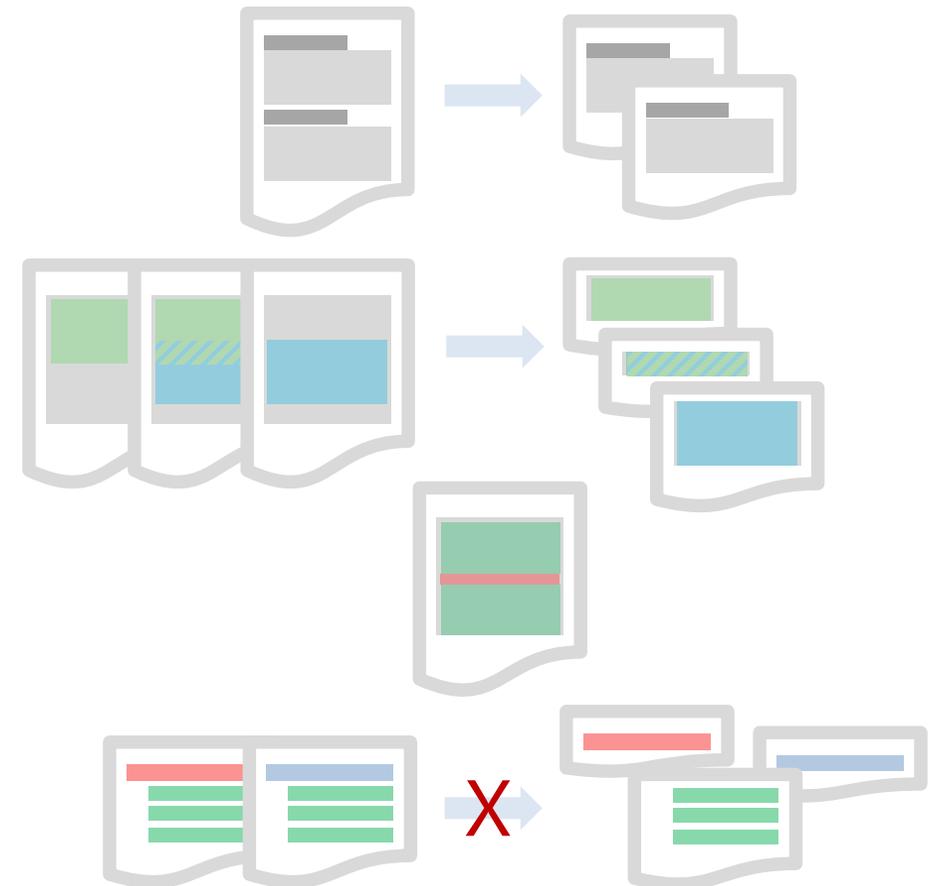


	B	A	N	A	N	E
0	1	2	3	4	5	6
A	1	1	1	2	3	4
N	2	2	2	1	2	3
A	3	3	2	2	1	2
N	4	4	3	2	2	1
A	5	5	4	3	2	2
S	6	6	5	4	3	3



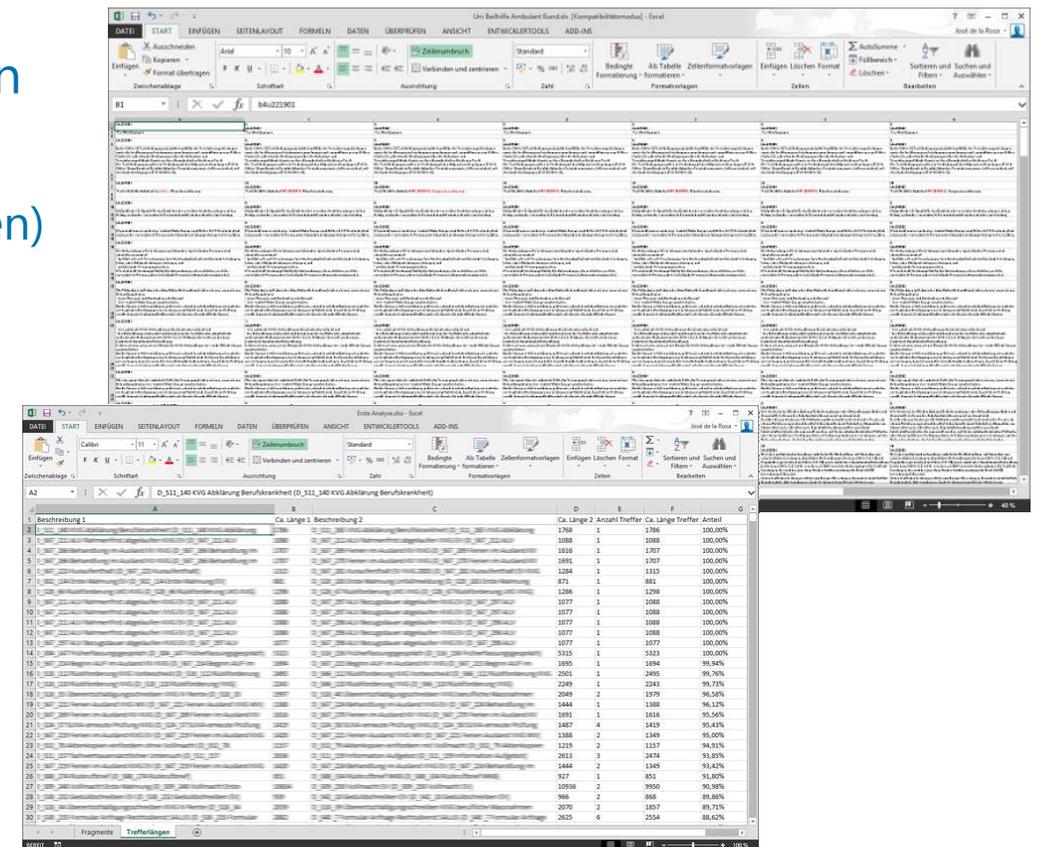
Analyse – fachliche Optimierungen zur Wiederverwendbarkeit

- Splitten großer Fragmente
z.B. anhand Überschriften, künftige Wiederverwendung
- Splitten von „Overlaps“
höhere Wiederverwendung
- Zusammenhalten kurze „Ausreißer“
durch Variablen aufzufangen
- Zusammenhalten fachlicher Einheiten
z.B. Einleitende Sätze und Aufzählungen



Analyse – Darstellung Ergebnisse

- Herausforderung: Tausende Ergebnisse visualisieren
- Frühe Entscheidung: Zweistufiges Verfahren
 - > Gruppenweise Konsolidierung (Tarif-, Produktgruppen)
 - > Nachträgliches 2. Konsolidieren Textbausteine
- Sorge ob fremder Formate (spitze Klammern...)
 - > Fachabteilung ohne technische Kenntnisse
 - > Konzept mit Excel (andere möglich)
- Visualisierung „Treffer“ als Zeilen im Sheet
 - > Nur Inhalt, keine Formatierung
 - > Angabe % Abweichung
- Zusätzlich Ranking Dokumentähnlichkeit möglich



Import und Analyse – Zusammenfassung

- Für Fachabteilung einfaches, effizientes Verfahren
- Import der Dokumente in Library Manager (LM)
- Analyse von Gruppen durch LM
- Bereit für weitere Verarbeitung: Erzeugung Textbausteine zur Konsolidierung

Kontrolle und Steuerung z.B. per Übersicht durch Excel-Sheet

The screenshot displays the 'Library Manager' application window. On the left, a tree view shows a hierarchy of document groups under 'AVB'. The main area shows a detailed view of a document, including 'Admin Informationen' and 'Angemeldete Benutzer'. An 'Analyse Dokumente (3/5)' dialog box is open, showing a progress bar. In the foreground, an Excel spreadsheet is overlaid, displaying a table with columns for document descriptions and numerical data.

