

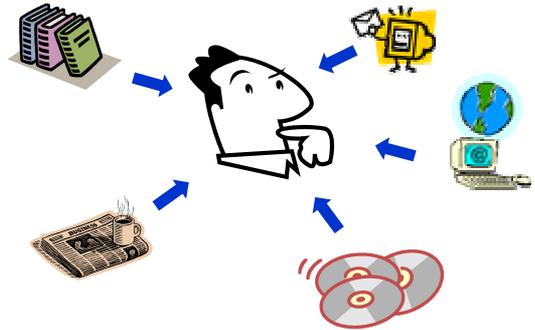
## Informationsextraktion aus Textdokumenten: Herausforderung, Technologie, Anwendungen

Dr. Roland Stuckardt  
IT-Beratung • Sprachtechnologie  
Frankfurt am Main

roland@stuckardt.de  
www.stuckardt.de

Veranstaltungsreihe des VDI – Frankfurt/Darmstadt, AK Kommunikationstechnik

## Das informationssuchende Individuum im 21. Jahrhundert



2

## Überflutung mit digitalen Texten

- Die Wissensproduktion nimmt stetig zu
- Der **textuellen** Informationsvermittlung kommt nach wie vor zentrale Bedeutung zu.
- Viele Texte liegen (auch) **als Computerdatei** bzw. **online** vor:
  - Internet: WWW, E-Mail, News
  - CD-ROM-Ausgaben klassischer Printmedien
  - digitale Bibliotheken
  - ...

3

## Software-Lösung für Informationssuchende

- informationssuchende Individuen verlangen nach
  - Informationen, die auf Ihre persönlichen Bedürfnisse zugeschnitten sind
  - effektivem Schutz vor informationeller Überflutung

... und somit nach **Software-Lösungen zur inhaltsorientierten Erschließung von Texten**

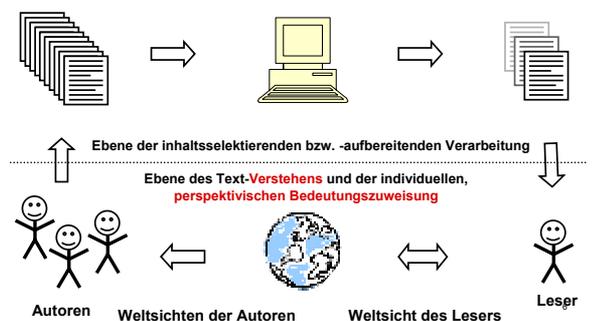
4

## Gliederung

- 1) Der Mensch als Informationssuchendes Individuum:
  - Herausforderung: Software-Lösungen zur Bewältigung der (textuellen) Informationsflut
- 2) Basisszenarien der inhaltsorientierten Texterschließung:
  - Text Retrieval
  - Textkategorisierung
  - Textzusammenfassung
  - **Informationsextraktion**
- 3) Informationsextraktion: Basisszenario und Basistechnologie
  - Aufgabenspezifikation und Gütekriterien
  - Leistungsfähigkeit aktueller Technologie
  - Software-Architektur und Inhaltserschließungstechniken
  - Beispiel: das FASTUS-System
- 4) Anwendungsszenario: Telefonüberwachung
  - IE-Technologie im Workflow einer forensischen Anwendung

5

## Rolle von Software-Lösungen für die inhaltsorientierte Erschließung von Texten



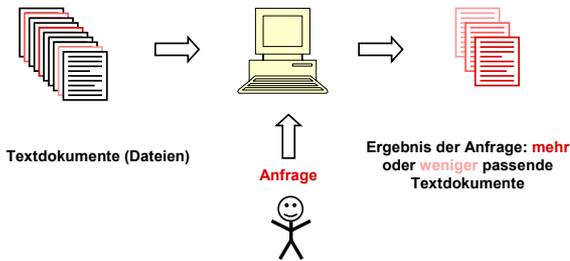
## Rolle von Software-Lösungen für die inhaltsorientierte Erschließung von Texten

- **Computer** selektiert und kategorisiert Texte bzw. strukturiert Textinhalte auf der Grundlage **bedeutungserhaltender** Operationen
- **Mensch versteht** Texte bzw. aufbereitete Inhalte vor dem Hintergrund seines perspektivischen Weltbezugs und **weist individuell Bedeutung** zu

