Context is the Key:

# Context-aware Corpus Annotation Using Subjective Content Descriptions

Colloquium

## Felix Kuhr

Institute of Information Systems

University of Lübeck

March 24, 2022

- <u>Information retrieval (IR)</u> is the task of finding documents that are <u>relevant</u> to a user's need for information

- <u>Information retrieval (IR)</u> is the task of finding documents that are <u>relevant</u> to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:

- <u>Information retrieval (IR)</u> is the task of finding documents that are <u>relevant</u> to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
  - Compare words in a query with content of documents

- <u>Information retrieval (IR)</u> is the task of finding documents that are <u>relevant</u> to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
    - Compare words in a query with content of documents

An information retrieval system can be characterized by:

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

- <u>Information retrieval (IR)</u> is the task of finding documents that are <u>relevant</u> to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
    - Compare words in a query with content of documents

An information retrieval system can be characterized by:
- **Corpus** of documents

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

- <u>Information retrieval (IR)</u> is the task of finding documents that are <u>relevant</u> to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
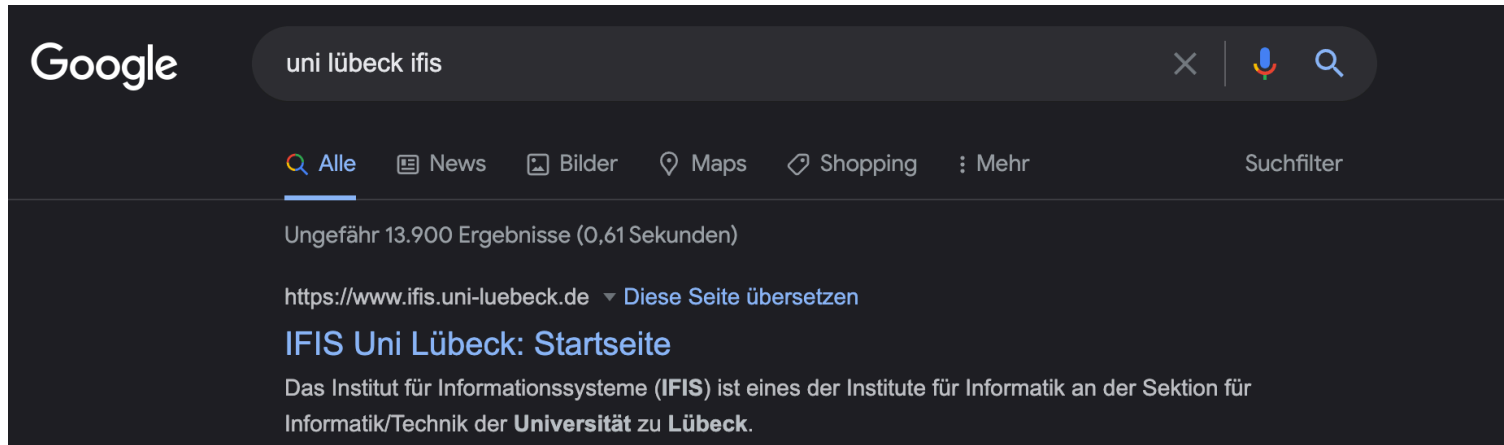  - Compare words in a query with content of documents

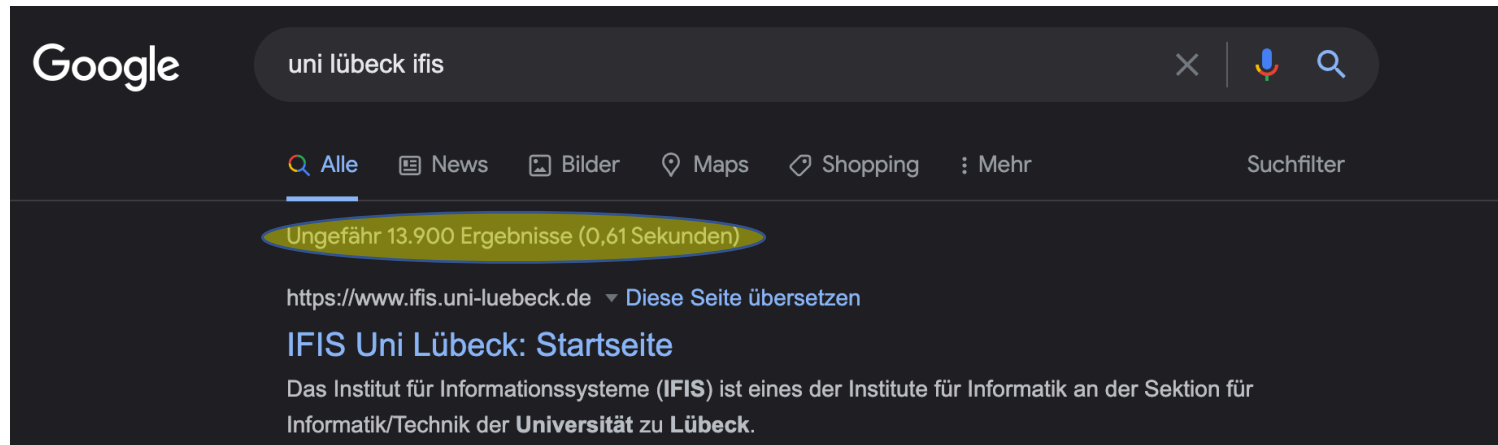An information retrieval system can be characterized by:
- **Corpus** of documents
- **Queries** posed in a query language

- <u>Information retrieval (IR)</u> is the task of finding documents that are <u>relevant</u> to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
  - Compare words in a query with content of documents

An information retrieval system can be characterized by:
- **Corpus** of documents
- **Queries** posed in a query language

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

- Information retrieval (IR) is the task of finding documents that are relevant to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
    - Compare words in a query with content of documents

An information retrieval system can be characterized by:

- **Corpus** of documents

- **Queries** posed in a query language

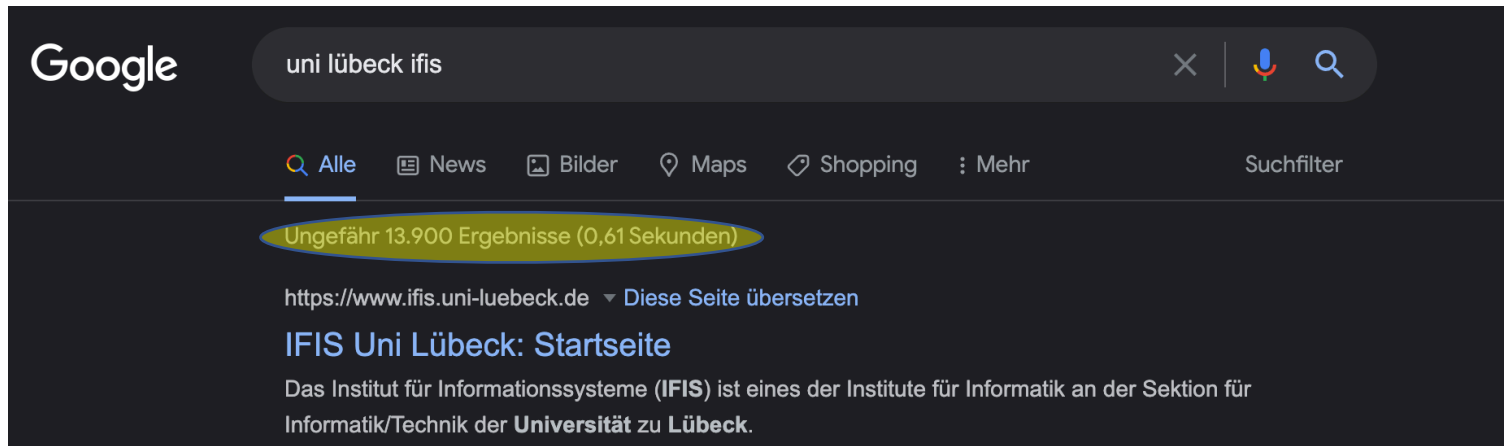UNIVERSITÄT ZU LÜBECK
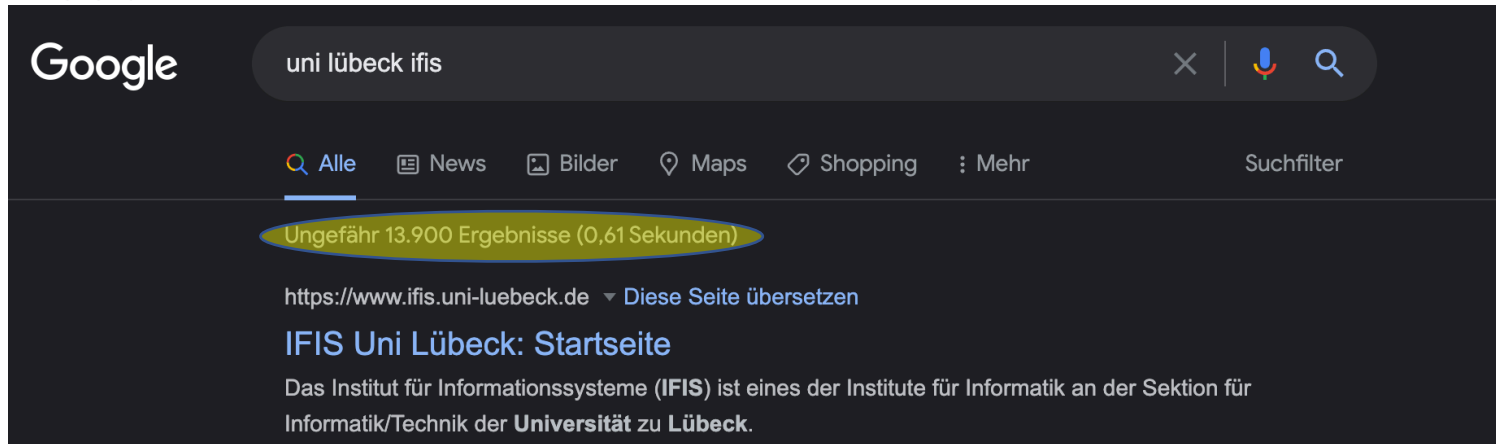INSTITUT FÜR INFORMATIONSSYSTEME

- Information retrieval (IR) is the task of finding documents that are relevant to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
    - Compare words in a query with content of documents

An information retrieval system can be characterized by:
- **Corpus** of documents
- **Queries** posed in a query language

Are *documents* in the result set relevant to the information need of a user?



Google

uni lübeck ifis

&#x1F50D; Alle    News    Bilder    Maps    Shopping    Mehr      Suchfilter

Ungefähr 13.900 Ergebnisse (0,61 Sekunden)

https://www.ifis.uni-luebeck.de   ▾ Diese Seite übersetzen

IFIS Uni Lübeck: Startseite

Das Institut für Informationssysteme (**IFIS**) ist eines der Institute für Informatik an der Sektion für Informatik/Technik der **Universität** zu **Lübeck**.

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

- Information retrieval (IR) is the task of finding documents that are relevant to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
  - Compare words in a query with content of documents

An information retrieval system can be characterized by:

- **Corpus** of documents
- **Queries** posed in a query language
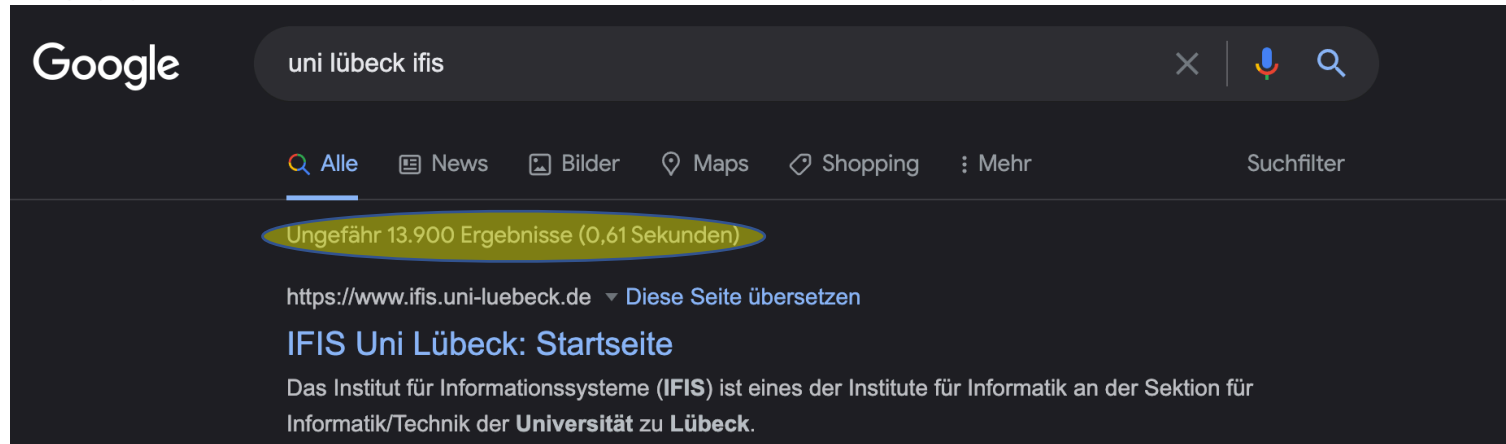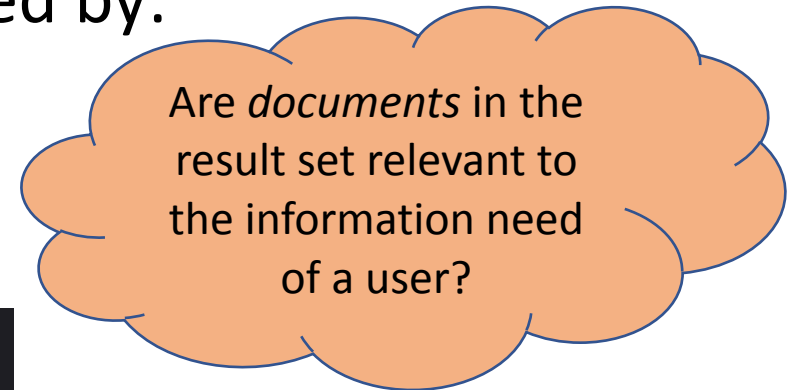- **Result** set of relevant documents

Are *documents* in the result set relevant to the information need of a user?



Google    uni lübeck ifis

Q Alle    News    Bilder    Maps    Shopping    Mehr        Suchfilter

Ungefähr 13.900 Ergebnisse (0,61 Sekunden)

https://www.ifis.uni-luebeck.de ▾ Diese Seite übersetzen
IFIS Uni Lübeck: Startseite
Das Institut für Informationssysteme (IFIS) ist eines der Institute für Informatik an der Sektion für Informatik/Technik der **Universität** zu **Lübeck**.

- Information retrieval (IR) is the task of finding documents that are relevant to a user's need for information
- Algorithms estimate relevance of displayed documents to searched queries:
  - Compare words in a query with content of documents

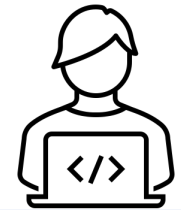An information retrieval system can be characterized by:
- **Corpus** of documents
- **Queries** posed in a query language
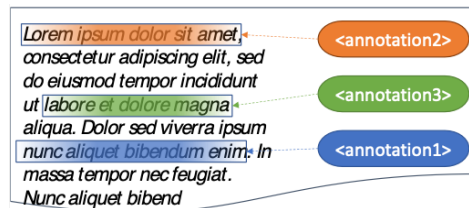- **Result** set of relevant documents
-

Are *documents* in the result set relevant to the information need of a user?

Google

uni lübeck ifis

Q Alle   📰 News   🖾 Bilder   ⦿ Maps   ⊘ Shopping   ⋮ Mehr       Suchfilter

Ungefähr 13.900 Ergebnisse (0,61 Sekunden)

https://www.ifis.uni-luebeck.de   ▾ Diese Seite übersetzen
IFIS Uni Lübeck: Startseite
Das Institut für Informationssysteme (**IFIS**) ist eines der Institute für Informatik an der Sektion für Informatik/Technik der **Universität** zu **Lübeck**.