Ca. Länge 1	Ca. Länge 2	Anzahl Treffer	Ca. Länge Treffer	Anteil
1. ...	1786	1	1786	100,00%
2. ...	1088	1	1088	100,00%
3. ...	1767	1	1767	100,00%
4. ...	1696	1	1696	100,00%
5. ...	1315	1	1315	100,00%
6. ...	861	1	861	100,00%
7. ...	1296	1	1296	100,00%
8. ...	1088	1	1088	100,00%
9. ...	1696	1	1696	100,00%
10. ...	1088	1	1088	100,00%
11. ...	1696	1	1696	100,00%
12. ...	1696	1	1696	100,00%
13. ...	1696	1	1696	100,00%
14. ...	1696	1	1696	100,00%
15. ...	1696	1	1696	100,00%
16. ...	1696	1	1696	100,00%
17. ...	1696	1	1696	100,00%
18. ...	1696	1	1696	100,00%
19. ...	1696	1	1696	100,00%
20. ...	1696	1	1696	100,00%
21. ...	1696	1	1696	100,00%
22. ...	1696	1	1696	100,00%
23. ...	1696	1	1696	100,00%
24. ...	1696	1	1696	100,00%
25. ...	1696	1	1696	100,00%
26. ...	1696	1	1696	100,00%
27. ...	1696	1	1696	100,00%
28. ...	1696	1	1696	100,00%
29. ...	1696	1	1696	100,00%
30. ...	1696	1	1696	100,00%
31. ...	1696	1	1696	100,00%
32. ...	1696	1	1696	100,00%
33. ...	1696	1	1696	100,00%
34. ...	1696	1	1696	100,00%
35. ...	1696	1	1696	100,00%
36. ...	1696	1	1696	100,00%
37. ...	1696	1	1696	100,00%
38. ...	1696	1	1696	100,00%
39. ...	1696	1	1696	100,00%
40. ...	1696	1	1696	100,00%
41. ...	1696	1	1696	100,00%
42. ...	1696	1	1696	100,00%
43. ...	1696	1	1696	100,00%
44. ...	1696	1	1696	100,00%
45. ...	1696	1	1696	100,00%
46. ...	1696	1	1696	100,00%
47. ...	1696	1	1696	100,00%
48. ...	1696	1	1696	100,00%
49. ...	1696	1	1696	100,00%
50. ...	1696	1	1696	100,00%
51. ...	1696	1	1696	100,00%
52. ...	1696	1	1696	100,00%
53. ...	1696	1	1696	100,00%
54. ...	1696	1	1696	100,00%
55. ...	1696	1	1696	100,00%
56. ...	1696	1	1696	100,00%
57. ...	1696	1	1696	100,00%
58. ...	1696	1	1696	100,00%
59. ...	1696	1	1696	100,00%
60. ...	1696	1	1696	100,00%
61. ...	1696	1	1696	100,00%
62. ...	1696	1	1696	100,00%
63. ...	1696	1	1696	100,00%
64. ...	1696	1	1696	100,00%
65. ...	1696	1	1696	100,00%
66. ...	1696	1	1696	100,00%
67. ...	1696	1	1696	100,00%
68. ...	1696	1	1696	100,00%
69. ...	1696	1	1696	100,00%
70. ...	1696	1	1696	100,00%
71. ...	1696	1	1696	100,00%
72. ...	1696	1	1696	100,00%
73. ...	1696	1	1696	100,00%
74. ...	1696	1	1696	100,00%
75. ...	1696	1	1696	100,00%
76. ...	1696	1	1696	100,00%
77. ...	1696	1	1696	100,00%
78. ...	1696	1	1696	100,00%
79. ...	1696	1	1696	100,00%
80. ...	1696	1	1696	100,00%
81. ...	1696	1	1696	100,00%
82. ...	1696	1	1696	100,00%
83. ...	1696	1	1696	100,00%
84. ...	1696	1	1696	100,00%
85. ...	1696	1	1696	100,00%
86. ...	1696	1	1696	100,00%
87. ...	1696	1	1696	100,00%
88. ...	1696	1	1696	100,00%
89. ...	1696	1	1696	100,00%
90. ...	1696	1	1696	100,00%
91. ...	1696	1	1696	100,00%
92. ...	1696	1	1696	100,00%
93. ...	1696	1	1696	100,00%
94. ...	1696	1	1696	100,00%
95. ...	1696	1	1696	100,00%
96. ...	1696	1	1696	100,00%
97. ...	1696	1	1696	100,00%
98. ...	1696	1	1696	100,00%
99. ...	1696	1	1696	100,00%
100. ...	1696	1	1696	100,00%