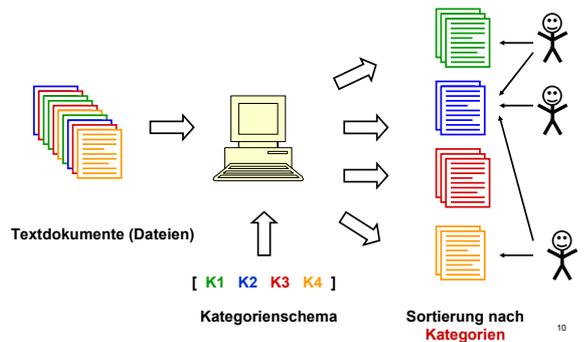
## Basisszenarien vs. Anwendungsszenarien

- **Basisszenarien:**
  - Finden **relevanter Texte** in der Informationsflut
  - **Ordnung** von Texten gemäß **Relevanzgrad**
  - **Kategorisierung** von Texten nach inhaltlicher Ähnlichkeit
  - **Zusammenfassung** von Texten
  - **Strukturierung** von Textinhalten
- **Anwendungsszenarien:**
  - Unterstützung von **Workflows**
  - Verfügbarkeit geeigneter **Benutzeroberflächen**
  - Abdeckung der jeweils relevanten **Textdateiformate**

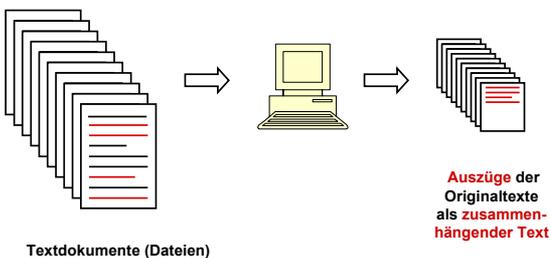
## Basisszenario: Text Retrieval



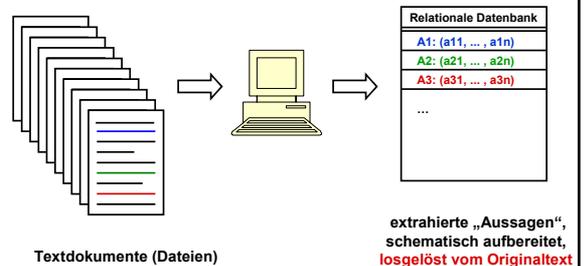
## Basisszenario: Textkategorisierung



## Basisszenario: Textzusammenfassung



## Basisszenario: Informationsextraktion



## Basisszenarien – Resümee:

	Ergebnis	Techniken	Status
Text Retrieval	relevanzsortierte Auswahl an Originaltexten	Analyse auf Wortebene, Statistik	Anwendungsreife
Textkategorisierung	Originaltexte kategorisiert	Wortebene, Statistik, Neuronetze	Anwendungsreife
Textzusammenfassung	Kernaussagen als zusammenhängender Text	inhaltliche Tiefenanalyse	Forschungsgegenstand
Informationsextraktion	Kernaussagen schematisch	<b>beschränkte</b> Analyse v. Satz- u. Textbedeutung	Anwendungsreife

13

## Fokus auf: Informationsextraktion

Behrens, Peter, \*1868 in Hamburg, +1940 in Berlin. Behrens entwickelte als einer der ersten Architekten des 20. Jahrhunderts eine architektonische Konzeption, die den Anforderungen der industrialisierten Zivilisation gerecht wurde - zu einer Zeit, in der die Gesellschaft noch in archaischen Vorstellungen dachte, gleichzeitig aber blind auf die überwältigenden Fortschritte der Technik vertraute. Behrens stand am Beginn der modernen Architektur in Deutschland, auf die er zwischen 1900 und 1914 einen entscheidenden Einfluss ausübte.

Ziel: Extraktion von Aussagen über künstlerische Aktivitäten, Schema:

[ **Künstler (K)**, **erschafft (S)**, **Objekt (O)** ]

Relationaler Charakter:

**entwickeln(Peter Behrens, architektonische Konzeption)**

14

## Warum das Problem schwierig ist

- sprachliche Ausdrucksvielfalt

Behrens errichtete die Turbinenhalle der AEG

Die Turbinenhalle der AEG wurde von Behrens errichtet

Die Errichtung der Turbinenhalle durch Behrens erfolgte 1911.

Die Turbinenhalle steht in Berlin. Behrens errichtete >sie< 1911.