# Annotation Systems

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Annotation Systems

Manual Annotation Systems

**+** Quality of annotations

User-centric annotations

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

<annotation2>
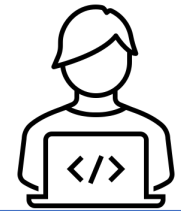<annotation3>
<annotation1>

# Annotation Systems

Manual Annotation Systems

**+**

Quality of annotations

User-centric annotations

**–**

Human annotation experts

High costs / time-consuming

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Annotation Systems

Manual Annotation Systems

Automatic Annotation Systems

**+** Quality of annotations

User-centric annotations

**+** Low costs

Fast annotation process

**−** Human annotation experts

High costs / time-consuming

# Annotation Systems


Manual Annotation Systems


Automatic Annotation Systems

**+** Quality of annotations

User-centric annotations

**+** Low costs

Fast annotation process

**−** Human annotation experts

High costs / time-consuming

**−** Quality of annotations

Missing explainability

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Annotation Systems


Manual Annotation Systems


Automatic Annotation Systems

**+** Quality of annotations

User-centric annotations

**+** Low costs

Fast annotation process

**−** Human annotation experts

High costs / time-consuming

**−** Quality of annotations

Missing explainability

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Context is the Key: Example (One Calais)*

January 6, 1984). Ivana became a naturalized United States citizen in 1988. By early 1990, Trump's troubled marriage to Ivana and affair with actress Marla Maples had been reported in the tabloid press. They were divorced in 1992.. Trump married his second wife, actress Marla Maples in 1993. They had one daughter together, Tiffany (born October 13, 1993). The couple were separated in 1997 and l                                                    n Slovene model Melania Knauss, who                                                     re married on January 22, 2005, at Be                                                   h, Florida. In 2006, Melania becan                                                    jave birth to their son, whom they                                                    y his doctor, Harold Bornstein M.D.,                                                   were in normal range. Trump says                                                     including marijuana. He also does not d                                                   olism. His BMI, according to his

Context-aware Corpus Annotation Using Subjective Content Descriptions

* https://permid.org/onecalaisViewer

# Context is the Key: Example (One Calais)*

January 6, 1984). Ivana became a naturalized United States citizen in 1988. By early 1990, Trump's troubled marriage to Ivana and affair with actress Marla Maples had been reported in the tabloid press. They were divorced in 1992.. Trump married his second wife, actress Marla Maples in 1993. They had one daughter together, Tiffany (born October 13, 1993). The couple were separated in 1997

and la...                                                    ...n Slovene model Melania Knauss,

who h...                                                    ...re married on January 22, 2005,

at Be...                                                    ...h, Florida. In 2006, Melania

becam...                                                    ...gave birth to their son, whom

they...                                                    ...y his doctor, Harold Bornstein

M.D.,...                                                    ...were in normal range. Trump

says t...                                                    ...ncluding marijuana. He also does

not d...                                                    ...olism. His BMI, according to his

**COMPANY**
**TIFFANY & CO.**

Relevance
**20%**

PermID
**4295905088** ☑

Continuous Relevance      3%
nationality               N/A
confidencelevel           0.891

Add meta data to text by linking extractable entities to external data

# Context is the Key: Example (One Calais)*

January 6, 1984). Ivana became a naturalized United States citizen in 1988. By early 1990, Trump's troubled marriage to Ivana and affair with actress Marla Maples had been reported in the tabloid press. They were divorced in 1992.. Trump married his second wife, actress Marla Maples in 1993. They had one daughter together, Tiffany (born October 13, 1993). The couple were separated in 1997

and l[...]n Slovene model Melania Knauss,

who [...]re married on January 22, 2005,

at Be[...]h, Florida. In 2006, Melania

beca[...]gave birth to their son, whom

they [...]his doctor, Harold Bornstein

M.D.,[...]were in normal range. Trump

says [...]ncluding marijuana. He also does

not d[...]olism. His BMI, according to his

Which annotation fulfill a user's information needs?

Add meta data to text by linking extractable entities to external data

Text related with Company Tiffany & Co.

Persons with names also used for companies

# Context is the Key: Example (One Calais)*

January 6, 1984). Ivana became a naturalized United States citizen in 1988. By early 1990, Trump's troubled marriage to Ivana and affair with actress Marla Maples had been reported in the tabloid press. They were divorced in 1992.. Trump married his second wife, actress Marla Maples in 1993. They had one daughter together, Tiffany (born October 13, 1993). The couple were separated in 1997

Slovene model Melania Knauss,

married on January 22, 2005,

Florida. In 2006, Melania

gave birth to their son, whom

his doctor, Harold Bornstein

were in normal range. Trump

including marijuana. He also does

olism. His BMI, according to his

Donald Trump purchased the building's air rights of Tiffany & Co. flagship store for $5 million in 1979 while he was developing the neighboring Trump Tower. Trump later named his daughter Tiffany Trump after the location

Add meta data to text by linking extractable entities to external data

Text related with Company Tiffany & Co.

Persons with names also used for companies

* https://permid.org/onecalaisViewer

# Context is the Key: Example (One Calais)*

January 6, 1984). Ivana became a naturalized United States citizen in 1988. By early 1990, Trump's troubled marriage to Ivana and affair with actress Marla Maples had been reported in the tabloid press. They were divorced in 1992.. Trump married his second wife, actress Marla Maples in 1993. They had one daughter together, Tiffany (born October 13, 1993). The couple were separated in 1997

and l[...] Slovene model Melania Knauss,

who [...] re married on January 22, 2005,

at Be[...]h, Florida. In 2006, Melania

beca[...]gave birth to their son, whom

they [...] his doctor, Harold Bornstein

M.D.,[...] were in normal range. Trump

says [...]ncluding marijuana. He also does

not d[...]olism. His BMI, according to his

Add meta data to text by linking extractable entities to external data

Text related with Company Tiffany & Co.

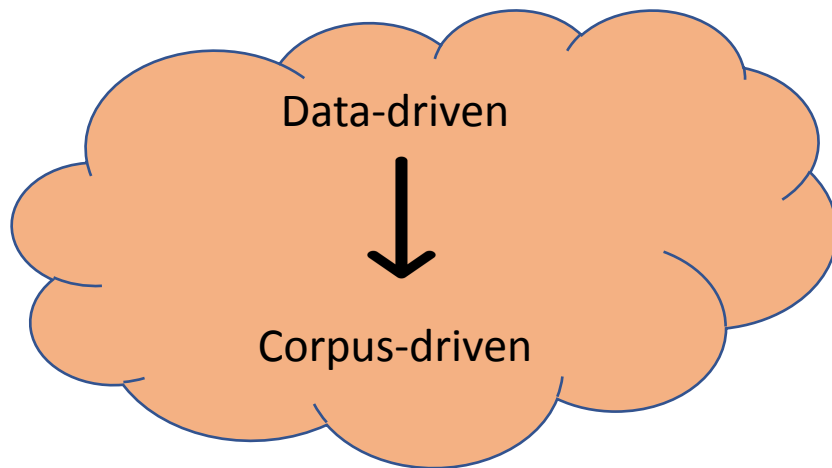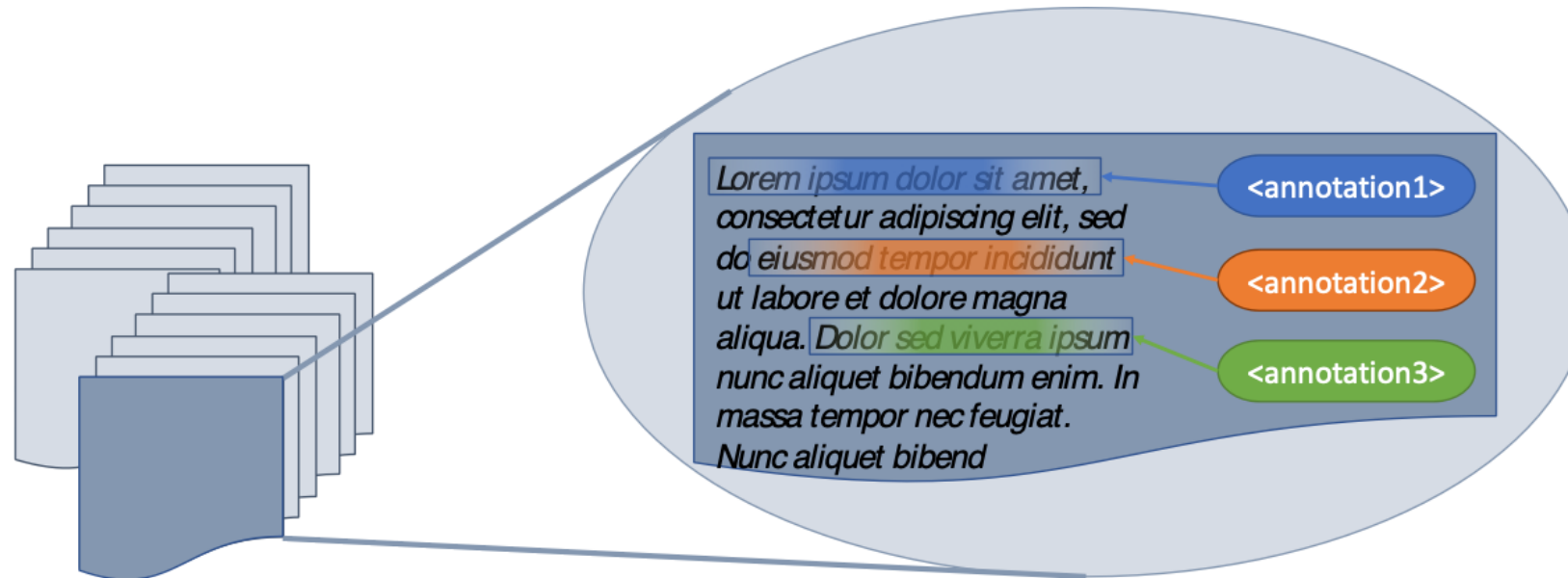Persons with names also used for companies

# Context is the Key: Example (One Calais)*

January 6, 1984). Ivana became a naturalized United States citizen in 1988. By early 1990, Trump's troubled marriage to Ivana and affair with actress Marla Maples had been reported in the tabloid press. They were divorced in 1992.. Trump married his second wife, actress Marla Maples in 1993. They had one daughter together, Tiffany (born October 13, 1993). The couple were separated in 1997

and l... ... n Slovene model Melania Knauss,

who ... ...re married on January 22, 2005,

at Be... ...h, Florida. In 2006, Melania

becam... ...gave birth to their son, whom

they ... ...y his doctor, Harold Bornstein

M.D., ... ...were in normal range. Trump

says ... ...ncluding marijuana. He also does

not d... ...olism. His BMI, according to his

Data-driven

↓

Corpus-driven

Add meta data to text by linking extractable entities to external data

Text related with Company Tiffany & Co.

Persons with names also used for companies

Context-aware Corpus Annotation Using Subjective Content Descriptions

* https://permid.org/onecalaisViewer

# Subjective Content Descriptions (SCDs)
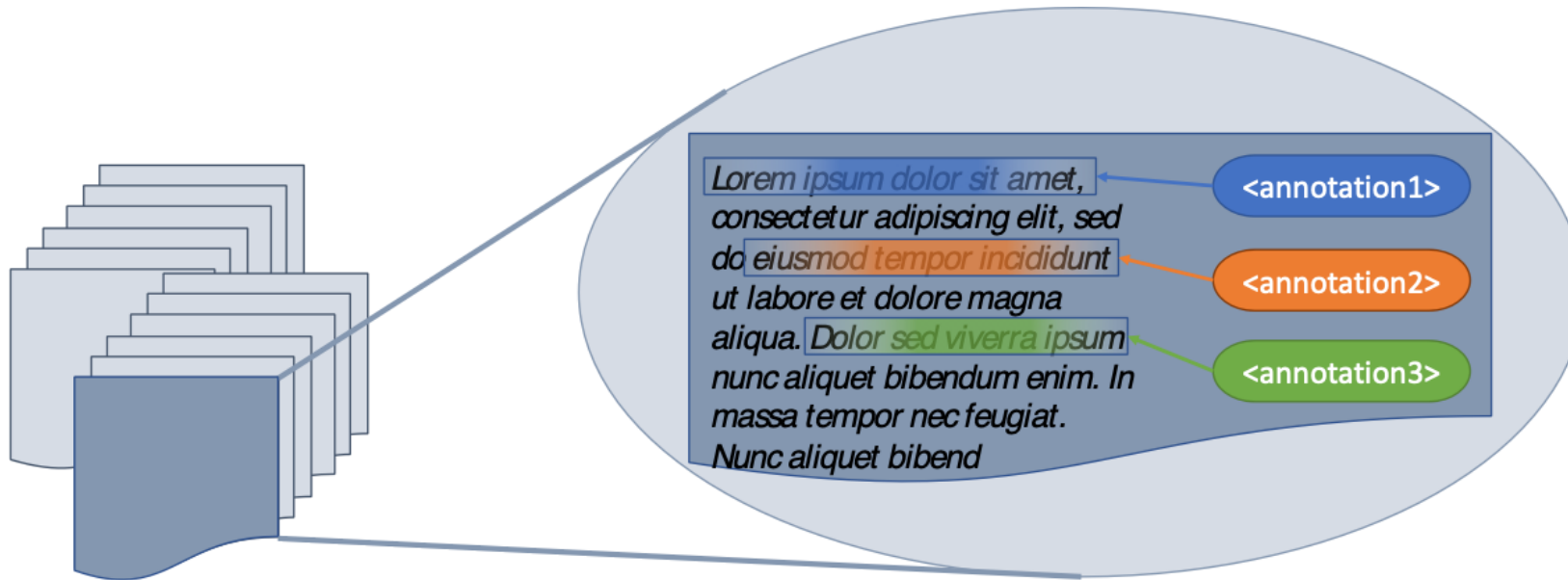
UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Subjective Content Descriptions (SCDs)
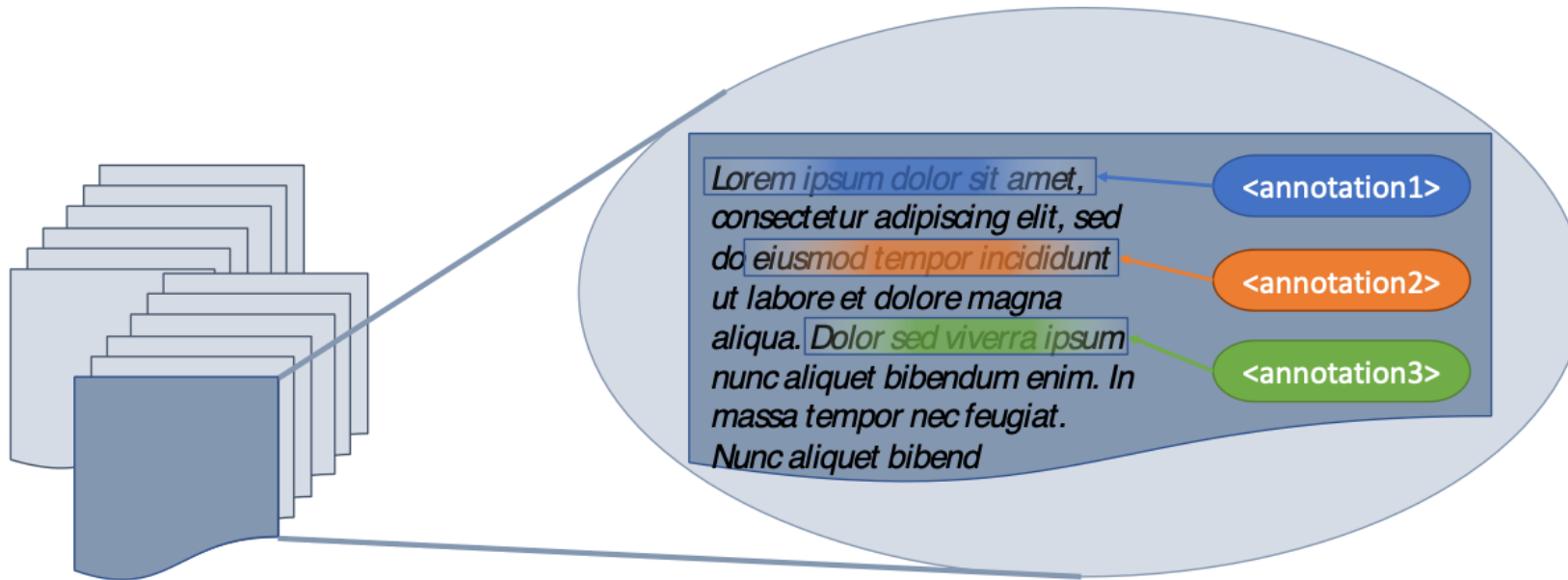
- … represent additional data for a document

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document

- … are associated with a text span (window) in a document

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document

- … are associated with a text span (window) in a document
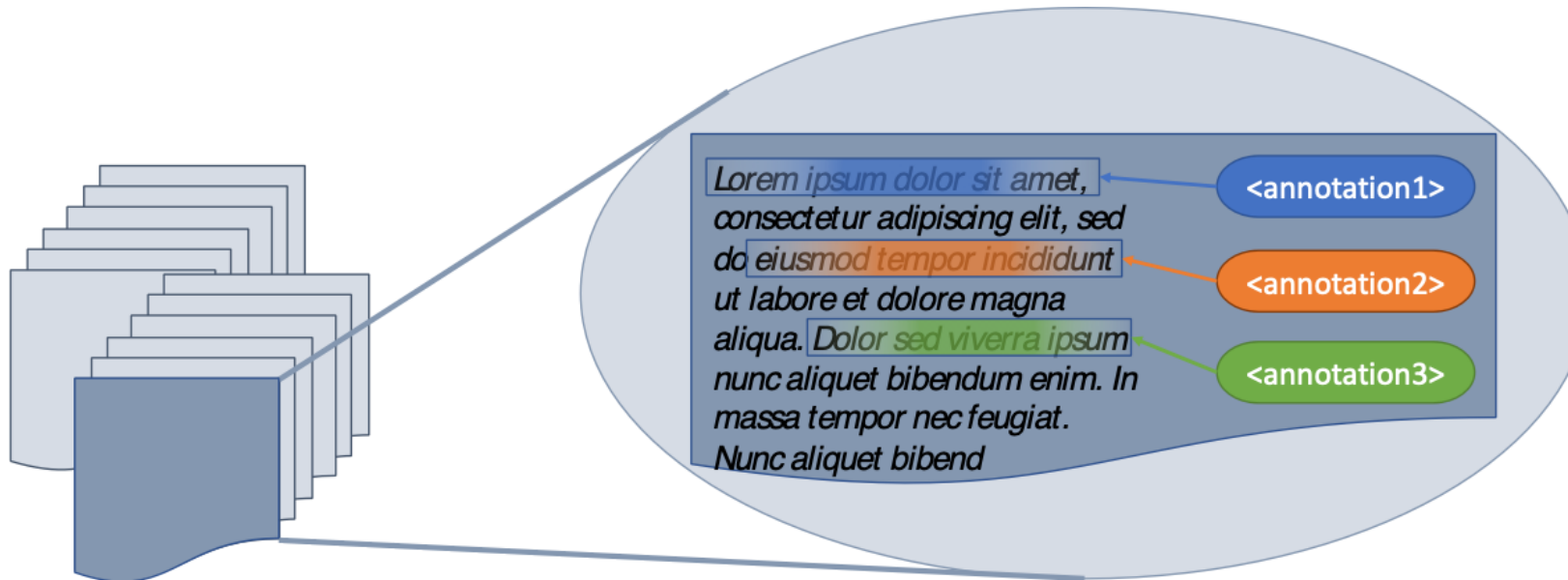
Different Typs of SCDs:

- Additional definitions



UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

Context-aware Corpus Annotation Using Subjective Content Descriptions

5

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document

- … are associated with a text span (window) in a document

Different Typs of SCDs:

- Additional definitions

- Links to external sources

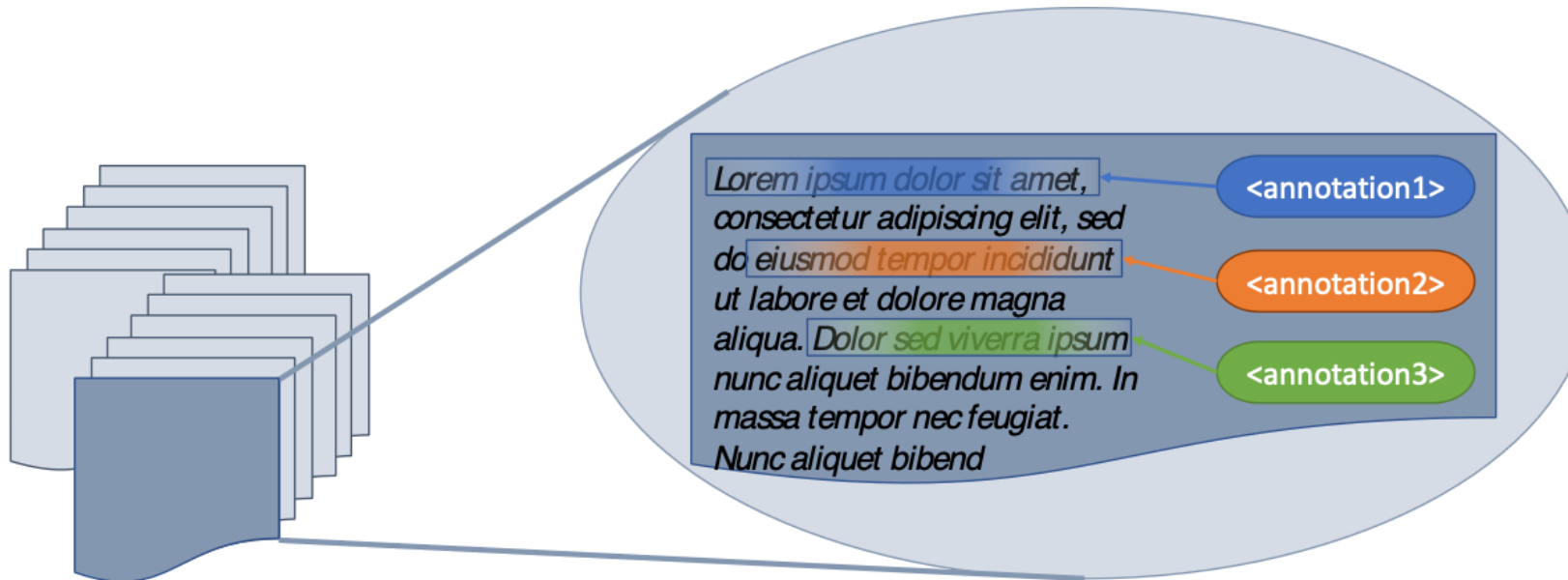UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document

- … are associated with a text span (window) in a document

Different Typs of SCDs:

- Additional definitions

- Links to external sources
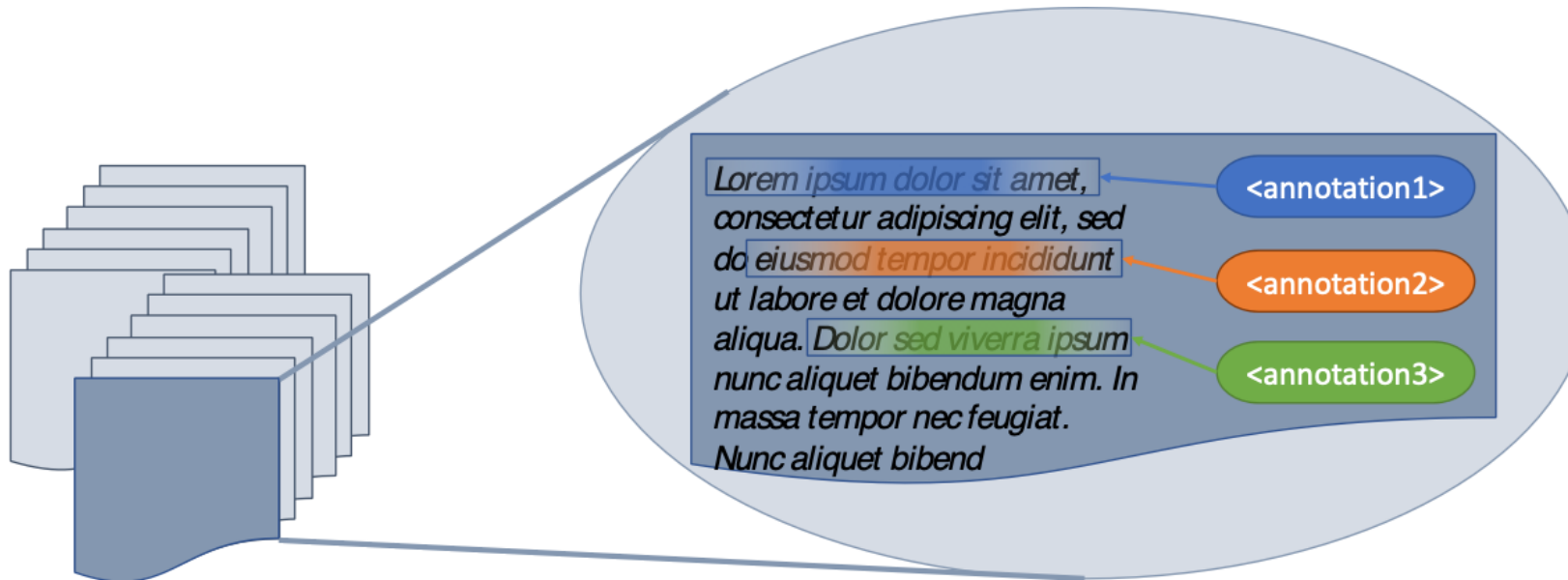
- Relational data to clarify dependencies between entities

*Lorem ipsum dolor sit amet,* — <annotation1>
*consectetur adipiscing elit, sed do eiusmod tempor incididunt* — <annotation2>
*ut labore et dolore magna aliqua. Dolor sed viverra ipsum* — <annotation3>
*nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend*

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document

- … are associated with a text span (window) in a document

- … can take any form

Different Typs of SCDs:

- Additional definitions

- Links to external sources
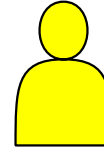
- Relational data to clarify dependencies between entities

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

`<annotation1>`

`<annotation2>`

`<annotation3>`

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document

- ... span

The emergence of mutations in the spread of the coronavirus (SARS-CoV-2) is a natural process. Mutations can develop during the process of copying the genetic make-up, when a cell splits. The more copies of a cell are made, the higher the likelihood that mutations will come into existence. In the case of the coronavirus, a much more infectious type (B.1.1.7) has been recorded in the United Kingdom (UK) since September 2020, and this virus variant is increasingly spreading in Germany. Scientists worry that this mutation is not only significantly more infectious, but can also cause more severe cases of illness. Most recently, B.1.617.2 (Delta) has been found in almost 100 percent of the new confirmed coronavirus (COVID-19) cases in Germany, which illustrates the rapid spread of this mutation.

| | |
|---|---|
| Virus mutation is a normal process | **interim name 2019-nCoV** |

| | |
|---|---|
| B.1.1.7 is a mutation of the coronavirus | **Date of first identified case of SARS-CoV-2: December 2019** |

| | |
|---|---|
| B.1.617.2 is a synonym for Delta mutation of the coronavirus | **Omikron replaces Delta as global dominant corona variant** |

Different Typs of SCDs:

- Additional definitions

- Links to external sources

- Relational data to clarify dependencies between entities

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document
- … are associated with a text span (window) in a document
- … can take any form
- Build SCD-word distribution matrix
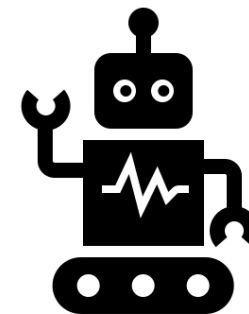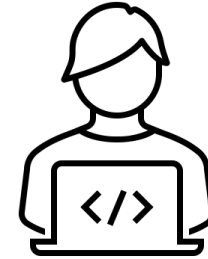- →Use SCD-word distribution for additional tasks

Different Typs of SCDs:

- Additional definitions
- Links to external sources
- Relational data to clarify dependencies between entities

# Subjective Content Descriptions (SCDs)

- … represent additional data for a document

- … are associated with a text span (window) in a document

- … can take any form

- Build SCD-word distribution matrix

- →Use SCD-word distribution for additional tasks

Different Typs of SCDs:

- Additional definitions

- Links to external sources

- Relational data to clarify dependencies between entities

# Lead Scenario: Information Retrieval (IR) Agent

UNIVERSITÄT ZU LÜBECK
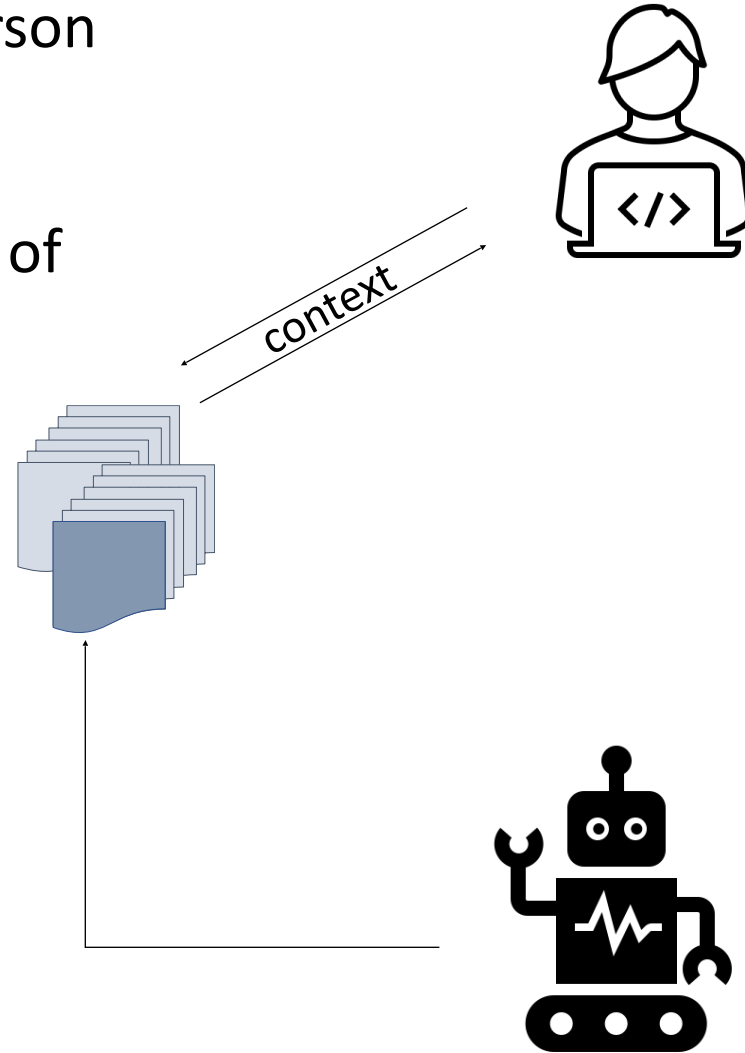INSTITUT FÜR INFORMATIONSSYSTEME

# Lead Scenario: Information Retrieval (IR) Agent

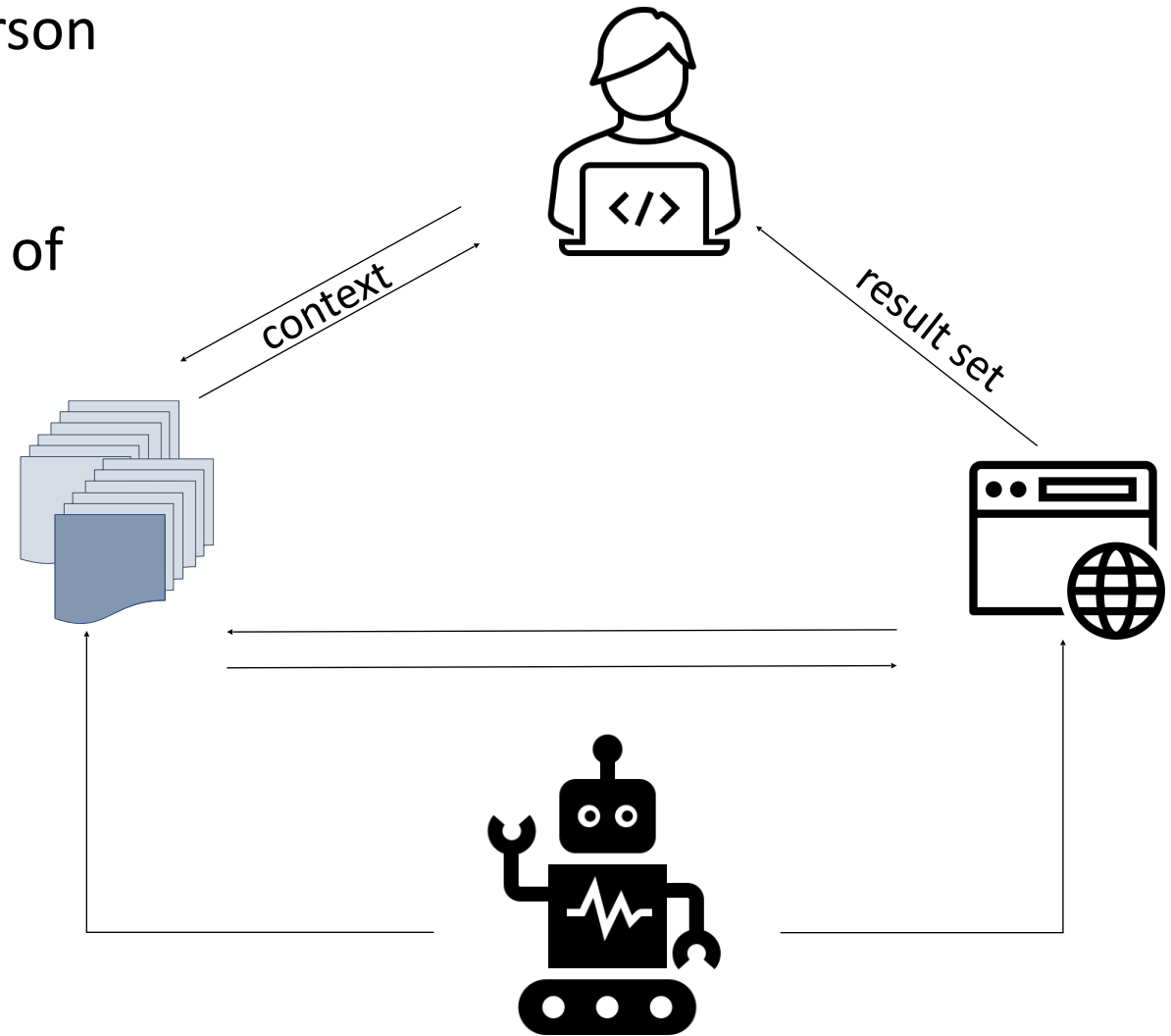- Agent's goal: Meet information need of a person