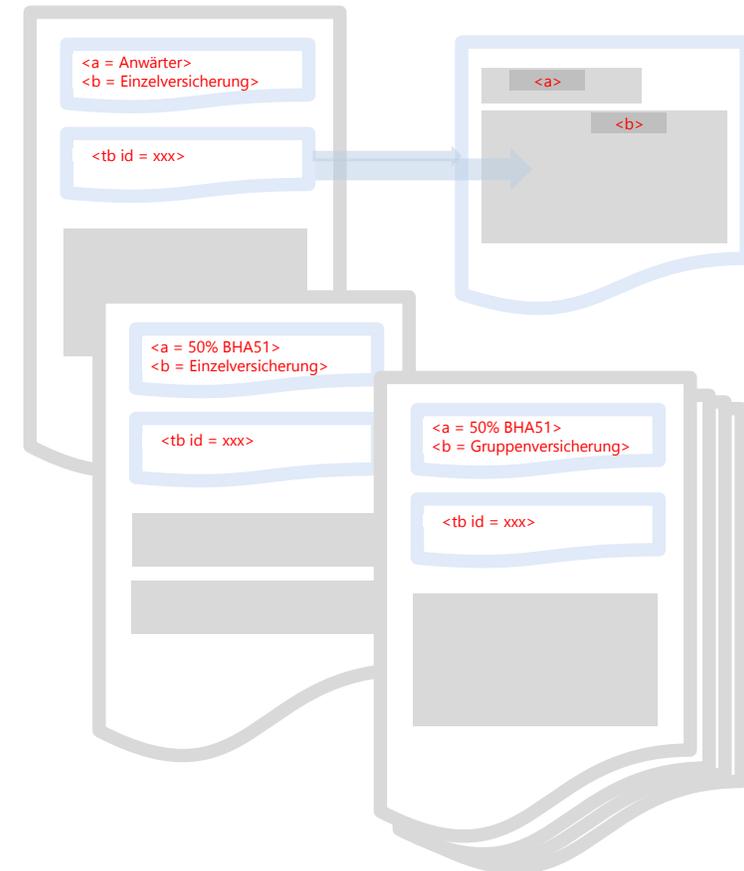
(Fallbeispiel bewusst unleserlich, Urheber- / Kundenschutz)