15

## Eingrenzung der Aufgabenstellung „Textuelle Informationsextraktion“

### Message Understanding Competitions (MUC), US - Department of Defense / DARPA, 1989-98:

- Definition der zu extrahierenden Inhalte (sog. **Tasks**)
- Definition entsprechender **Gütekriterien**, um Systeme **evaluieren** und vergleichen zu können
- durch den Menschen erstellte **Referenzdaten** dienen als Vergleichsmaßstab

16

## Named Entity Task

Behrens, Peter, \*1868 in Hamburg, +1940 in Berlin. Behrens entwickelte als einer der ersten Architekten des 20. Jahrhunderts eine architektonische Konzeption, die den Anforderungen der industrialisierten Zivilisation gerecht wurde - zu einer Zeit, in der die Gesellschaft noch in archaischen Vorstellungen dachte, gleichzeitig aber blind auf die überwältigenden Fortschritte der Technik vertraute. Behrens stand am Beginn der modernen Architektur in Deutschland, auf die er zwischen 1900 und 1914 einen entscheidenden Einfluss ausübte.

### Ergebnis:

Behrens, Peter	1, 14
Hamburg	27, 32
Berlin	44, 49
Behrens	51, 57
Behrens	399, 405

17

## Scenario Template Task = Kernaufgabe der Informationsextraktion

Behrens, Peter, \*1868 in Hamburg, +1940 in Berlin. Behrens entwickelte als einer der ersten Architekten des 20. Jahrhunderts eine architektonische Konzeption, die den Anforderungen der industrialisierten Zivilisation gerecht wurde - zu einer Zeit, in der die Gesellschaft noch in archaischen Vorstellungen dachte, gleichzeitig aber blind auf die überwältigenden Fortschritte der Technik vertraute. Behrens stand am Beginn der modernen Architektur in Deutschland, auf die er zwischen 1900 und 1914 einen entscheidenden Einfluss ausübte.

### Ergebnis:

[ Schaffens-Akt: entwickeln, 60, 70  
Künstler: Behrens, 51, 57  
Gegenstand: architekton. Konzeption, 131, 157 ]

[ Schaffens-Akt: Einfluss ausüben, 519, 534  
Künstler: Behrens, 399, 405  
Gegenstand: moderne Architektur..., 427, 461 ]

18

# Coreference Task

Behrens, Peter, \*1868 in Hamburg, +1940 in Berlin. Behrens entwickelte als einer der ersten Architekten des 20. Jahrhunderts eine architektonische Konzeption, die den Anforderungen der industrialisierten Zivilisation gerecht wurde - zu einer Zeit, in der die Gesellschaft noch in archaischen Vorstellungen dachte, gleichzeitig aber blind auf die überwältigenden Fortschritte der Technik vertraute. Behrens stand am Beginn der modernen Architektur in Deutschland, auf die er zwischen 1900 und 1914 einen entscheidenden Einfluss ausübte.

Ergebnis: Koreferenz-Klassen

architekton. Konzeption  
die

Zeit  
die

moderne Architektur in ...  
die

Behrens, Peter  
Behrens  
e. d. ersten Architekten ...  
Behrens  
er

# Gütekriterien

Die Leistung der Softwaresysteme wird in Bezug auf **Qualität** und **Quantität** evaluiert:

$$Precision = \frac{\#(System \cap Referenz)}{\#System}$$

$$Recall = \frac{\#(System \cap Referenz)}{\#Referenz}$$

# Qualität, Quantität für Named Entity Task

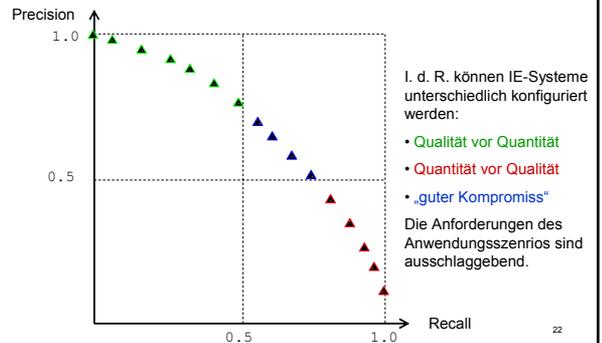
System	Referenz
Turbinenhalle	30,43
AEG	52,55
Behrens	71,78
AEG	52,55
Peter Behrens	65,78
Berlin	96,103
Hennigsdorf	114,125