# Lead Scenario: Information Retrieval (IR) Agent

- Agent's goal: Meet information need of a person

- Agent is working on an IR-corpus that represents a model for the information need of a person
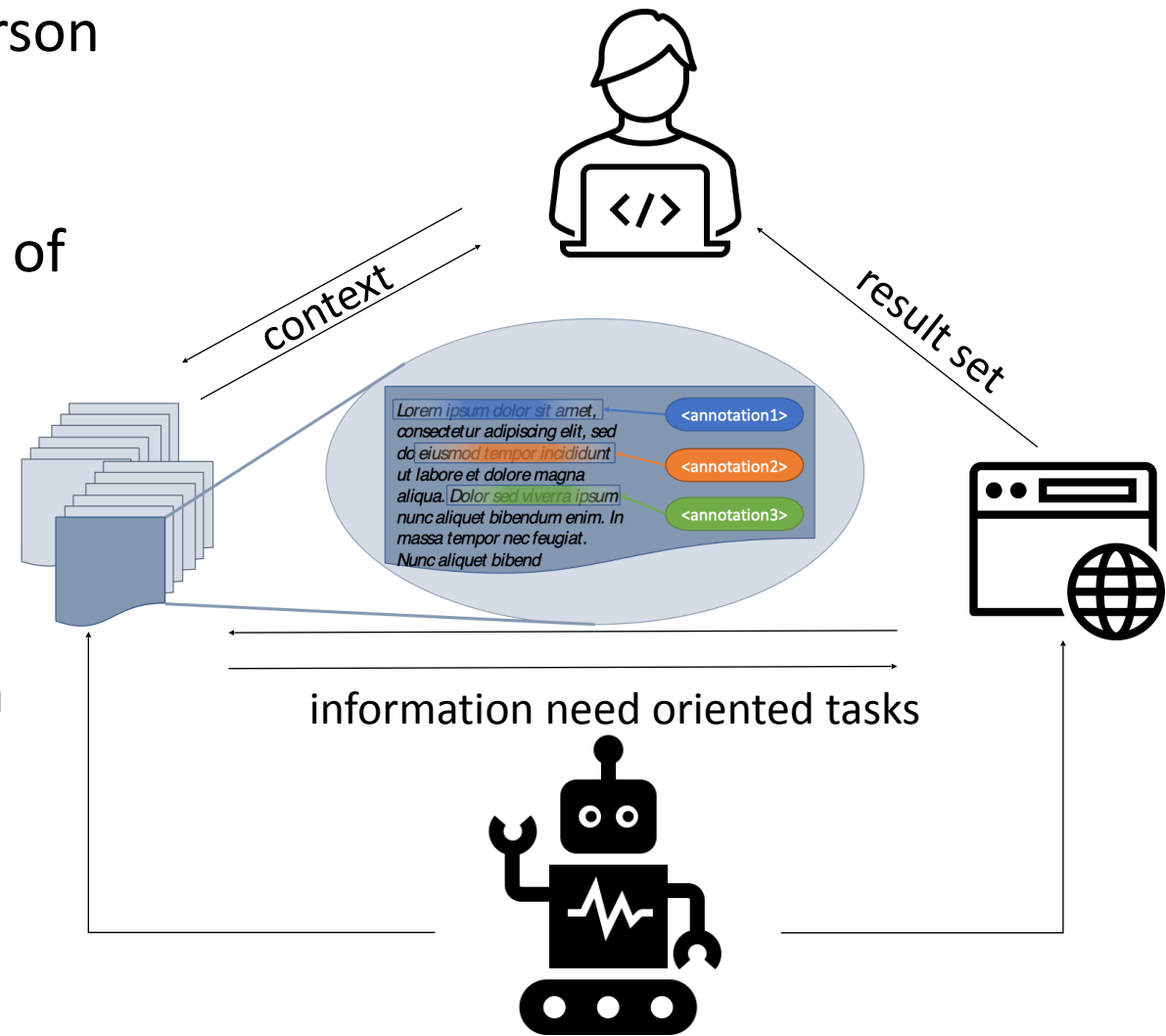
context

# Lead Scenario: Information Retrieval (IR) Agent

- Agent's goal: Meet information need of a person

- Agent is working on an IR-corpus that represents a model for the information need of a person

- Agent optimizes the model to meet the information need

context

result set

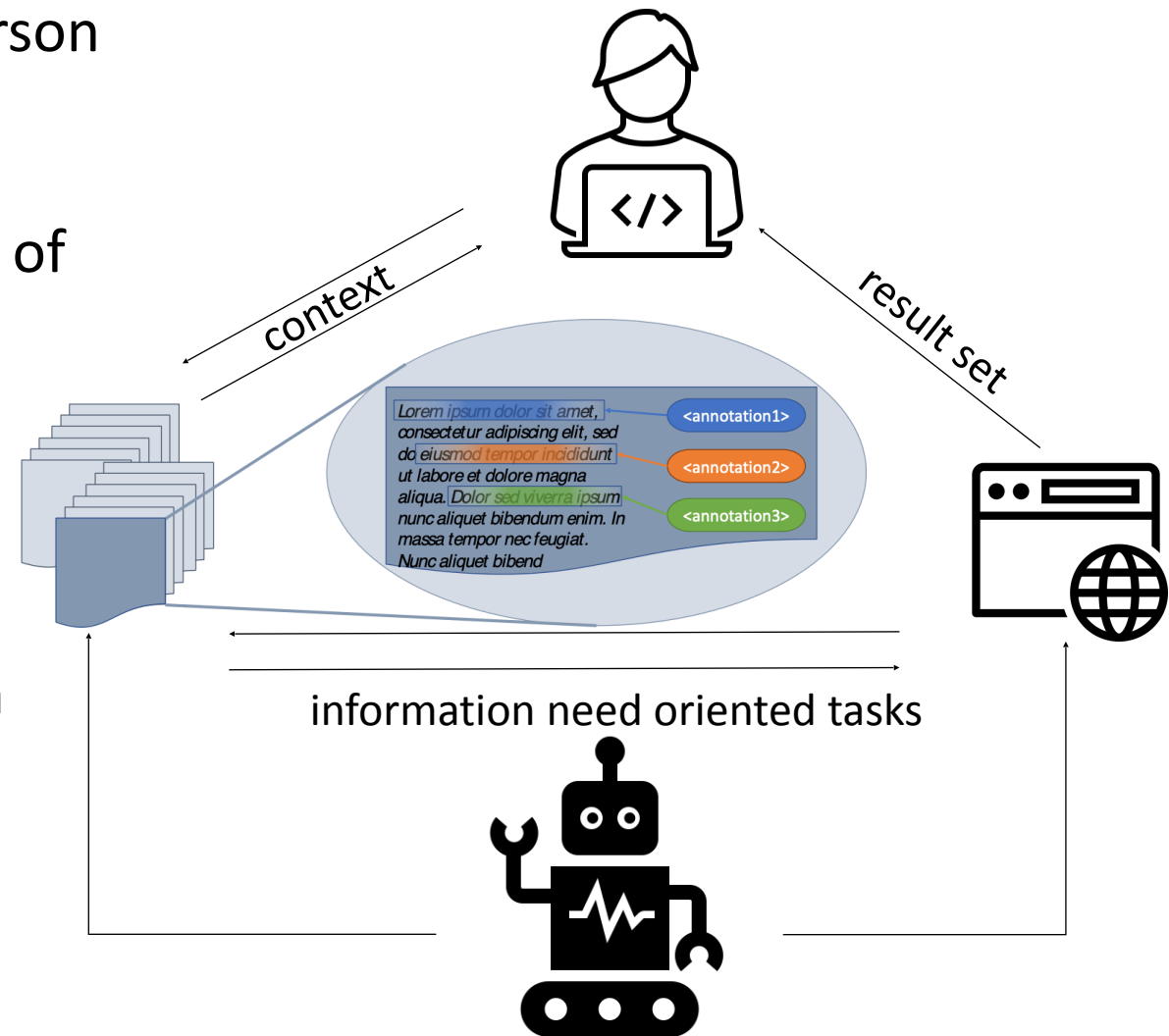# Lead Scenario: Information Retrieval (IR) Agent

- Agent's goal: Meet information need of a person

- Agent is working on an IR-corpus that represents a model for the information need of a person

- Agent optimizes the model to meet the information need

- Documents in the corpus are associated with annotations



context

result set

information need oriented tasks

Lorem *ipsum dolor sit amet*, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

\<annotation1\>
\<annotation2\>
\<annotation3\>

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Lead Scenario: Information Retrieval (IR) Agent

- Agent's goal: Meet information need of a person

- Agent is working on an IR-corpus that represents a model for the information need of a person

- Agent optimizes the model to meet the information need

- Documents in the corpus are associated with annotations

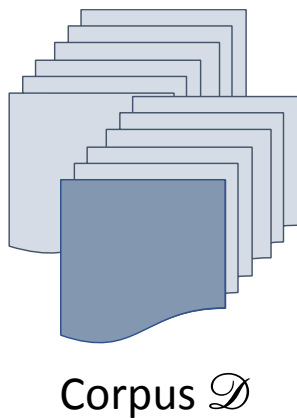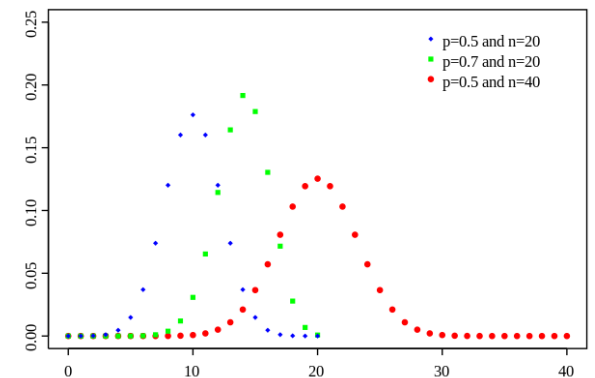→ Subjective Content Descriptions (SCDs)

context

result set

*Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend*

<annotation1>
<annotation2>
<annotation3>

information need oriented tasks

# Foundations: SCD-Word Distribution [2]

- SCD-word distribution results from SCDs associated with *windows* in documents

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment. Proceedings of AI 2019: Advances in Artificial Intelligence, 2019

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME
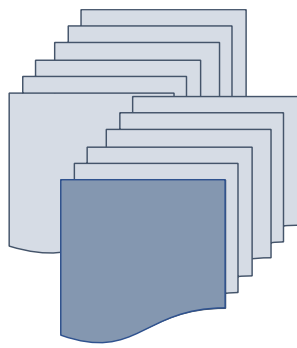
# Foundations: SCD-Word Distribution [2]

- SCD-word distribution results from SCDs associated with *windows* in documents

- For each SCD estimate relative weighted frequency of words

- Binomial distribution to represent weights



Corpus $\mathscr{D}$

$$\delta(\mathscr{D}) = \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} \begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{bmatrix}$$

$w_1 \quad w_2 \quad w_3 \quad \cdots \quad w_n$



p=0.5 and n=20
p=0.7 and n=20
p=0.5 and n=40

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment.  Proceedings of AI 2019: Advances in Artificial Intelligence, 2019

# Foundations: SCD-Word Distribution[2]

- SCD-word distribution results from SCDs associated with *windows* in documents

- For each SCD estimate relative weighted frequency of words

- Binomial distribution to represent weights

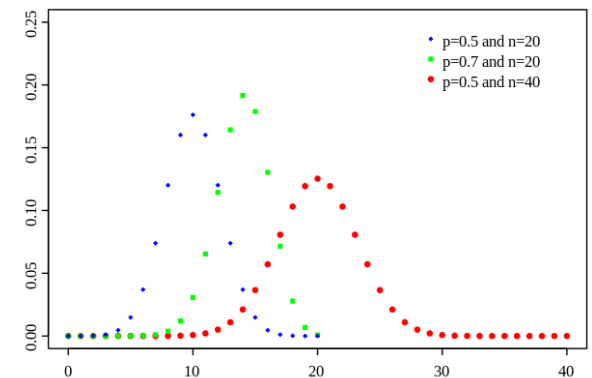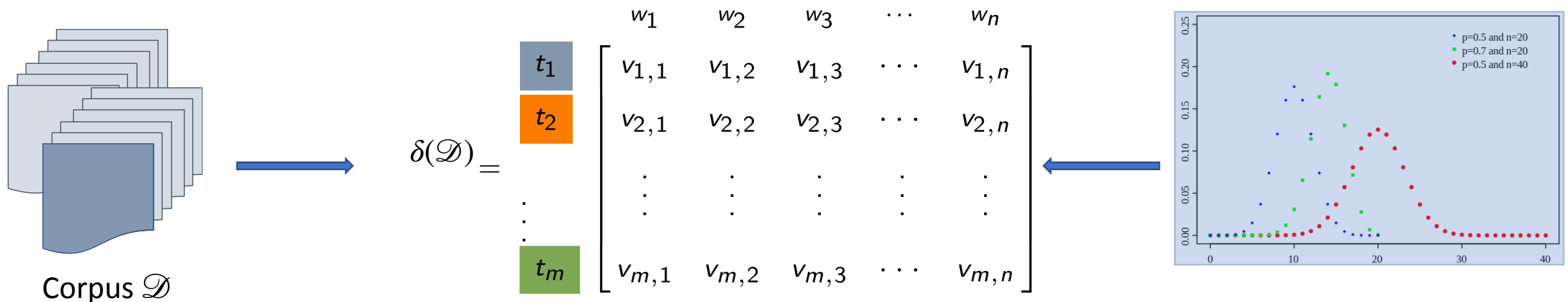**Algorithm 1** Forming SCD-word probability distribution matrix $\delta(\mathcal{D})$

1: **function** BUILDMATRIX(Corpus $\mathcal{D}$)
2:     **Input**: Corpus $\mathcal{D}$
3:     **Output**: SCD-word probability distribution matrix $\delta(\mathcal{D})$
4:     Initialize an $m \times V$ matrix $\delta(\mathcal{D})$ with zeros
5:     **for** each $d \in \mathcal{D}$ **do**
6:         **for** each $t \in T(d)$ **do**
7:             **for** $\rho$ of $t$ **do**
8:                 **for** each $w \in win_{d,\rho}$ **do**
9:                     $\delta(\mathcal{D})[t][w] \mathrel{+}= I(w, win_{d,\rho})$
10:     Normalize $\delta(\mathcal{D})[t]$
11:     **return** $\delta(\mathcal{D})$



Corpus $\mathcal{D}$

$$\delta(\mathscr{D}) = \begin{array}{cc} & \begin{array}{ccccc} w_1 & w_2 & w_3 & \cdots & w_n \end{array} \\ \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} & \left[ \begin{array}{ccccc} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{array} \right] \end{array}$$



p=0.5 and n=20
p=0.7 and n=20
p=0.5 and n=40

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment.  Proceedings of AI 2019: Advances in Artificial Intelligence, 2019

# Foundations: SCD-Word Distribution[2]

- SCD-word distribution results from SCDs associated with *windows* in documents

- For each SCD estimate relative weighted frequency of words

- Binomial distribution to represent weights

**Algorithm 1** Forming SCD-word probability distribution matrix $\delta(\mathcal{D})$

1: **function** BUILDMATRIX(Corpus $\mathcal{D}$)
2:   **Input**: Corpus $\mathcal{D}$
3:   **Output**: SCD-word probability distribution matrix $\delta(\mathcal{D})$
4:   Initialize an $m \times V$ matrix $\delta(\mathcal{D})$ with zeros
5:   **for** each $d \in \mathcal{D}$ **do**
6:     **for** each $t \in T(d)$ **do**
7:       **for** $\rho$ of $t$ **do**
8:         **for** each $w \in win_{d,\rho}$ **do**
9:           $\delta(\mathcal{D})[t][w]$ += $I(w, win_{d,\rho})$
10:   Normalize $\delta(\mathcal{D})[t]$
11:   **return** $\delta(\mathcal{D})$



Corpus $\mathcal{D}$

$$\delta(\mathcal{D}) = \begin{array}{c c} & \begin{array}{ccccc} w_1 & w_2 & w_3 & \cdots & w_n \end{array} \\ \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} & \left[ \begin{array}{ccccc} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{array} \right] \end{array}$$

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Most-Probably Suited SCDs (MPSCDs) [2]

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment. Proceedings of AI 2019: Advances in Artificial Intelligence, 2019

# Most-Probably Suited SCDs (MPSCDs)

- <u>Given:</u> Sequence of text windows in a (new) document

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Most-Probably Suited SCDs (MPSCDs)

- <u>Given:</u> Sequence of text windows in a (new) document

- Words in window define word distribution ($t'$)

*Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend*
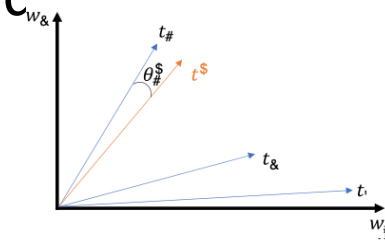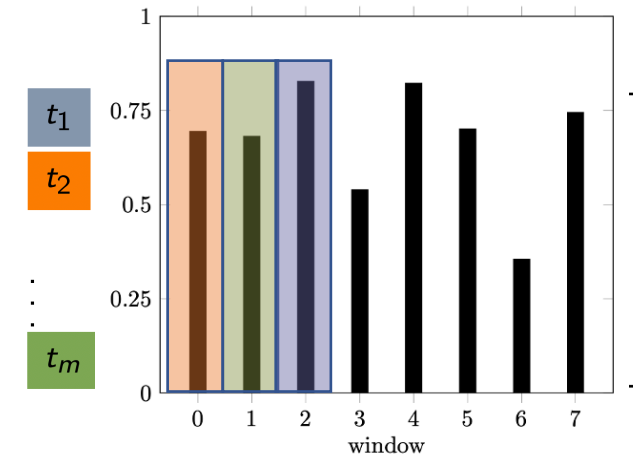
*[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment. Proceedings of AI 2019: Advances in Artificial Intelligence, 2019*

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Most-Probably Suited SCDs (MPSCDs) [2]



- <u>Given:</u> Sequence of text windows in a (new) document

- Words in window define word distribution ($t'$)

- <u>Goal:</u> Estimate MPSCDs for windows in document

# Most-Probably Suited SCDs (MPSCDs) [2]

- <u>Given:</u> Sequence of text windows in a (new) document

- Words in window define word distribution ($t'$)

- <u>Goal:</u> Estimate MPSCDs for windows in document

- Each vector $v_i$ defined in SCD-word matrix $\delta(D)$ defines an angle with window word vector $t'$

$$
\begin{array}{c}
\begin{array}{ccccc} w_1 & w_2 & w_3 & \cdots & w_n \end{array} \\
\begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array}
\begin{bmatrix}
v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\
v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n}
\end{bmatrix}
\end{array}
$$

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment.  Proceedings of AI 2019: Advances in Artificial Intelligence, 2019

# Most-Probably Suited SCDs (MPSCDs) [2]



- <u>Given</u>: Sequence of text windows in a (new) document

- Words in window define word distribution ($t'$)