Automatisierung der Redaktionsarbeit

Maschinelle Konsolidierung

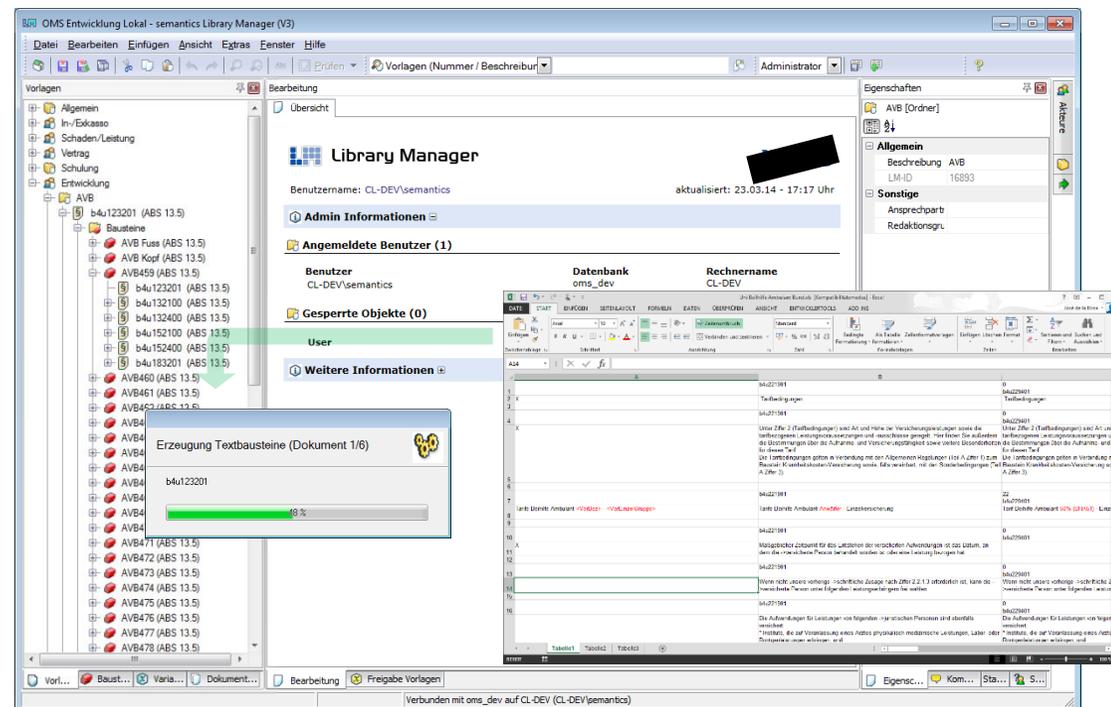
Konsolidierung – Vorgang

- Bausteine anlegen
 - > Inhalt aus Dokument kopieren
 - > Variablen einbetten
- In Dokument Texte durch Bausteinreferenzen ersetzen
- In Dokument Werte für Variablen eintragen
- Ersetzungen und Werte in allen Dokumenten
- Aufgrund Masse maschinelle Durchführung notwendig
- Fachliche Steuerung und Entscheidung gefordert
dabei weiterhin ohne technische Kenntnisse



Konsolidierung – Technik

- Steuerung per Excel
- Kennzeichnung auszulagernder Fragmente
- Bezeichnung gewünschter Variablen
- Punktueller Ausschluss möglich
- Anschließend automatische Verarbeitung durch LM
 - > Fragmente aus Dokument bekannt
 - > Position Variablen aus Excel
 - > Variablenwerte aus Dokumenten gemäß Position in Excel