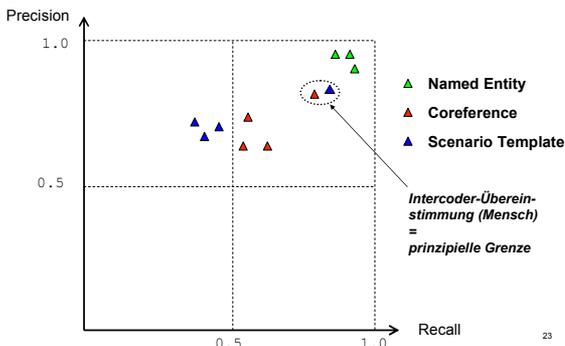
Precision = 2/3

Recall = 2/4

# Austauschverhältnis Qualität, Quantität



# Ergebnisse aktueller Technologie



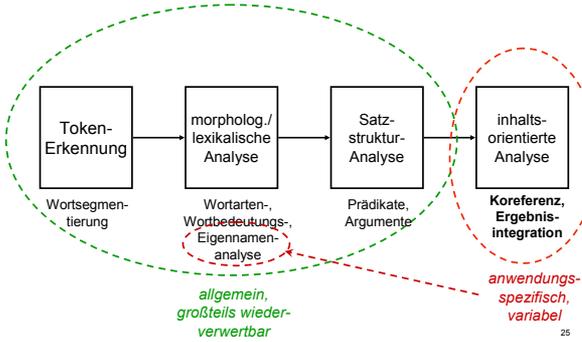
# Scenario Template Task, MUC-3: Terroraktivitäten in Lateinamerika

BOGOTA, 7 SEP 89 (INRAVISION TELEVISION CADENA 1) -- [REPORT] MARIBEL OSORIO [TEXT] MEDELLIN CONTINUES TO LIVE THROUGH A WAVE OF TERROR. FOLLOWING LAST NIGHT'S ATTACK ON A BANK, WHICH CAUSED A LOT OF DAMAGE, A LOAD OF DYNAMITE WAS HURLED AGAINST A POLICE STATION. FORTUNATELY NO ONE WAS HURT. HOWEVER, AT APPROXIMATELY 1700 TODAY A BOMB EXPLODED INSIDE A FAST-FOOD RESTAURANT. A MEDIUM-SIZED BOMB EXPLODED SHORTLY BEFORE 1700 AT THE PRESTO INSTALLATIONS LOCATED ON [WORDS INDISTINCT] AND PLAYA AVENUE. [...]

... sehr komplex, bis zu 18 (!) Attribute pro Vorfall - u. a.:

- DATE OF INCIDENT
- TYPE OF INCIDENT
- INSTRUMENT: TYPE (S)
- LOCATION OF INCIDENT
- PERPETRATOR: ID OF INDIV(S)
- PERPETRATOR: ID OF ORG(S)

# Architektur von IE-Systemen



25

# Morphologische Analyse

Künste	->	Kunst
Glasmalerei	->	Glas+Malerei
eingeführt	->	einführen
fürhte ... ein	->	einführen

- fürs Deutsche: **schwer**
  - zusammengesetzte Nomen
  - deshalb jedoch gerade essenziell
  - Softwarelösung: GerTwoL, LingSoft Helsinki
- fürs Englische: **leicht**
  - vollständige Auflistung aller Formen machbar

26

# Wortarten, Wortbedeutungen

- in Fällen von Mehrdeutigkeit ist der **Kontext** entscheidend:

schönen	Wortart: <b>Verb</b> oder <b>Adjektiv</b> ?
An den Zahlen gibt es nichts zu <b>schönen</b> .	
Am Morgen danach hatte er einen <b>schönen</b> Kater.	
Kater	Bedeutung: <b>Tier</b> oder <b>Brummschädel</b> ?
Ein schwarzer <b>Kater</b> kreuzte Stefans Weg.	
Am Morgen danach hatte er einen <b>schönen</b> <b>Kater</b> .	