- <u>Goal</u>: Estimate MPSCDs for windows in document

- Each vector $v_i$ defined in SCD-word matrix $\delta(D)$ defines an angle with window word vector $t'$

- Define function $MPSCD(M, t')$

  - Provides SCD $t_i$ that is associated with SCD-word matrix vector $v_i$ with smallest angle to window word vector $t'$

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment. Proceedings of AI 2019: Advances in Artificial Intelligence, 2019

# Most-Probably Suited SCDs (MPSCDs) [2]

- <u>Given:</u> Sequence of text windows in a (new) document

- Words in window define word distribution ($t'$)

- <u>Goal:</u> Estimate MPSCDs for windows in document

- Each vector $v_i$ defined in SCD-word matrix $\delta(D)$ defines an angle with window word vector $t'$

- Define function $MPSCD(M, t')$
  - Provides SCD $t_i$ that is associated with SCD-word matrix vector $v_i$ with smallest angle to window word vector $t'$

- MPSCD with similarity measure is applied to each window of a document

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment. Proceedings of AI 2019: Advances in Artificial Intelligence, 2019

# Context-specific Corpus Enrichment

- <u>Goal:</u> Add new documents to IR corpus with an initial set of SCDs already associated with documents in the corpus

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[3] Felix Kuhr, Tanya Braun, Ralf Möller: Augmenting and Automating Corpus Enrichment. Proceedings of the 14th IEEE International Conference on Semantic Computing (ICSC-20), 2020

# Context-specific Corpus Enrichment [3]

- <u>Goal:</u> Add new documents to IR corpus with an initial set of SCDs already associated with documents in the corpus

- Different Categories: sim, unrel, rev, ext



similar document



unrelated document
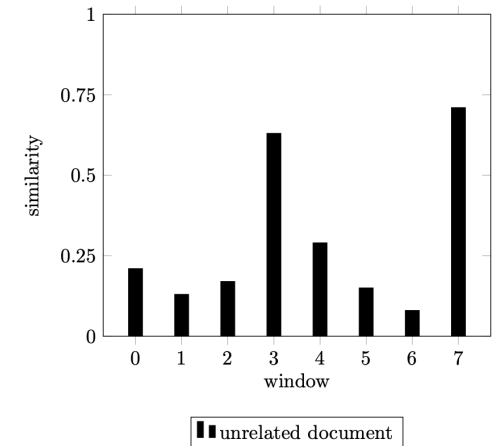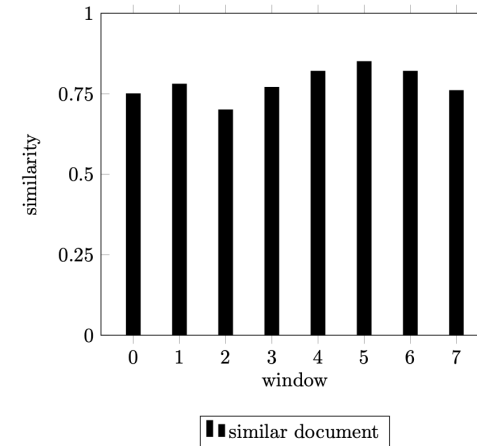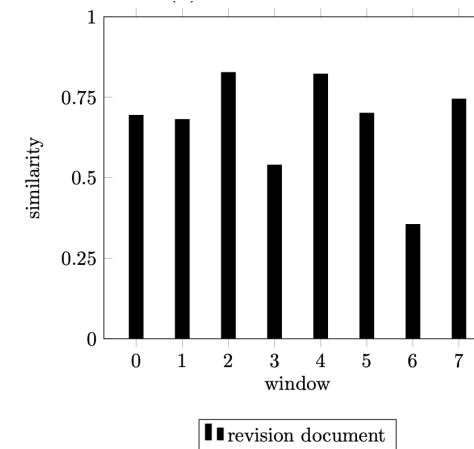


revision document



extending document

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[3] Felix Kuhr, Tanya Braun, Ralf Möller: Augmenting and Automating Corpus Enrichment. Proceedings of the 14th IEEE International Conference on Semantic Computing (ICSC-20), 2020
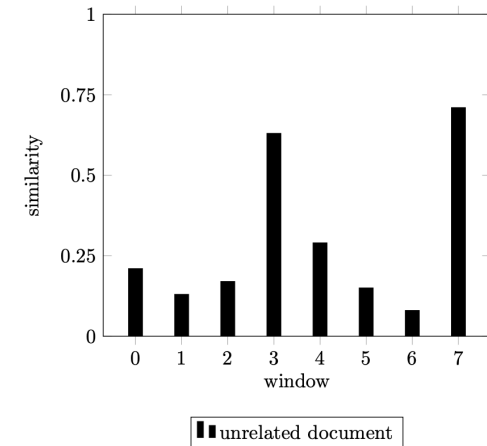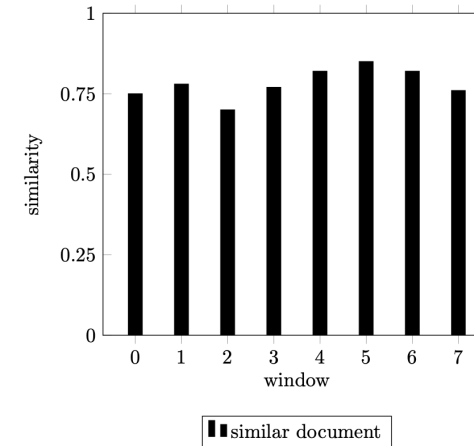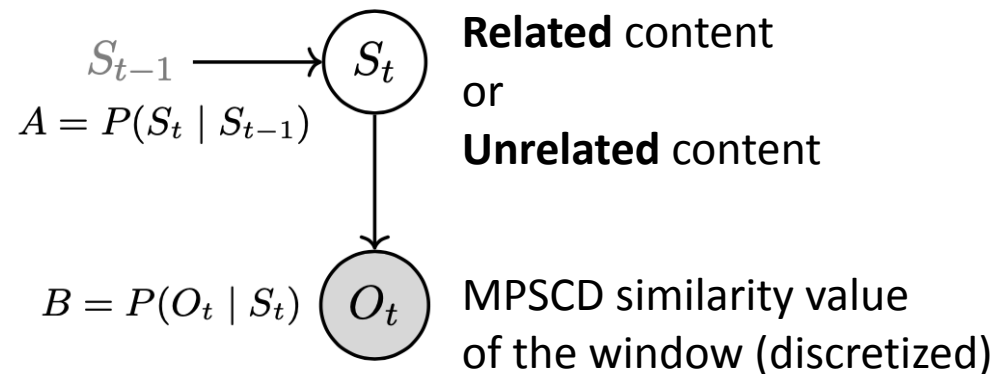
# Context-specific Corpus Enrichment [3]

- <u>Goal:</u> Add new documents to IR corpus with an initial set of SCDs already associated with documents in the corpus

- Different Categories: sim, unrel, rev, ext

- Given: 4 category HMMs, each associated with a category label

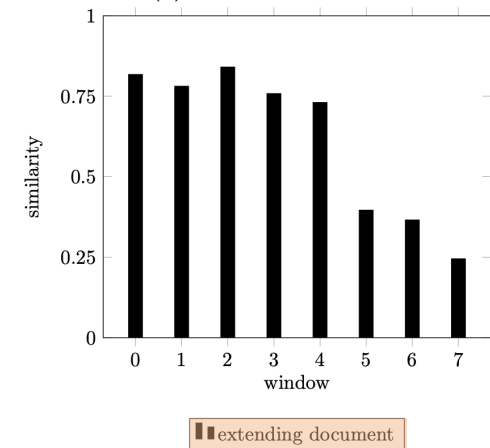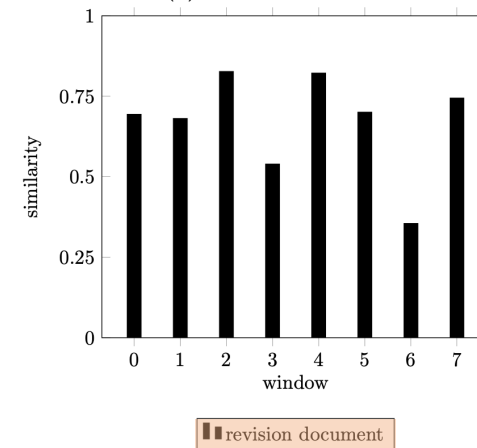[3] Felix Kuhr, Tanya Braun, Ralf Möller: Augmenting and Automating Corpus Enrichment. Proceedings of the 14th IEEE International Conference on Semantic Computing (ICSC-20), 2020

UNIVERSITÄT ZU LÜBECK
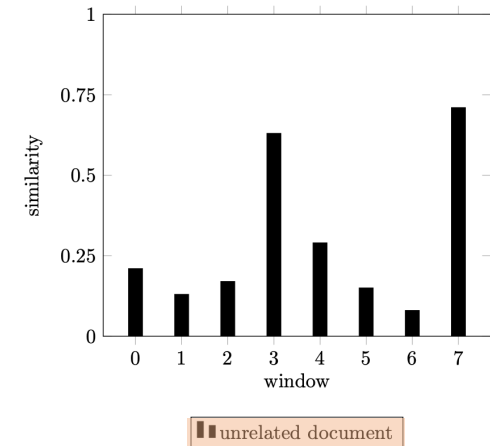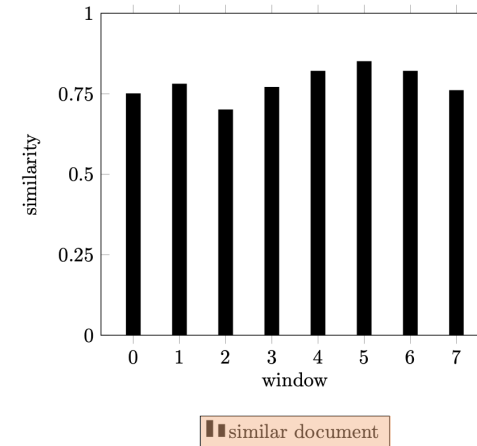INSTITUT FÜR INFORMATIONSSYSTEME

# Context-specific Corpus Enrichment [3]

- <u>Goal:</u> Add new documents to IR corpus with an initial set of SCDs already associated with documents in the corpus

- Different Categories: <span style="color:green">sim</span>, <span style="color:red">unrel</span>, <span style="color:green">rev</span>, <span style="color:green">ext</span>

- Given: 4 category HMMs, each associated with a category label

- HMM Learning by using Baum-Welch Algorithm



$S_{t-1} \longrightarrow S_t$ **Related** content
or
$A = P(S_t \mid S_{t-1})$ **Unrelated** content

$B = P(O_t \mid S_t)$  $O_t$  MPSCD similarity value
of the window (discretized)

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[3] Felix Kuhr, Tanya Braun, Ralf Möller: Augmenting and Automating Corpus Enrichment. Proceedings of the 14th IEEE International Conference on Semantic Computing (ICSC-20), 2020
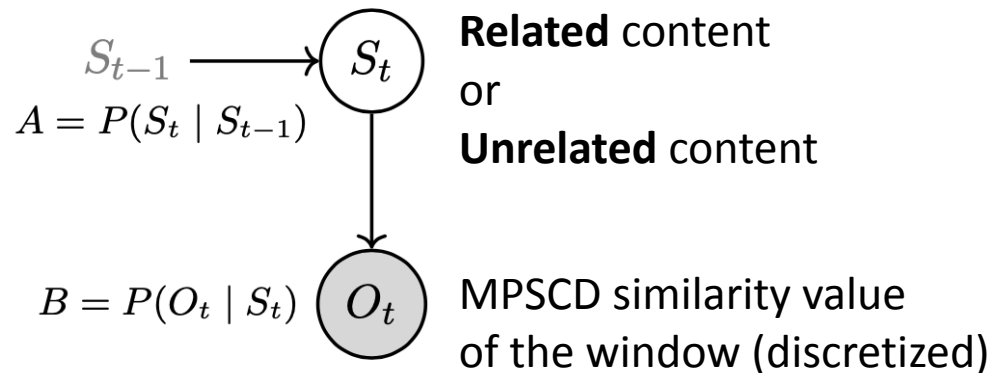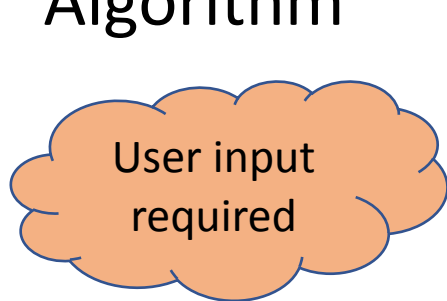
# Context-specific Corpus Enrichment [3]

- <u>Goal:</u> Add new documents to IR corpus with an initial set of SCDs already associated with documents in the corpus

- Different Categories: sim, unrel, rev, ext

- Given: 4 category HMMs, each associated with a category label

- HMM Learning by using Baum-Welch Algorithm

User input required

$S_{t-1} \rightarrow S_t$

$A = P(S_t \mid S_{t-1})$

**Related** content
or
**Unrelated** content

$B = P(O_t \mid S_t)$ $O_t$

MPSCD similarity value
of the window (discretized)



similar document



unrelated document



revision document



extending document

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[3] Felix Kuhr, Tanya Braun, Ralf Möller: Augmenting and Automating Corpus Enrichment. Proceedings of the 14th IEEE International Conference on Semantic Computing (ICSC-20), 2020

# Context-specific Corpus Enrichment - Decision Making Process

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend
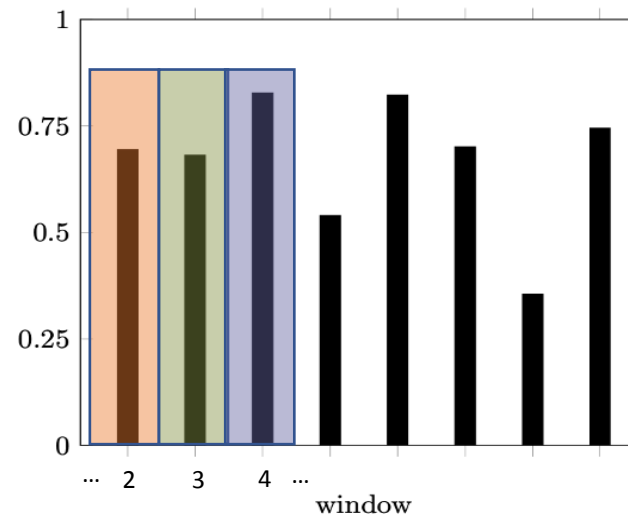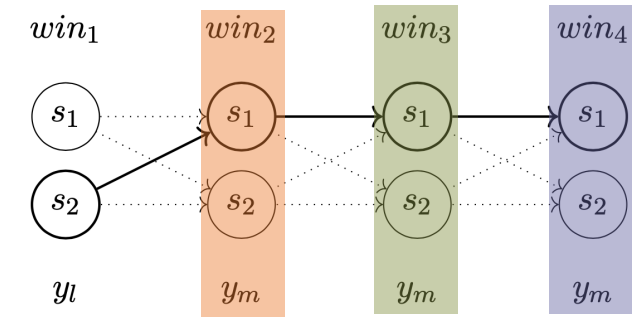
Given: new document

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

<annotation2>

<annotation3>

<annotation1>

Determine the MPSCD sequence for the window sequence of the new document based on available SCD-word distribution

$$
\begin{array}{c}
\phantom{t_1} \\[0.2em]
t_1 \\[0.6em]
t_2 \\[0.6em]
\vdots \\[0.6em]
t_m
\end{array}
\begin{array}{cccccc}
w_1 & w_2 & w_3 & \cdots & w_n \\
\left[\begin{array}{ccccc}
v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\
v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n}
\end{array}\right]
\end{array}
$$

Given: SCD-word distribution of IR corpus

# Context-specific Corpus Enrichment - Decision Making Process



Given: new document

Given: SCD-word distribution of IR corpus

Determine the MPSCD sequence for the window sequence of the new document based on available SCD-word distribution

Discretize similarity values:
$$y_l: 0 - 0.3, \quad y_m: 0.3 - 0.75, \quad y_h: 0.75 - 1$$

# Context-specific Corpus Enrichment - Decision Making Process



Given: new document

$$
\begin{array}{c}
\phantom{t_1}\quad w_1 \quad\ w_2 \quad\ w_3 \quad\ \cdots \quad\ w_n \\
\begin{array}{c}
t_1 \\ t_2 \\ \vdots \\ t_m
\end{array}
\begin{bmatrix}
v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\
v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n}
\end{bmatrix}
\end{array}
$$

Given: SCD-word distribution of IR corpus

Determine the MPSCD sequence for the window sequence of the new document based on available SCD-word distribution



Discretize similarity values:
$y_l : 0 - 0.3, \ y_m : 0.3 - 0.75, \ y_h : 0.75 - 1$

Determine MPE sequence w.r.t. each category HMM on sequence of MPSCD similarity values.

# Context-specific Corpus Enrichment - Decision Making Process



Given: new document

Given: SCD-word distribution of IR corpus

Determine the MPSCD sequence for the window sequence of the new document based on available SCD-word distribution

Discretize similarity values:
$y_l: 0 - 0.3, y_m: 0.3 - 0.75, y_h: 0.75 - 1$

Determine MPE sequence w.r.t. each category HMM on sequence of MPSCD similarity values.