Ergebnis

Erwartungen und Ungeplantes

Erwartungen übertroffen

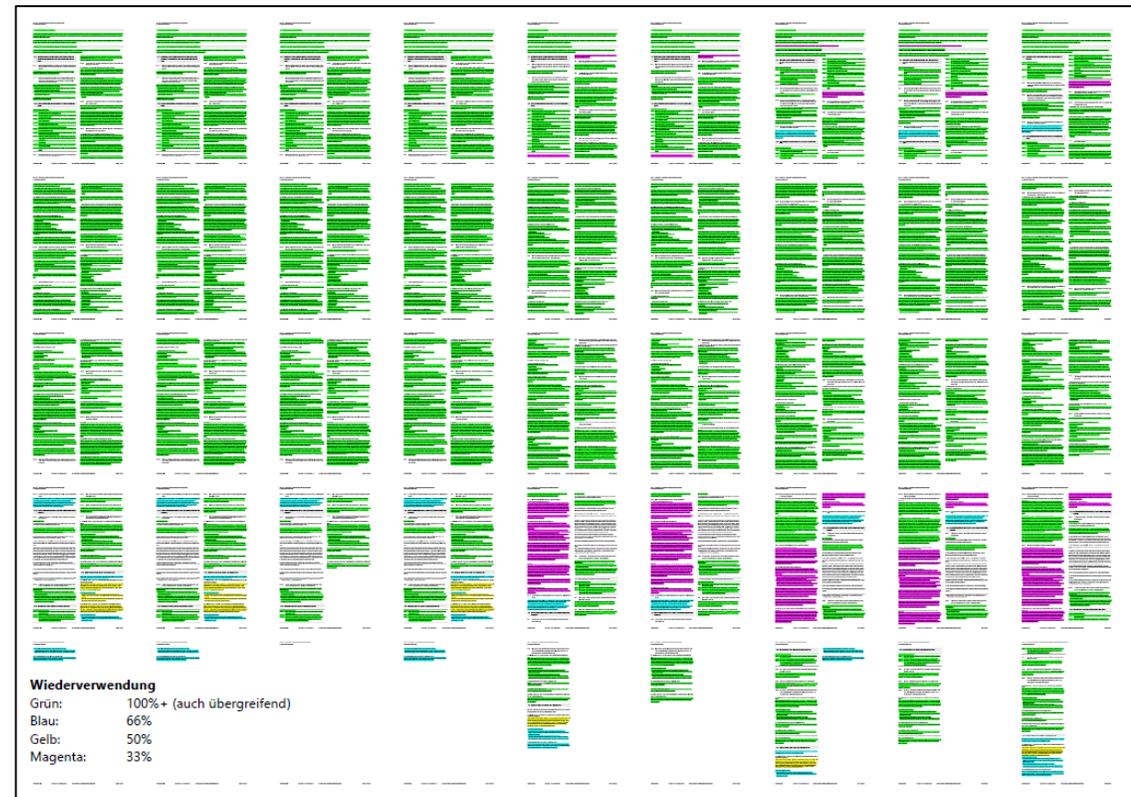
- 99% Abdeckung
- 15.000 Textbausteine (TB) erzeugt
- 180.000 Referenzierungen TB
- 400fache Wiederverwendung relevante TB (*VVG-Reform, ... zudem zahlreiche TB in allen Dokumenten*)
- 300 Variablen, bis 330 Ausprägungen
- Ersparnis vs. manuellem Projekt ca. **20 Personenjahre**

Analyse 1 PT / 10 Dokumente \approx 250 PT

Erzeugung je TB 10 Min. \approx 315 PT

Einbettung je Verwendung 10 Min. \approx 3.750 PT

⇒ *Projekt ohne Technik nicht realisierbar*



Ungeplante (positive) Nebeneffekte

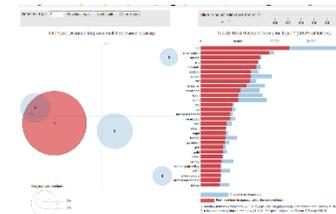
- „Selbsteilung“ großer Textkorpora: Massenvergleich spürt kleine Fehler auf
 - > Interpunktion, Orthographie, Doppellungen, Inkonsistenzen („oder“ statt „bzw.“), ...
- Einsparungen und Minimierung Fehlerpotential in Linie durch Library Manager
 - > Einführung LM als strukturiertes Content- und Textsystem inkl. Darstellung Abhängigkeiten, Redaktionsprozessen, Dokumentation, Annotationen, Versions- und Historienführung, Variantenlogik
 - > Ablösung manueller Bereitstellungsverfahren durch fachliche und technische Freigabeprozesse
- Zentrale Bereitstellung AVB für heterogene Systeme
 - > Single-Source Bereitstellung für Altsystem, Kundenportal, Produktives OMS, ...
- Belobigung durch Abteilung „Interne Revision“
 - > Turnusmäßige Untersuchung führte zu besonderer Auszeichnung der Arbeitsweise der verantwortlichen Abteilung

Status Quo und Ausblick

Anwendungen und Entwicklungen

Ausblick – Entwicklungen

- Neue Herausforderungen: Analyse von 24.000 Dokumenten
- Clustering der Dokumente z.B. mit „Latent Semantischer Analyse“ (LSA) oder „Latent Dirichlet Allocation“ (LDA), ...
- Kombination LSA und Thesaurus
 - > Semantische Ähnlichkeitssuche
 - > Operativer „Suchindex“ in Redaktionssystem
 - > Suche nach Sekundärbegriffen
- Reverse Engineering der Fachlogik
 - > Auswertung Zusammenhänge höherer Ordnung