27

# Statistische Wortart-Disambiguierung: N-Gram-Modell des Kontexts

- z. B. N=3, Gruppen von 3 Wörtern anschauen (Trigramme):  
An den Zahlen gibt es **nichts zu schönen**.  
 $Wsk(„schönen_V“ \text{ nach } „nichts \text{ zu“}) = 0.03$   
 $Wsk(„schönen_Adj“ \text{ nach } „nichts \text{ zu“}) = 0.0005$   
Am Morgen danach hatte **er einen schönen** Kater.  
 $Wsk(„schönen_V“ \text{ nach } „er \text{ einen“}) = 0.00001$   
 $Wsk(„schönen_Adj“ \text{ nach } „er \text{ einen“}) = 0.05$
- Grundgedanke: Anwendung **partieller** Kontextinformation  
betrachtetes Wort:  $w_k$   
bisher gelesener Text:  $T = w_1 w_2 \dots w_{k-1}$  = **Kontext**  
es werden aber nur die vorherigen N-1 Wörter betrachtet
- Kontextwahrscheinlichkeiten lassen sich per statistischer Analyse großer Textsammlungen **automatisch** ermitteln

28

# Eigennamenerkennung

- wichtiger Bestandteil der **anwendungsspezifischen** Analyse
- vollständige Auflistung oft unmöglich
- **Beispiel: Firmennamen** in Wirtschaftstexten
  - Neugründungen
  - Fusionen
  - Änderung der Rechtsform



- **Erkennungsregeln:**  
Firma := Token<sup>+</sup> { „AG“ | „GmbH“ | „KG“ | ... | „OHG“ }
- nicht immer derart einfach – **Produktamen:**  
„I can't believe it's not butter“ (amerikan. Margarine)

29

# Analyse der Satzstruktur

- = **syntaktische** Analyse  
**Peter Behrens**, der wie Gropius ein berühmter Architekt der Moderne war, **errichtete die Turbinenhalle der AEG**.  
**Die Turbinenhalle der AEG** wurde von dem berühmten Architekten 1911 **erbaut**.
- erkennt
  - eingeschobenen Relativsatz
  - Aktiv- und Passivsätze
  - ...
- bringt **Prädikate** (Verben) und zugehörige **Argumente** (Subjekt, Objekte, ...) zusammen

30

## Koreferenz- bzw. Pronomenresolution

- Argumente können **pronominal** realisiert sein:

errichten ( Subjekt: **er**,  
Objekt: Turbinenhalle der AEG)

- Ziel: Ermittlung informativerer Ausdrücke, die dasselbe Objekt bezeichnen (**Koreferenz**)

errichten ( Subjekt: **Peter Behrens**,  
Objekt: Turbinenhalle der AEG)

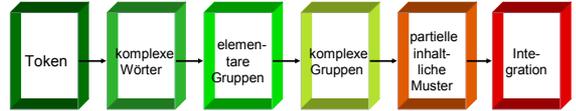
- Wie?** Hierzu später mehr ...

31

## Das FASTUS-System

- Kaskade sog. **Endlicher Automaten**:
  - einheitliches Analysemodell für (fast) alle Phasen
  - seichte Satzstruktur-Analyse
  - schnell, robust, relativ leicht zu pflegen und zu portieren

- Sechs Phasen**:



32

## Die Phasen von FASTUS

- Beispiel: Analyse von **Wirtschaftsmeldungen**

Bridgestone Sports Co. said Friday it has set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be shipped to Japan. The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990 with production of 20,000 iron and metal wood clubs a month.

- Ziel: Aussagen über **Joint Ventures** extrahieren, Attribute

- beteiligte Firmen
- Joint-Venture-Firma, Eigenkapital
- Produkt(e)
- ab wann?

33

## Die Phasen von FASTUS

- Phase 1: Erkennung von Token

```

<Bridgestone> <Sports> <Co> <.> <said> <Friday> <it> <has>
<set> <up> <a> <joint> <venture> <in> <Taiwan> <with> <a>
<local> <concern> <and> <a> <Japanese> <trading> <house> <to>
<produce> <golf> <clubs> <to> <be> <shipped> <to> <Japan> <.>
  
```

- Phase 2: „Multiwörter“, Namen, Zahlen

```

<Bridgestone> <Sports> <Co> <.> <said> <Friday> <it> <has>
<set> <up> <a> <joint> <venture> <in> <Taiwan> <with> <a>
<local> <concern> <and> <a> <Japanese> <trading> <house> <to>
<produce> <golf> <clubs> <to> <be> <shipped> <to> <Japan> <.>
  