Take category of HMM with most-likely MPE sequence as classification

# Context-specific Corpus Enrichment - Decision Making Process



Given: new document

$$
\begin{array}{c}
\phantom{t_1} \\
\boxed{t_1} \\
\boxed{t_2} \\
\vdots \\
\boxed{t_m}
\end{array}
\begin{array}{ccccc}
w_1 & w_2 & w_3 & \cdots & w_n \\
\begin{bmatrix}
v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\
v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n}
\end{bmatrix}
\end{array}
$$

Given: SCD-word distribution of IR corpus

Determine the MPSCD sequence for the window sequence of the new document based on available SCD-word distribution

Discretize similarity values:
$$y_l : 0 - 0.3, \ y_m : 0.3 - 0.75, \ y_h : 0.75 - 1$$

Determine MPE sequence w.r.t. each category HMM on sequence of MPSCD similarity values.

Take category of HMM with most-likely MPE sequence as classification

Extend corpus based on document category and transfer SCDs above a threshold

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Context-specific Corpus Enrichment - Decision Making Process

Given: new document

Given: SCD-word distribution of IR corpus

$$
\begin{array}{c c c c c c}
 & w_1 & w_2 & w_3 & \cdots & w_n \\
t_1 & v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\
t_2 & v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\
\vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
t_m & v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n}
\end{array}
$$

Determine the MPSCD sequence for the window sequence of the new document based on available SCD-word distribution

<annotation2>

<annotation3>

<annotation1>

$win_1$  $win_2$  $win_3$  $win_4$

$s_1$  $s_2$

$y_l$  $y_m$  $y_m$  $y_m$

Focus on SCD similarity values

Discretize similarity values:
$y_l: 0 - 0.3, \quad y_m: 0.3 - 0.75, \quad y_h: 0.75 - 1$

Determine MPE sequence w.r.t. each category HMM on sequence of MPSCD similarity values.

Take category of HMM with most-likely MPE sequence as classification

Focus on content of SCDs

Extend corpus based on document category and transfer SCDs above a threshold

# Corpus-Driven Document Enrichment using SCDs

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Corpus-Driven Document Enrichment using SCDs

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

# Corpus-Driven Document Enrichment using SCDs

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

# Corpus-Driven Document Enrichment using SCDs

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

# Corpus-Driven Document Enrichment using SCDs

Goal: Enrich a document with _relevant_ SCDs associated with other documents in an IR-corpus.

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Corpus-Driven Document Enrichment using SCDs

**Goal:** Enrich a document with *relevant* SCDs associated with other documents in an IR-corpus.

*Fixed-point iteration procedure:*

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

# Corpus-Driven Document Enrichment using SCDs



Iterative SCD Enrichment of documents

IR-corpus

$related-documents(d_i, \ IR-corpus)$

- Subset of $IR-corpus$
  - $topic-similar$ documents whose
  - SCDs are $SCD-similar$ to $d_i$

Goal: Enrich a document with *relevant* SCDs associated with other documents in an IR-corpus.

*Fixed-point iteration procedure:*

- determine the expected related documents in IR-corpus $D$,

$d_i$ - related documents

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Corpus-Driven Document Enrichment using SCDs



Goal: Enrich a document with *relevant* SCDs associated with other documents in an IR-corpus.

*Fixed-point iteration procedure:*

- determine the expected related documents in IR-corpus $D$,
- determine the set of SCDs $T$ from $D$ that are newly added to $d$, then

# Corpus-Driven Document Enrichment using SCDs



Topic similarity
SCD similarity
Frequency

$d_i$ - related
documents

**Goal:** Enrich a document with _relevant_ SCDs associated with other documents in an IR-corpus.

_Fixed-point iteration procedure:_

- determine the expected related documents in IR-corpus $D$,
- determine the set of SCDs $T$ from $D$ that are newly added to $d$, then

_expected−relevance$(t, d_i)$_

- estimates relevance of $t$ w.r.t. d by document $d_i$:

  Mean topic similarity of related documents containing SCD $t$

  Mean SCD similarity to related documents containing SCD $t$

  Number of related documents in which SCD $t$ occurs

_mean−expected−relevance$(d_i)$_
average expected relevance value of SCDs in $d_i$ -related documents
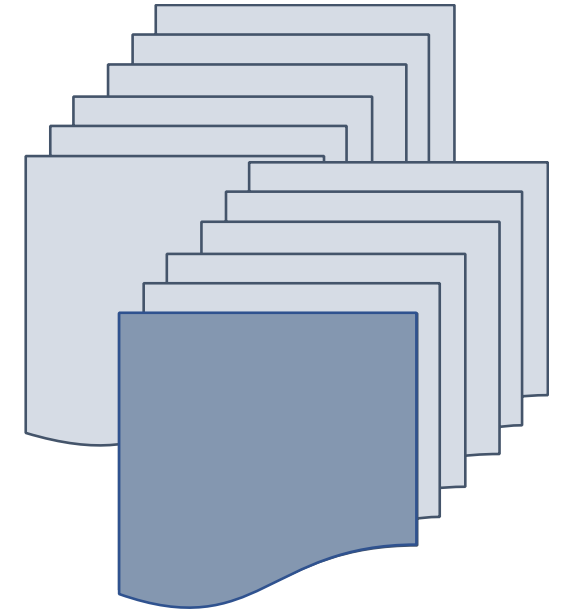
UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Corpus-Driven Document Enrichment using SCDs



**Goal:** Enrich a document with *relevant* SCDs associated with other documents in an IR-corpus.

*Fixed-point iteration procedure:*

- determine the expected related documents in IR-corpus $D$,
- determine the set of SCDs $T$ from $D$ that are newly added to $d$, then

# Corpus-Driven Document Enrichment using SCDs

**Goal:** Enrich a document with *relevant* SCDs associated with other documents in an IR-corpus.

*Fixed-point iteration procedure:*

- determine the expected related documents in IR-corpus $D$,
- determine the set of SCDs $T$ from $D$ that are newly added to $d$, then
- determine the expected related documents $D$ again, and so one

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Corpus-Driven Document Enrichment using SCDs



$d_i$ - related documents

**Goal:** Enrich a document with *relevant* SCDs associated with other documents in an IR-corpus.

*Fixed-point iteration procedure:*

- determine the expected related documents in IR-corpus $D$,
- determine the set of SCDs $T$ from $D$ that are newly added to $d$, then
- determine the expected related documents $D$ again, and so one
- until no more SCDs are assigned to document $d$.

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Corpus-Driven Document Enrichment using SCDs

$d_i$ - related documents

enrich($d_i$, $IR-corpus$)

- Add SCD $t$ to $d_i$ if

$$expected-relevance(t, d_i) > mean-expected-relevance(d_i)$$

- Iterative enrichment process
  - Related documents changes with enriched SCDs
- Terminating enrichment process
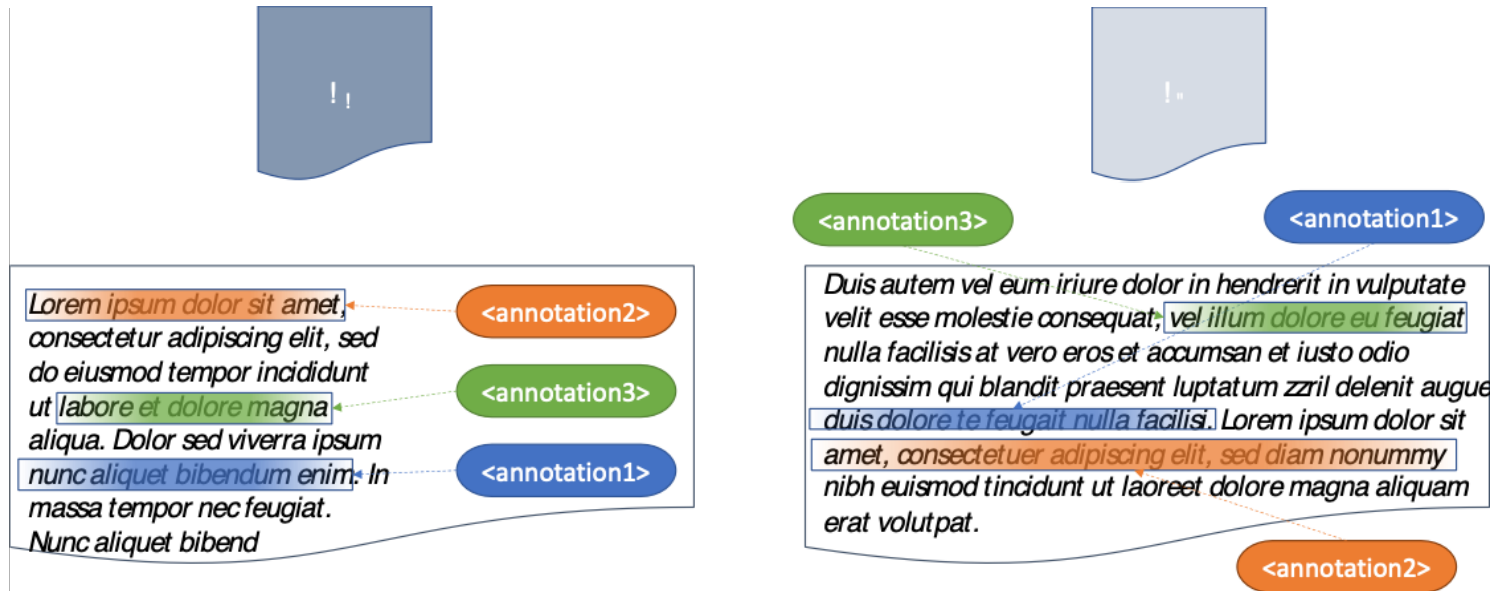  - Value of SCD similarity of $d_i$ to related documents increases in a negligible way
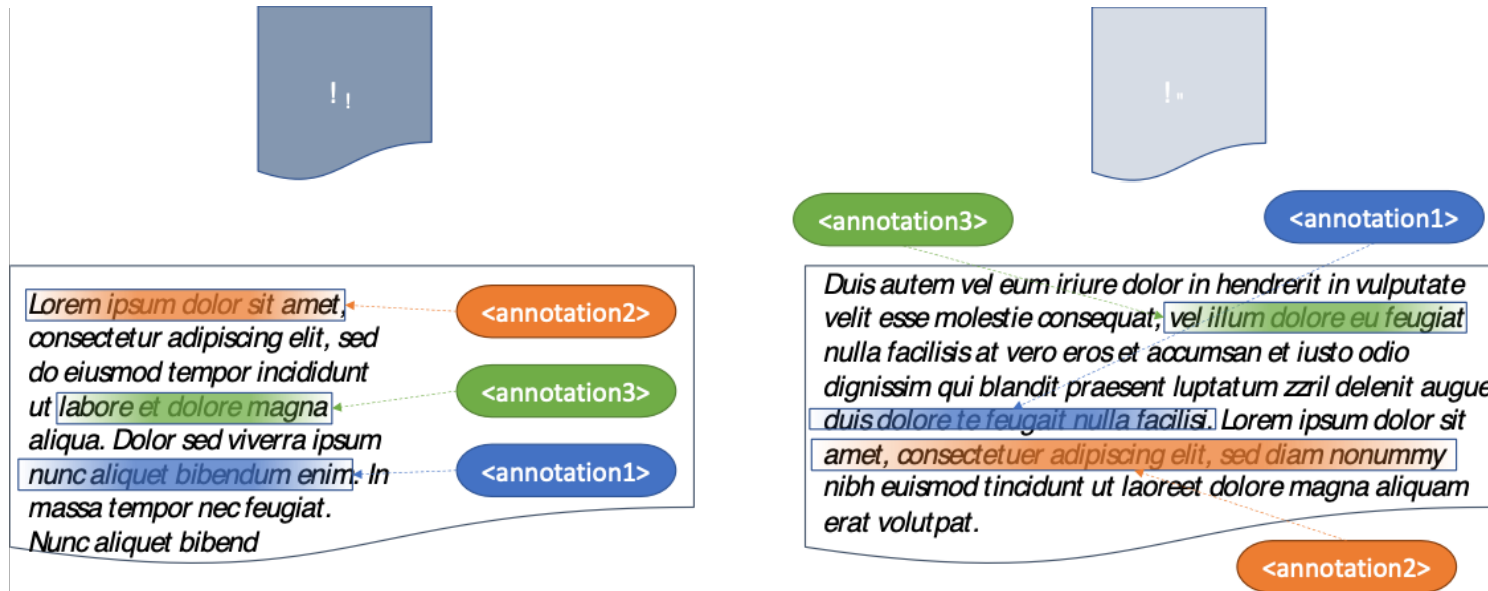
Goal: Enrich a document with _relevant_ SCDs associated with other documents in an IR-corpus.

_Fixed-point iteration procedure:_

- determine the expected related documents in IR-corpus $D$,
- determine the set of SCDs $T$ from $D$ that are newly added to $d$, then
- determine the expected related documents $D$ again, and so one
- until no more SCDs are assigned to document $d$.

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019

# Expected Relevance Value

# Expected Relevance Value

- <u>Given:</u> document $d_j$ from IR-corpus $D$

- <u>Question:</u> What is the expected relevance value of an SCD associated to a related document?

# Expected Relevance Value

- Given: document $d_j$ from IR-corpus $D$

- Question: What is the expected relevance value of an SCD associated to a related document?

Relevance of an SCD depends on the information need of a user



*Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend*

<annotation1>
<annotation2>
<annotation3>

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Expected Relevance Value

- <u>Given:</u> document $d_j$ from IR-corpus $D$

- <u>Question:</u> What is the expected relevance value of an SCD associated to a related document?

- Some ways to adjust performance:
  - Similarity between documents
  - Similarity between SCDs
  - Frequency of SCDs

Relevance of an SCD depends on the information need of a user

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

&lt;annotation1&gt;

&lt;annotation2&gt;

&lt;annotation3&gt;

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Bootstrap Approaches for SCDs

**Inline SCDs**[5]

[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Bootstrap Approaches for SCDs

**Inline SCDs** [5]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective
Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)

# Bootstrap Approaches for SCDs

**Inline SCDs** [5]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values
    - Small similarity values → different content → new inline-SCDs
    - Inline-SCD = Content of window
    - Inline-SCD represent new row in SCD word matrix
    - HMMs given as user input

$w_1^d \ w_2^d \ w_3^d \ w_4^d \ w_5^d \ w_6^d \ w_7^d \ w_8^d \ {\color{red}w_9^d \ w_{10}^d w_{11}^d w_{12}^d w_{13}^d} w_{14}^d w_{15}^d w_{16}^d w_{17}^d w_{18}^d w_{19}^d w_{20}^d$



$wind_{d,t,\rho_i}$

[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Bootstrap Approaches for SCDs

## Inline SCDs [5]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values
  - Small similarity values → different content → new inline-SCDs
  - Inline-SCD = Content of window
  - Inline-SCD represent new row in SCD word matrix
  - HMMs given as user input

[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Bootstrap Approaches for SCDs

## Inline SCDs [5]

## Adapting SCD word distribution from another IR-corpus [6]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values

    - Small similarity values → different content → new inline-SCDs

    - Inline-SCD = Content of window

    - Inline-SCD represent new row in SCD word matrix

    - HMMs given as user input

*[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*
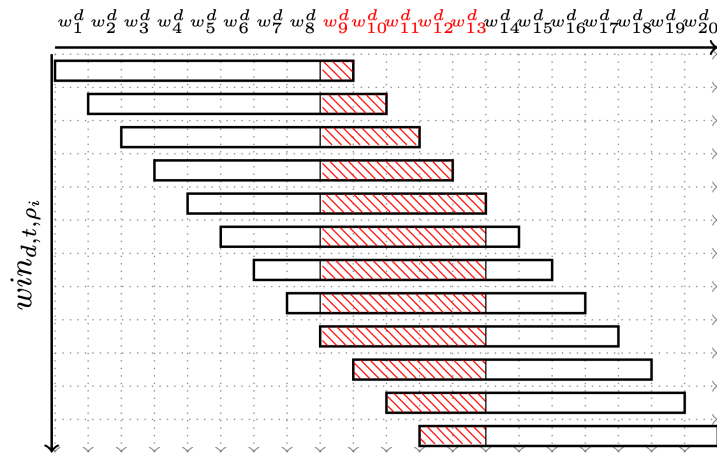
*[6] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Bootstrap Approaches for SCDs

## Inline SCDs [5]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values
  - Small similarity values → different content → new inline-SCDs
  - Inline-SCD = Content of window
  - Inline-SCD represent new row in SCD word matrix
  - HMMs given as user input

$w_1^d\ w_2^d\ w_3^d\ w_4^d\ w_5^d\ w_6^d\ w_7^d\ w_8^d\ \textcolor{red}{w_9^d\ w_{10}^d w_{11}^d w_{12}^d w_{13}^d} w_{14}^d w_{15}^d w_{16}^d w_{17}^d w_{18}^d w_{19}^d w_{20}^d$

Similarity value    0.9
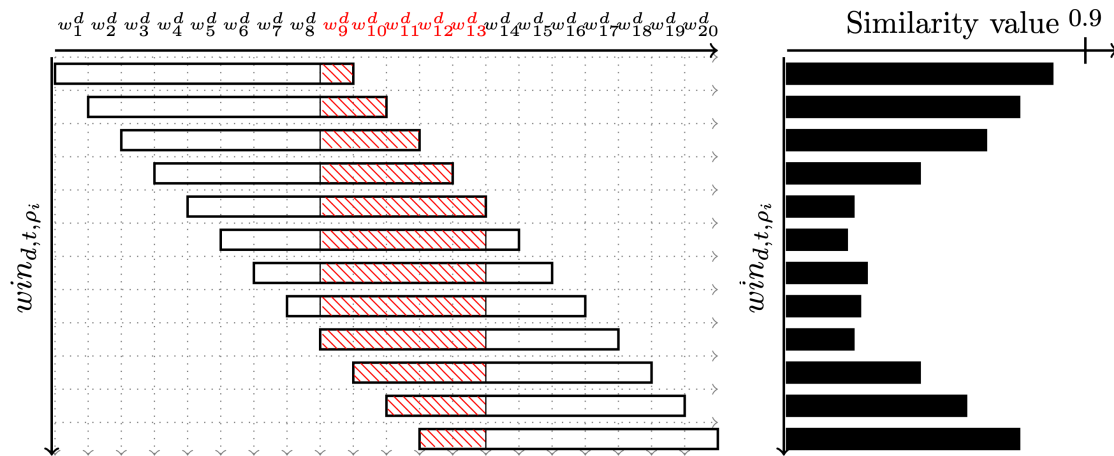
$win_{d,t,\rho_i}$