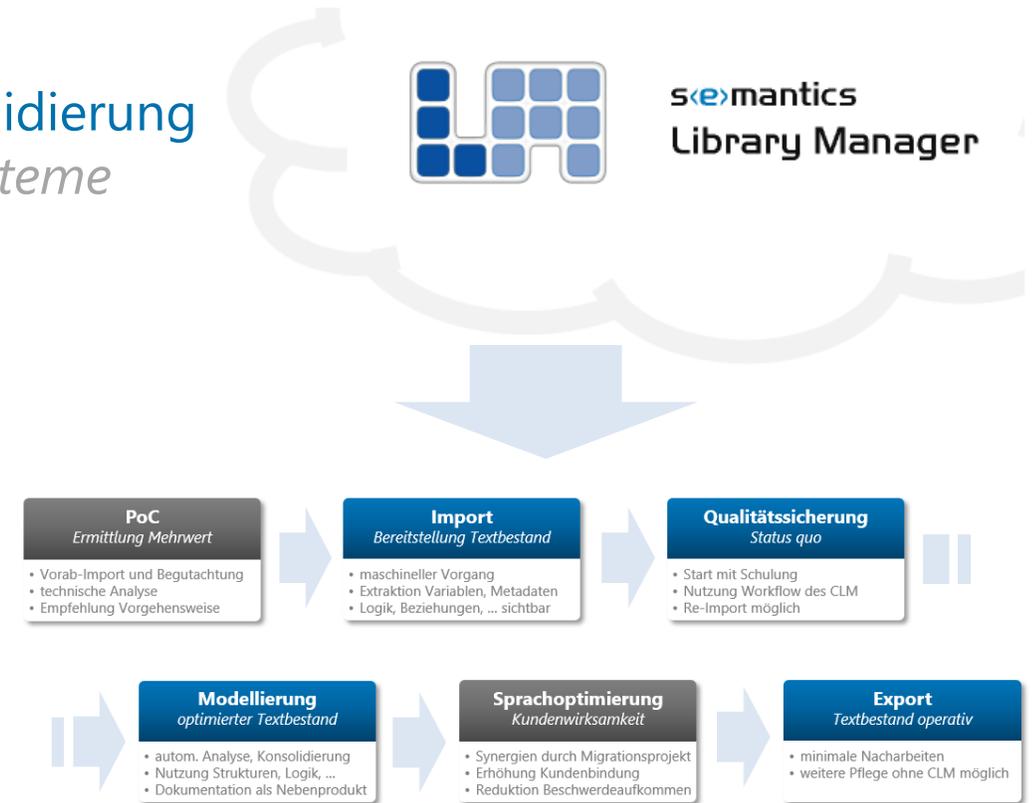


	doc1	doc2	doc3	doc4
angabe	2		1	
anschrift	1	1	1	1
arzt	3	1	2	1
bitte	1	1	1	
datum	1	1	1	1
fahrzeug				3
für	3	2	2	3
ja	2		1	
name	2	2	3	2
nein		1		1
ort	1	1	1	1
pensionsfond	2		2	
person	2	1		
seit		3		
tod			1	
überangebot				
unterschrift		1		
versorgungsb	1	1	1	
versorgungsv...	3	2		
vertragspartner	1	2		
wann	1	3	1	2
welche			3	2
werde			2	1
wurde	1	1	1	2
zahl			2	2



Status Quo – Anwendungen

- Standard Dienstleistung: Migration und/oder Konsolidierung
Migrations- und Konsolidierungstool für beliebige Systeme
- Bereitstellung „as a Service“: RemoteApps, Cloud
Keine Installation, Infrastruktur etc. notwendig
- Temporärer Einsatz zur Konsolidierung möglich
Dienstleistungen semantics inklusive



Fragen und Antworten



Dipl.-Ing. José Manuel de la Rosa Govantes

semantics Kommunikationsmanagement GmbH
Viktoriaallee 45
52066 Aachen

+49 241 89 49 89 29
j.delarosa@semantics.de
www.semantics.de

Besten Dank für Ihre Aufmerksamkeit!