```

34

## Die Phasen von FASTUS

- Phase 3: elementare Gruppen

```

<Bridgestone> <Sports> <Co> <.> <said> <Friday> <it> <has> <set>
<up> <a> <joint> <venture> <in> <Taiwan> <with> <a> <local>
<concern> <and> <a> <Japanese> <trading> <house> <to> <produce>
<golf> <clubs> <to> <be> <shipped> <to> <Japan> <.>
  
```

Unternehmen  
Ortsangaben  
Verbgruppen  
Nominalgruppen

Präpositionen  
Konjunktionen

35

## Die Phasen von FASTUS

- Erkennung von Adjektiv-Gruppen / -Sequenzen

Grammatikregeln als *reguläre Ausdrücke*:

```

AdjP -> Ordinal
| [ { Q-er | Q-est } ] { Adj | Vparticle }+
| N<sing,!Time-NP> - Vparticle
| Number [-] { month | day | year } [-] old
  
```

```

Adjs -> AdjP [ [ , | [,] Conj ] { AdjP | Vpart } ]*
  
```

Beispiele:

```

<fifth> <Symphony>
<most> <hungry> <guys>
<five> <-> <year> <-> <old> <boy>
<big> <,> <bad>, <and> <ugly> <car>
  
```

36

## Die Phasen von FASTUS

- Phase 4: Erkennung komplexer Nominal- und Verbgruppen

z.B. koordinierte Nominalgruppen:

```
... <The> <joint> <venture> <,> <Bridgestone> <Sports>
<Taiwan> <Co> <.> <,> <capitalized> <at> <20> <million>
<new> <Taiwan> <dollars>, <will> <start> <production> <in>
<January> <1990> <with> <production> <of> <20,000> <iron>
<and> <metal> <wood> <clubs> <a> <month> <.>
```

37

## Die Phasen von FASTUS

- Phase 5: partielle inhaltliche Muster

(1)	<b>Relationship:</b>	<b>TIE-UP</b>
	<b>Entities:</b>	„Bridgestone Sports Co.“ „a local concern“ „a Japanese trading house“
(2)	<b>Activity:</b>	<b>PRODUCTION</b>
	<b>Product:</b>	„golf clubs“
(3)	<b>Relationship:</b>	<b>TIE-UP</b>
	<b>JV Company:</b>	„Bridgestone Sports Taiwan Co.“
	<b>Amount:</b>	NT\$2000000
(4)	<b>Activity:</b>	<b>PRODUCTION</b>
	<b>Company:</b>	„Bridgestone Sports Taiwan Co.“
	<b>Start Date:</b>	DURING: January 1990
(5)	<b>Activity:</b>	<b>PRODUCTION</b>
	<b>Product:</b>	„iron and metal wood clubs“

38

## Die Phasen von FASTUS

- Phase 6: Integration in IE-Zielstruktur

### TIE-UP-1:

**Relationship:** TIE-UP  
**Entities:** „Bridgestone Sports Co.“  
 „a local concern“  
 „a Japanese trading house“  
**JV Company:** „Bridgestone Sports Taiwan Co.“  
**Activity:** ACTIVITY-1  
**Amount:** NT\$2000000

### ACTIVITY-1:

**Company:** „Bridgestone Sports Taiwan Co.“  
**Product:** „iron and metal wood clubs“  
**Start Date:** DURING: January 1990

39

## Koreferenz- bzw. Pronomenresolution

Behrens ist berühmt. Er baute die Turbinenhalle. Der Architekt ...

- Facetten:

- Ermittlung der Klassen **koreferenter** Ausdrücke
- IE-Ziel: Ermittlung möglichst **informativer** Ausdrücke
- aus algorithmischer Sicht: zwei Schritte
  - Identifikation von Ausdrücken, die Objekte referenzieren (<Behrens>, <er>, <die Turbinenhalle>, <der Architekt>)
  - Identifikation koreferenter **Antezedens**-Ausdrücke (insbesondere für Pronomen)
 

(<Behrens> ← <er>, <Behrens> ← <der Architekt>)

## Strategien zur Koreferenzresolution

- Ermittlung koreferenter Antezedenten für **Pronomen**:

- Restriktion: Kongruenz in Numerus und Genus  
 Behrens sagte **seiner Kollegin**, dass **er sie** schätze.
- Restriktion: syntaktisch-konfigurationale Zugänglichkeit  
 Behrens verlangt, dass **der Friseur sich** rasiert.  
 Behrens verlangt, dass **der Friseur ihn** rasiert.
- Präferenzkriterium: Bevorzugung des Subjekts  
 Der **Friseur** betrat **den Salon**. **Er** bediente den Kunden.
- ...