$win_{d,t,\rho_i}$

*[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*

## Adapting SCD word distribution from another IR-corpus [6]

$D_i$

$$\delta(D_i) = \begin{array}{c} \\ t_1 \\ t_2 \\ \vdots \\ \vdots \\ t_m \end{array} \begin{array}{ccccc} w_1 & w_2 & w_3 & \cdots & w_n \\ \left[\begin{array}{ccccc} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{array}\right] \end{array}$$
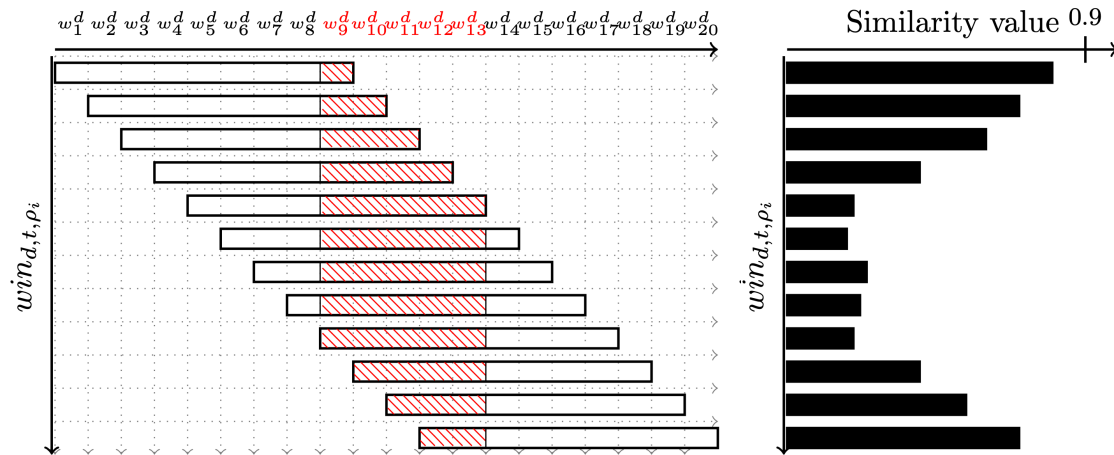
*[6] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

Context-aware Corpus Annotation Using Subjective Content Descriptions

# Bootstrap Approaches for SCDs

## Inline SCDs [5]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values

  - Small similarity values → different content → new inline-SCDs
  - Inline-SCD = Content of window
  - Inline-SCD represent new row in SCD word matrix
  - HMMs given as user input



## Adapting SCD word distribution from another IR-corpus [6]



$$\delta(D_i) = \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} \begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{bmatrix}$$

$$\delta(D_j) = \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} \begin{bmatrix} v'_{1,1} & v'_{1,2} & v'_{1,3} & \cdots & v'_{1,n} \\ v'_{2,1} & v'_{2,2} & v'_{2,3} & \cdots & v'_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v'_{m,1} & v'_{m,2} & v'_{m,3} & \cdots & v'_{m,n} \end{bmatrix}$$

*[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*

*[6] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*

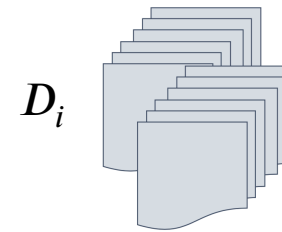UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Bootstrap Approaches for SCDs

## Inline SCDs [5]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values
  - Small similarity values → different content → new inline-SCDs
  - Inline-SCD = Content of window
  - Inline-SCD represent new row in SCD word matrix
  - HMMs given as user input



## Adapting SCD word distribution from another IR-corpus [6]



adapted version

$$\delta(D_i) = \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} \begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{bmatrix}$$

$$\delta(D_j) = \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} \begin{bmatrix} v'_{1,1} & v'_{1,2} & v'_{1,3} & \cdots & v'_{1,n} \\ v'_{2,1} & v'_{2,2} & v'_{2,3} & \cdots & v'_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v'_{m,1} & v'_{m,2} & v'_{m,3} & \cdots & v'_{m,n} \end{bmatrix}$$

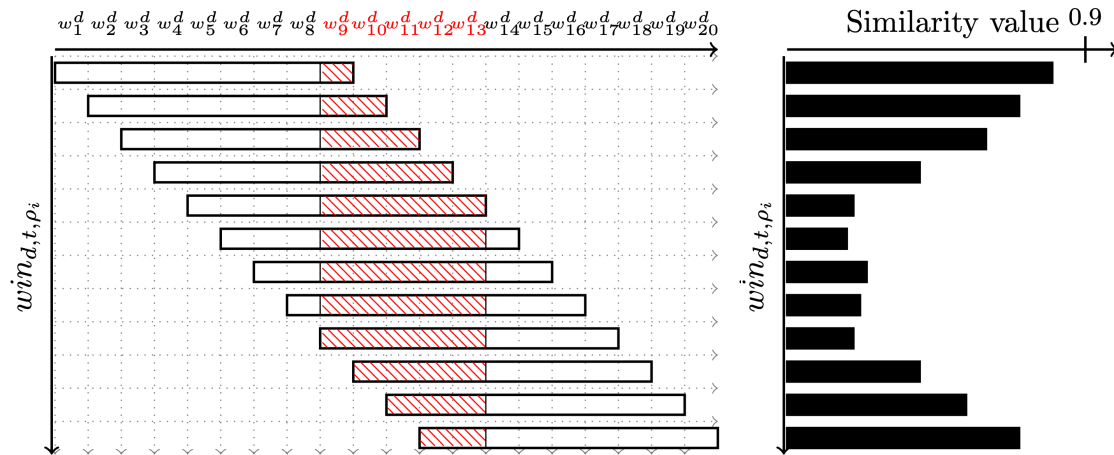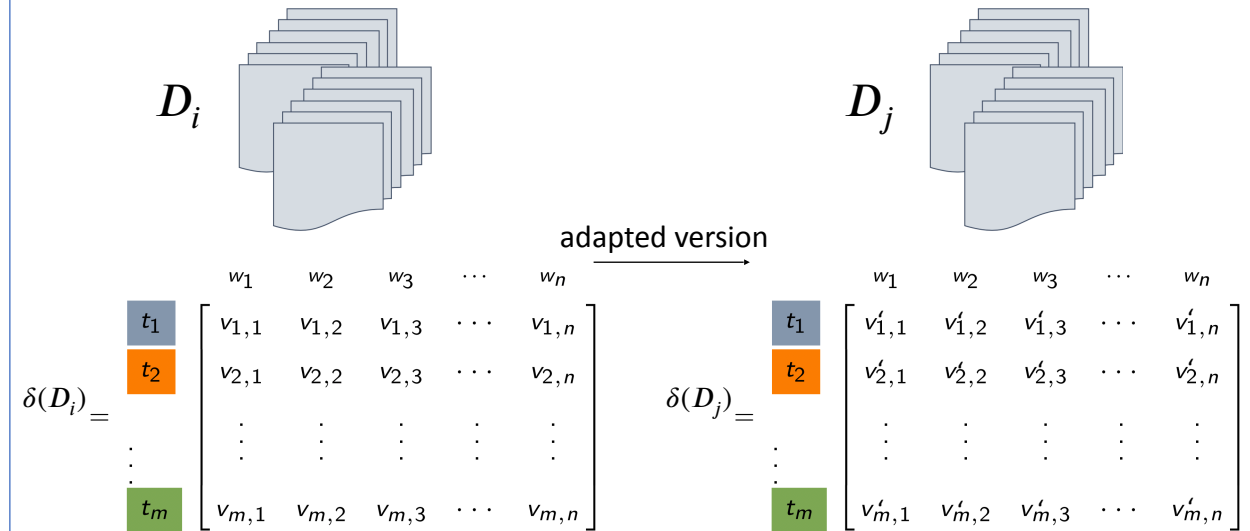- Adapt SCD word distribution from IR-corpus $D_i$ to documents in $D_j$

[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)
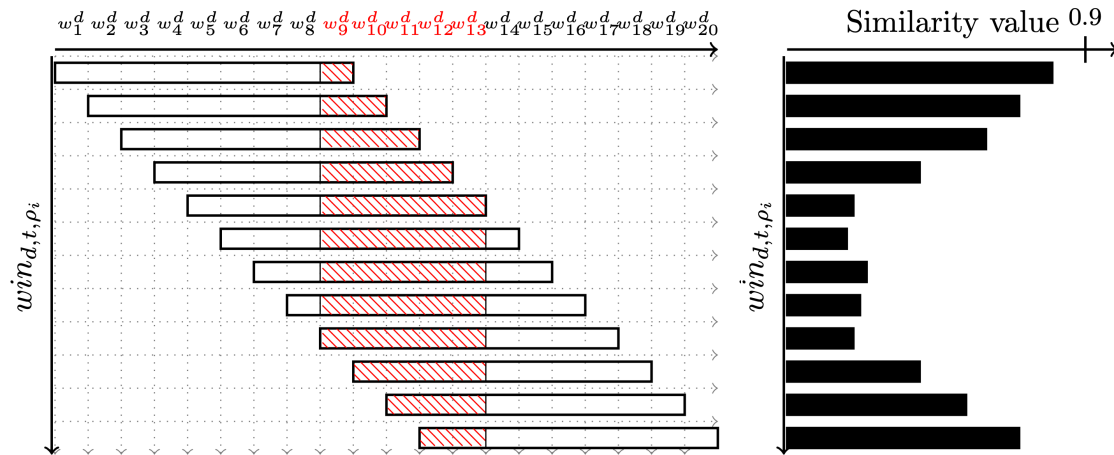
[6] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

Context-aware Corpus Annotation Using Subjective Content Descriptions

# Bootstrap Approaches for SCDs

## Inline SCDs [5]

## Adapting SCD word distribution from another IR-corpus [6]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values
  - Small similarity values → different content → new inline-SCDs
  - Inline-SCD = Content of window
  - Inline-SCD represent new row in SCD word matrix
  - HMMs given as user input



$D_i$    $D_j$

adapted version

$$\delta(D_i) = \begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{bmatrix}$$

$$\delta(D_j) = \begin{bmatrix} v'_{1,1} & v'_{1,2} & v'_{1,3} & \cdots & v'_{1,n} \\ v'_{2,1} & v'_{2,2} & v'_{2,3} & \cdots & v'_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v'_{m,1} & v'_{m,2} & v'_{m,3} & \cdots & v'_{m,n} \end{bmatrix}$$

- Adapt SCD word distribution from IR-corpus $D_i$ to documents in $D_j$
  - Analyze difference in word distributions of documents in corpus $D_i$ and $D_j$

*[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*
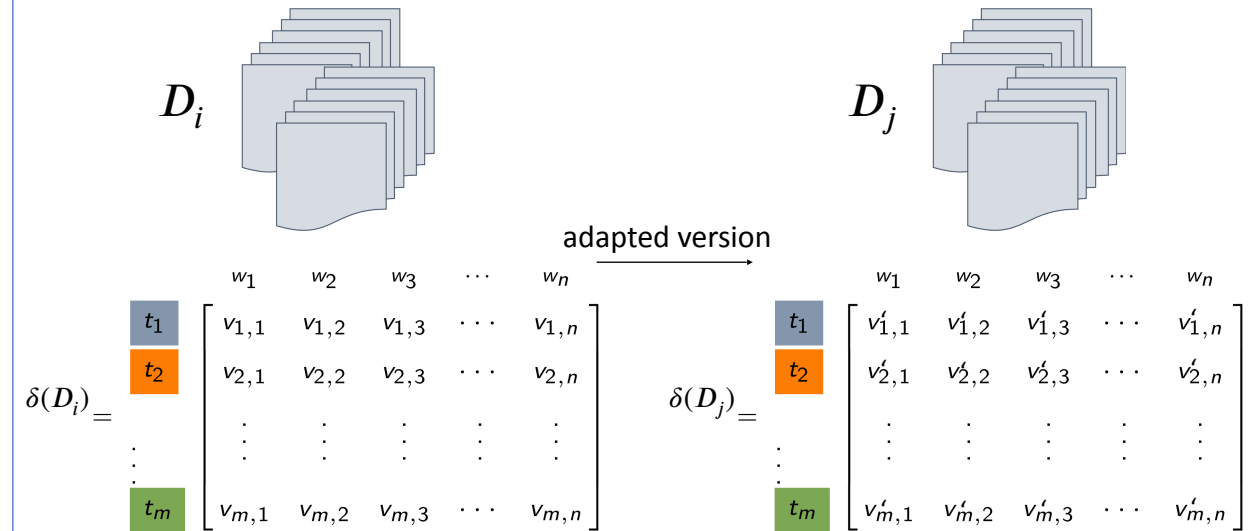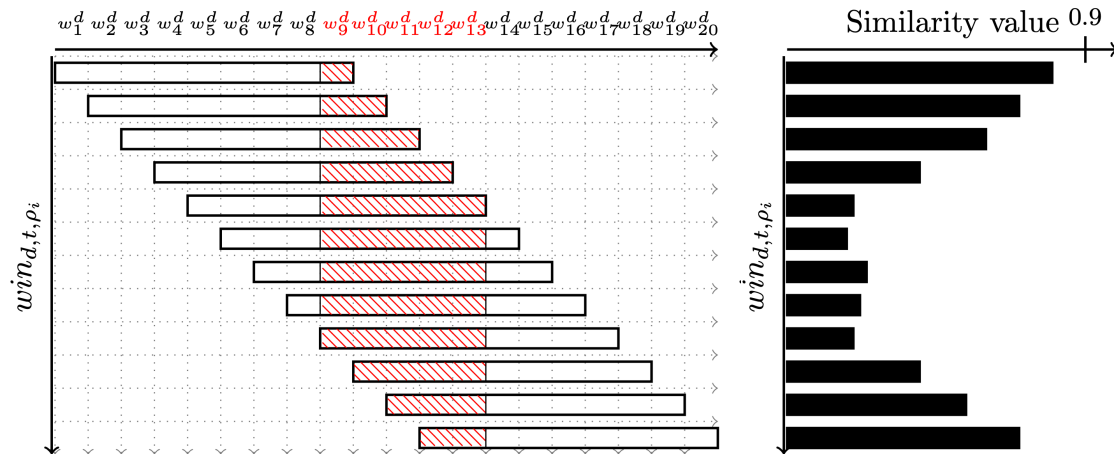
*[6] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)*
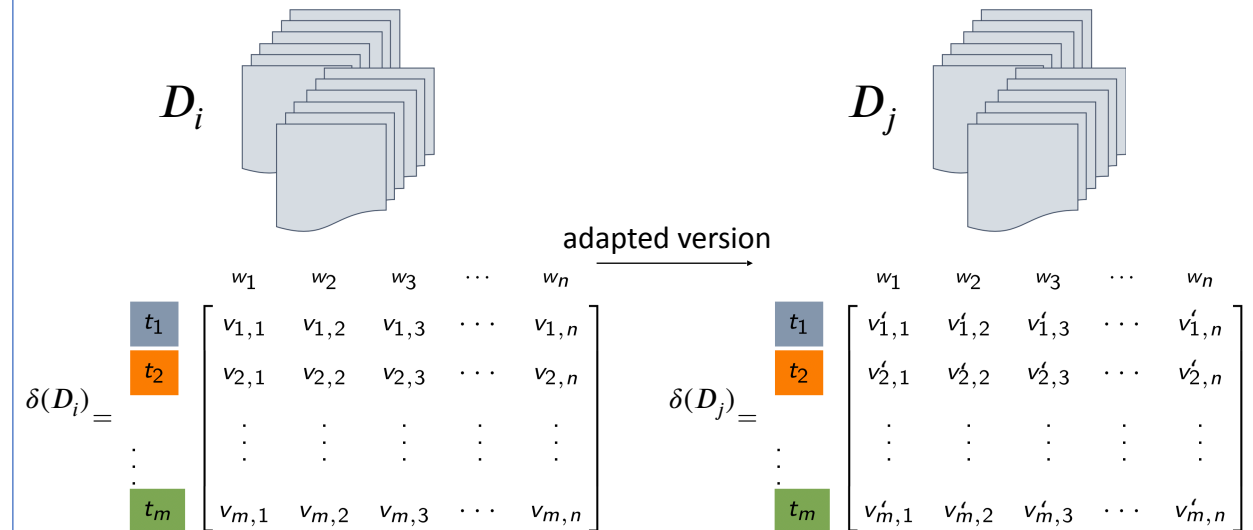
UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Bootstrap Approaches for SCDs

**Inline SCDs** [5]

**Adapting SCD word distribution from another IR-corpus** [6]

- <u>Given:</u> SCD word distribution, trained HMM to detect *inline* SCDs in text