41

## Anwendungsfall: Telefonüberwachung

Wie können aufgezeichnete Gespräche sinnvoll inhaltlich exploriert werden?

- Es liegen zwei Kategorien von Informationen vor
  - Strukturierte Daten:** beteiligte Anschlüsse/ Teilnehmer, Datum/Uhrzeit, Gesprächsdauer, Provider
  - Audiodaten:** das eigentliche Gespräch
- sehr große Menge an Information – über 1000 Gespräche pro Maßnahme sind die Regel

42

## Audiodaten: Analyse problematisch

Fast immer in fremden Sprachen

- **mehrsprachiges** Anwendungsszenario
- **Verschriftung durch Übersetzer** zwingend

- Auch verschriftete Gespräche sind noch nicht ausreichend strukturiert

⇒ **Informationsextraktion**

43

## Ziel: Strukturierte Repräsentation der Gesprächstexte

Identifizierung bestimmter Klassen besonders relevanter Objekte:

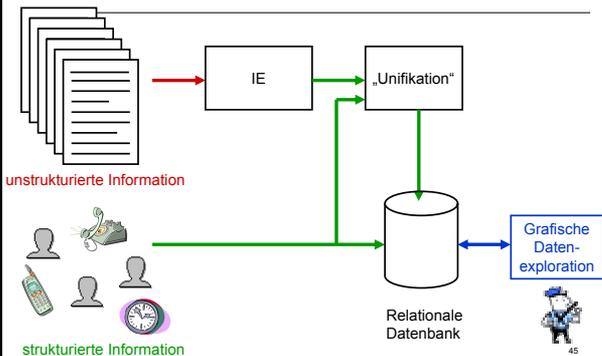
- Personen, Firmen, Orte, Datum, Waren, Transaktionen etc

Überführung erkannter Objekte in eine relationale Datenbank:

- **Zusammenführung mit strukturierten Daten**
- einheitliches grafisches Explorationswerkzeug für beide Informationskategorien

44

## Workflow der IE-Anwendung



45

## Zusammenfassung

- Informationsextraktion stellt eine zentrale **Herausforderung** im Zeitalter der digitalen Informationsflut dar
- flexible Software-**Basistechnologie** - zuschneidbar auf unterschiedliche Anwendungsszenarien - ist verfügbar
- Zur Entwicklung von **Anwendungen** textueller Informationsextraktion bedarf es Expertenwissen in den Gebieten
  - Informatik / Software Engineering
  - Linguistik
  - Computerlinguistik / Natural Language Engineering
- **Projektgeschäft** – (noch) kein fertiges Produkt verfügbar
- **Kompetente Beratung** in der Analyse potenzieller Anwendungsszenarien, in der Auswahl geeigneter Technologien und in der Projektdurchführung ist essenziell!

46

## Folien mit ergänzender Information

47

## Basisszenarien: unterscheidbar hinsichtlich

- **Input:** Dokumentenbestand typischerweise
  - statisch (weitgehend unveränderlich)?
  - dynamisch (fortwährend neue Dokumente)?
- **Output:** Welche Ergebnisse erhält der Nutzer?
  - Originaltexte?
  - Auszüge von Originaltexten?
  - strukturiert aufbereitete „Inhalte“?
- **Techniken** zur Inhaltserschließung?
- **Beispiele** möglicher Anwendungsszenarien?

48

## Basisszenario: Text Retrieval

- **Input:** sehr großer Bestand an Textdokumenten, größtenteils statisch
- **Output:** Auswahl an relevanten Textdokumenten, geordnet nach Relevanz
- **Anwender:** erhält vollständige Texte
- **Techniken:** Schlüsselwortsuche, statistische Häufigkeitsanalysen, Synonymlexika
- **Anwendungen:** Internet-Suchmaschinen, Suchwerkzeuge für CD-ROM-Textarchive

49

## Basisszenario: Textkategorisierung

- **Input:** große Anzahl von Textdokumenten, typischerweise dynamisch
- **Output:** Originaldokumente, geordnet nach Kategorien
- **Anwender:** erhält vollständige Texte
- **Techniken:** wie Information Retrieval; auch Neuronale Netze
- **Anwendungen:** Spamfilter; Verteiler für E-Mails, Pressemeldungen; Wissensmanagement in Unternehmen

50

## Basisszenario: Textzusammenfassung

- **Input:** als relevant bekannte, oft umfangreiche Textdokumente; statisch oder dynamisch
- **Output:** Kernaussagen der Originaldokumente als zusammenhängender Text
- **Anwender:** erhält Ausschnitte der Texte
- **Techniken:** inhaltliche Tiefenanalyse
- **Anwendungen:** nachgeschalteter Schritt in Text-Retrieval-Anwendungen, E-Mail-Verteilern

51

## Basisszenario: Informationsextraktion

- **Input:** größerer Bestand an Textdokumenten, statisch oder dynamisch
- **Output:** inhaltlich relevante „Aussagen“, schematisch aufbereitet (keine Textform)
- **Anwender:** erhält schematisch dargestellte Textinhalte (zunächst) losgelöst vom Originaltext
- **Techniken:** Analyse v. Satzstruktur u. Textbedeutung Erschließung relationaler inhaltlicher Strukturen
- **Anwendungen:** E-Mail-Antwortmanagement (CRM), forensische Analyse von Telefonatverschriftungen

52

## Methodologische Probleme der Entwicklung von IE-Software

- erschließungsrelevante „Aussagen“?
- unklare Erfolgskriterien
- mangelnde Vergleichbarkeit unterschiedlicher Verfahren



- bis in die 90er Jahre hinein wenig Fortschritt

53

## Stochastische Wortart-Disambiguierung: N-Gram-Modell des Kontexts

- z. B. N=3, Gruppen von 3 Wörtern anschauen (Trigramme):  
An den Zahlen gibt es nichts zu schönen.  

$$P(\text{schönen}_v | \text{nichts zu}) = 0.03$$

$$P(\text{schönen}_{Adj} | \text{nichts zu}) = 0.0005$$
 Am Morgen danach hatte er einen schönen Kater.  

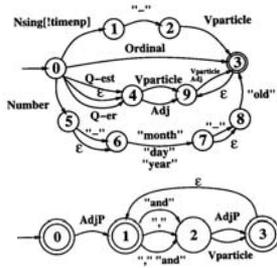
$$P(\text{schönen}_v | \text{er einen}) = 0.00001$$

$$P(\text{schönen}_{Adj} | \text{er einen}) = 0.05$$
- Grundgedanke: Anwendung partieller Kontextinformation  
 aktuell betrachtetes Wort:  $w_k$   
 bisher gelesener Text:  $T = w_1 w_2 \dots w_{k-1}$  = Kontext  
 Approximation der bedingten Wahrscheinlichkeit der Wortart (WA) von  $w_k$  per Begrenzung des berücksichtigten Kontexts:  

$$P(WA(w_k) | WA(w_1) \dots WA(w_{k-1})) \sim P(WA(w_k) | WA(w_{k-N+1}) \dots WA(w_{k-1}))$$
- bedingte Wahrscheinlichkeiten lassen sich per statistischer Analyse großer Textsammlungen automatisch ermitteln

## Die Phasen von FASTUS

- Erkennung von Adjektiv-Gruppen / -Sequenzen  
aus den Grammatikregeln erzeugte *Finite State Transducer*:



55

## Strategien zur Koreferenzresolution

- Strategien für **Nichtpronomen**?
- vielfältige anaphorische Relationen, nicht unbedingt Koreferenz  
Behrens hielt seinen **Wagen** an. Das **Kühlwasser** kochte.
- schwer zu
  - definieren / standardisieren
  - evaluieren
  - implementieren**
- Ansatz:
  - WordNet**: große Datenbank lexikalischer Beziehungen (*Teil-Ganzes, Element-von, Über- und Unterbegriffe ...*)
  - Suche in WordNet nach Beziehungen jenseits Koreferenz

## Strategien zur Koreferenzresolution

- Implementierung: **das ROSANA-System**
- ROSANA** := **RO**buste **S**yntaxbasierte **ANA**phern-Interpretation
- Evaluation in vier Teildisziplinen (Precision, Recall):
  - referenzierende Ausdrücke: (0.94,0.96)
  - unmittelbare Antezedenten:\* (0.71,0.71), (0.76,0.76)
  - informative Substitute**:\* (0.68,0.67), (0.66,0.66)
  - Koreferenz-Klassen (0.81,0.68)

\* = für Dritte-Person-Pronomen (Nichtpossessiva bzw. Possessiva)

57