- Estimate MPSCDs and use trained HMM to analyse sequence of corresponding SCD similarity values
  - Small similarity values → different content → new inline-SCDs
  - Inline-SCD = Content of window
  - Inline-SCD represent new row in SCD word matrix
  - HMMs given as user input



$$\delta(D_i) = \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} \begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & \cdots & v_{1,n} \\ v_{2,1} & v_{2,2} & v_{2,3} & \cdots & v_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v_{m,1} & v_{m,2} & v_{m,3} & \cdots & v_{m,n} \end{bmatrix}$$

$$\delta(D_j) = \begin{array}{c} t_1 \\ t_2 \\ \vdots \\ t_m \end{array} \begin{bmatrix} v'_{1,1} & v'_{1,2} & v'_{1,3} & \cdots & v'_{1,n} \\ v'_{2,1} & v'_{2,2} & v'_{2,3} & \cdots & v'_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ v'_{m,1} & v'_{m,2} & v'_{m,3} & \cdots & v'_{m,n} \end{bmatrix}$$

- Adapt SCD word distribution from IR-corpus $D_i$ to documents in $D_j$
  - Analyze difference in word distributions of documents in corpus $D_i$ and $D_j$
  - Reweight word distribution for each SCD in $\delta(D_i)$ s.t. distribution fits for $D_j$
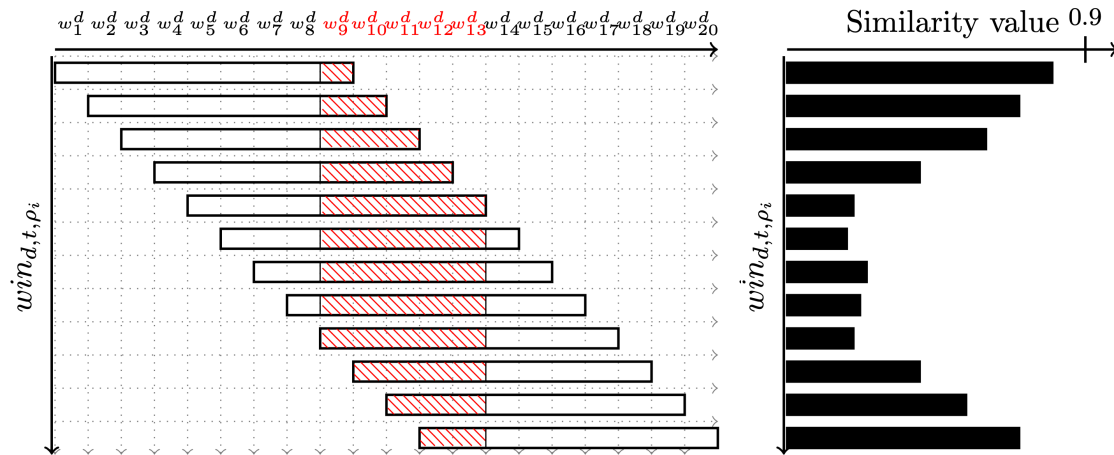
[5] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)

[6] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. 15th IEEE International Conference on Semantic Computing, (ICSC 2021)
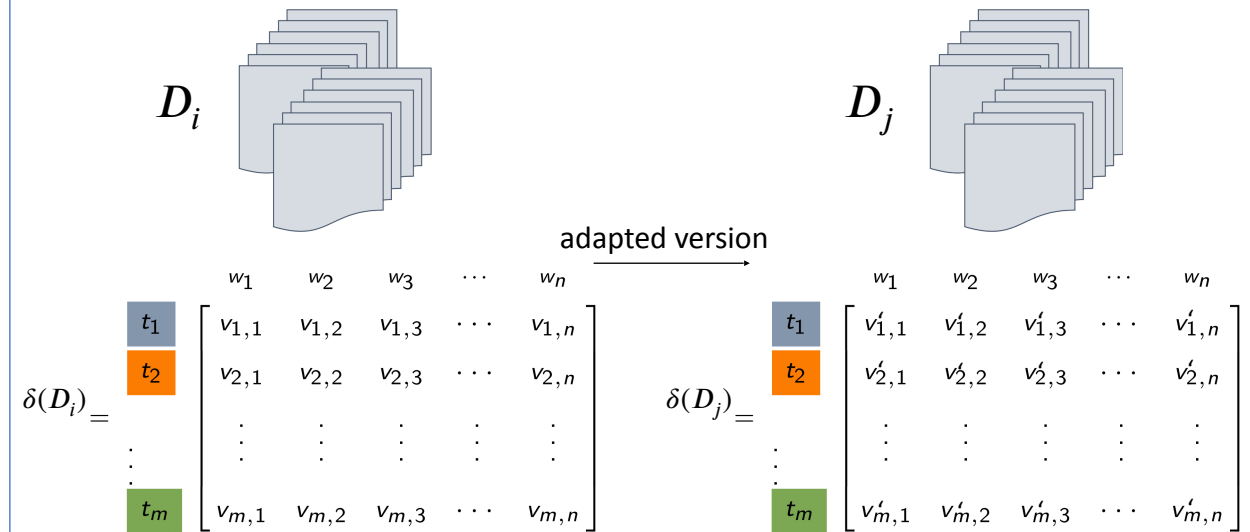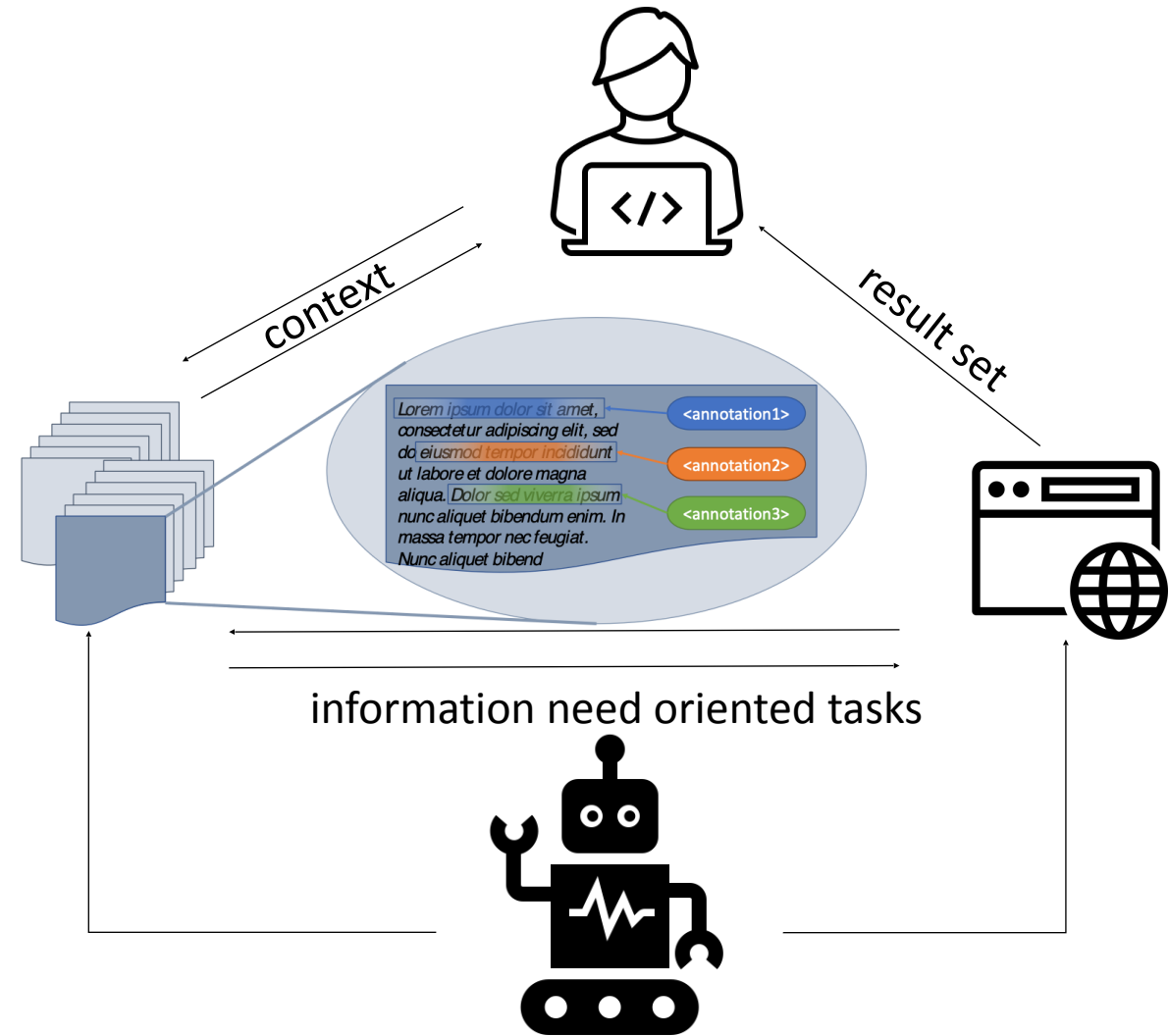
UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

# Conclusion

- Human-aware information retrieval considering <u>not only</u> content of documents and queries

- Fully automated annotation approach considering the human information need represented by a corpus and SCDs

- Approach for the bootstrap problem considering inline-SCDs

context

result set

Lorem *ipsum dolor sit amet*, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor *sed viverra ipsum* nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

<annotation1>

<annotation2>

<annotation3>

information need oriented tasks

# Conclusion

- Human-aware information retrieval considering <u>not only</u> content of documents and queries

- Fully automated annotation approach considering the human information need represented by a corpus and SCDs

- Approach for the bootstrap problem considering inline-SCDs

<span style="color:red">Future Work:</span>

- Focus on Bootstrap mechanisms to generate **new SCDs**
- Deliver a human-aware annotation service

context

result set

Lorem *ipsum dolor sit* amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

<annotation1>

<annotation2>

<annotation3>
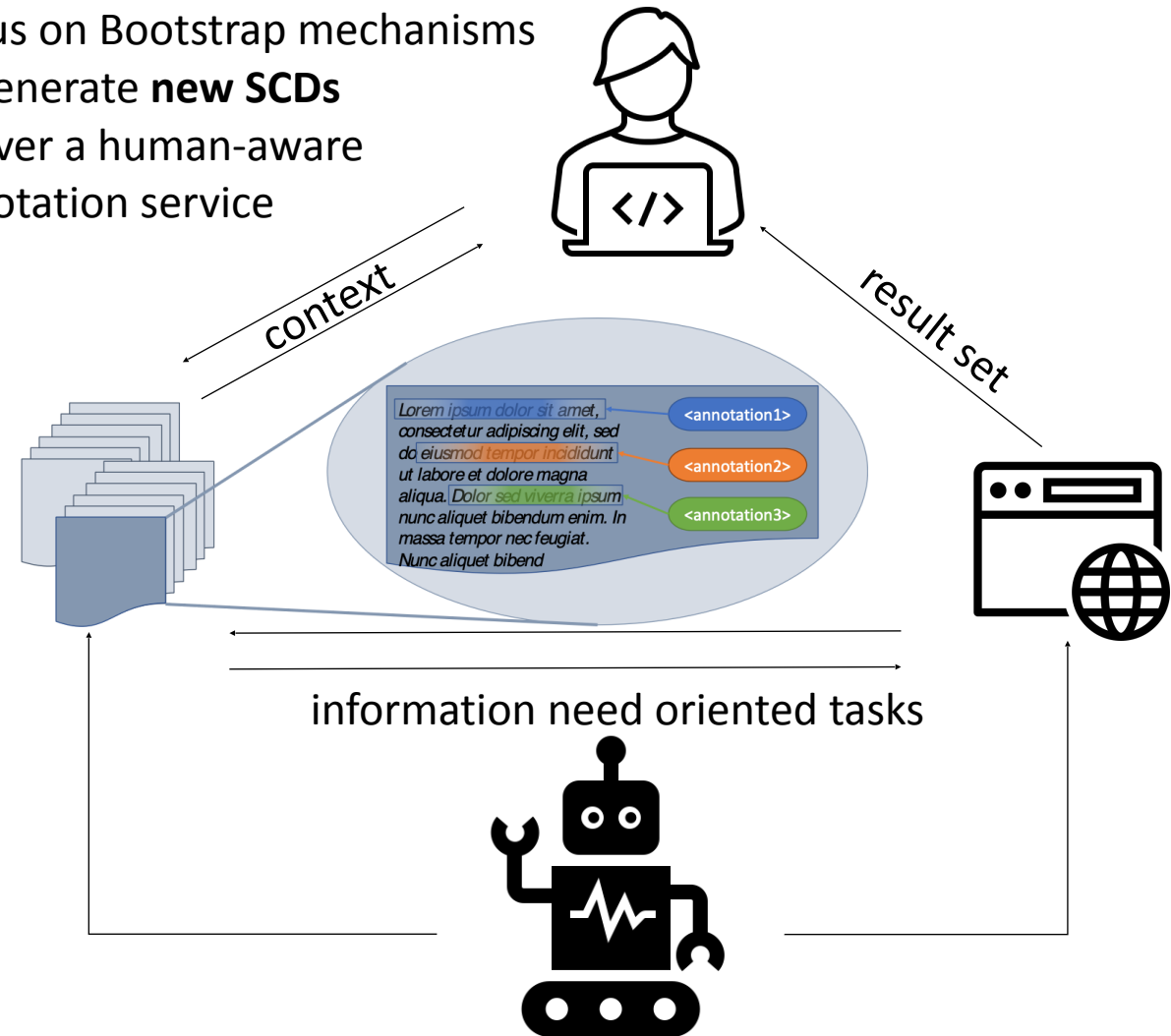
information need oriented tasks

# Conclusion

- Human-aware information retrieval considering <u>not only</u> content of documents and queries

- Fully automated annotation approach considering the human information need represented by a corpus and SCDs

- Approach for the bootstrap problem considering inline-SCDs

<u>Focus on human-aware AI approaches:</u>

**Future Work:**

- Focus on Bootstrap mechanisms to generate **new SCDs**
- Deliver a human-aware annotation service

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

<annotation1>

<annotation2>

<annotation3>
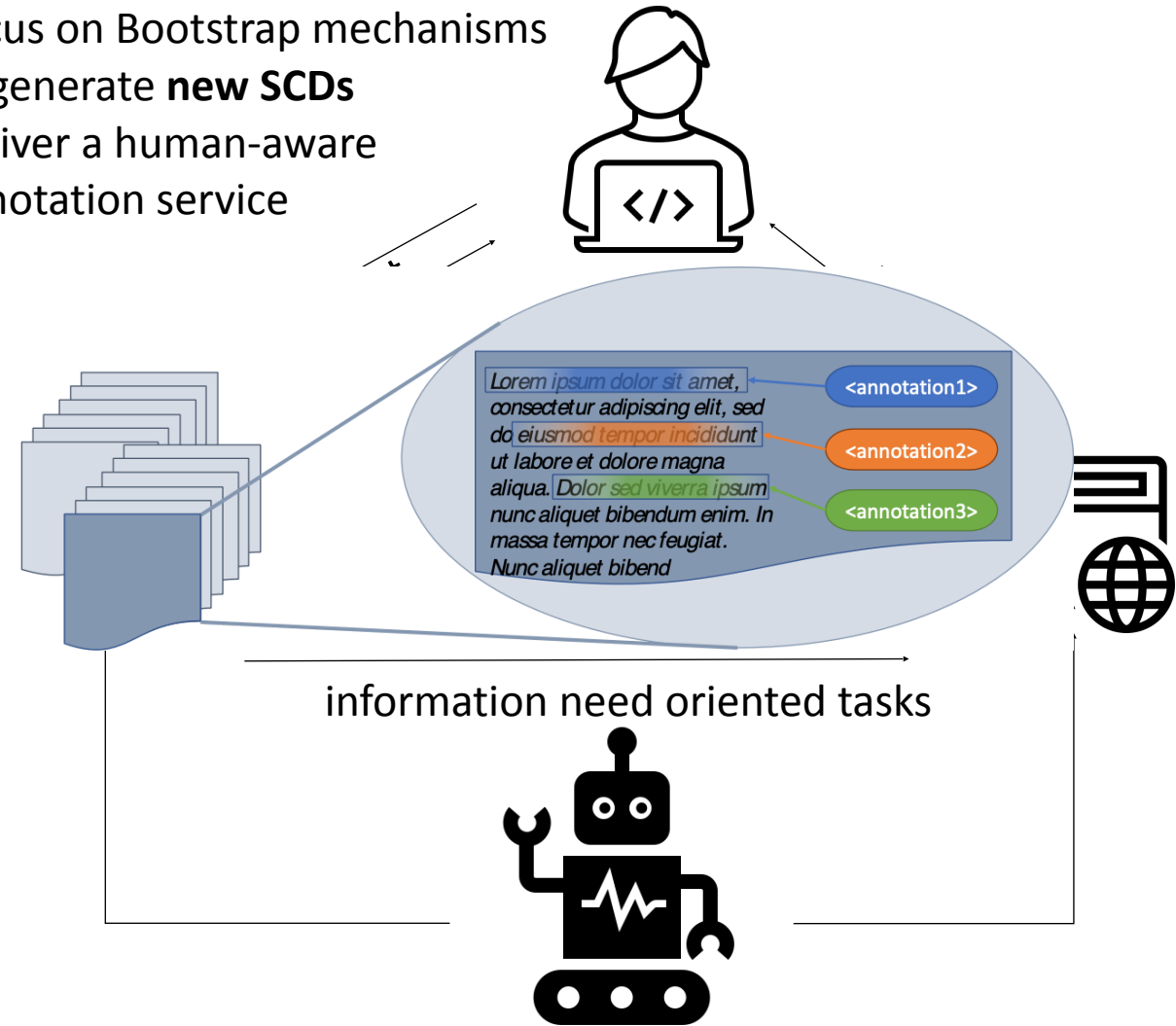
information need oriented tasks
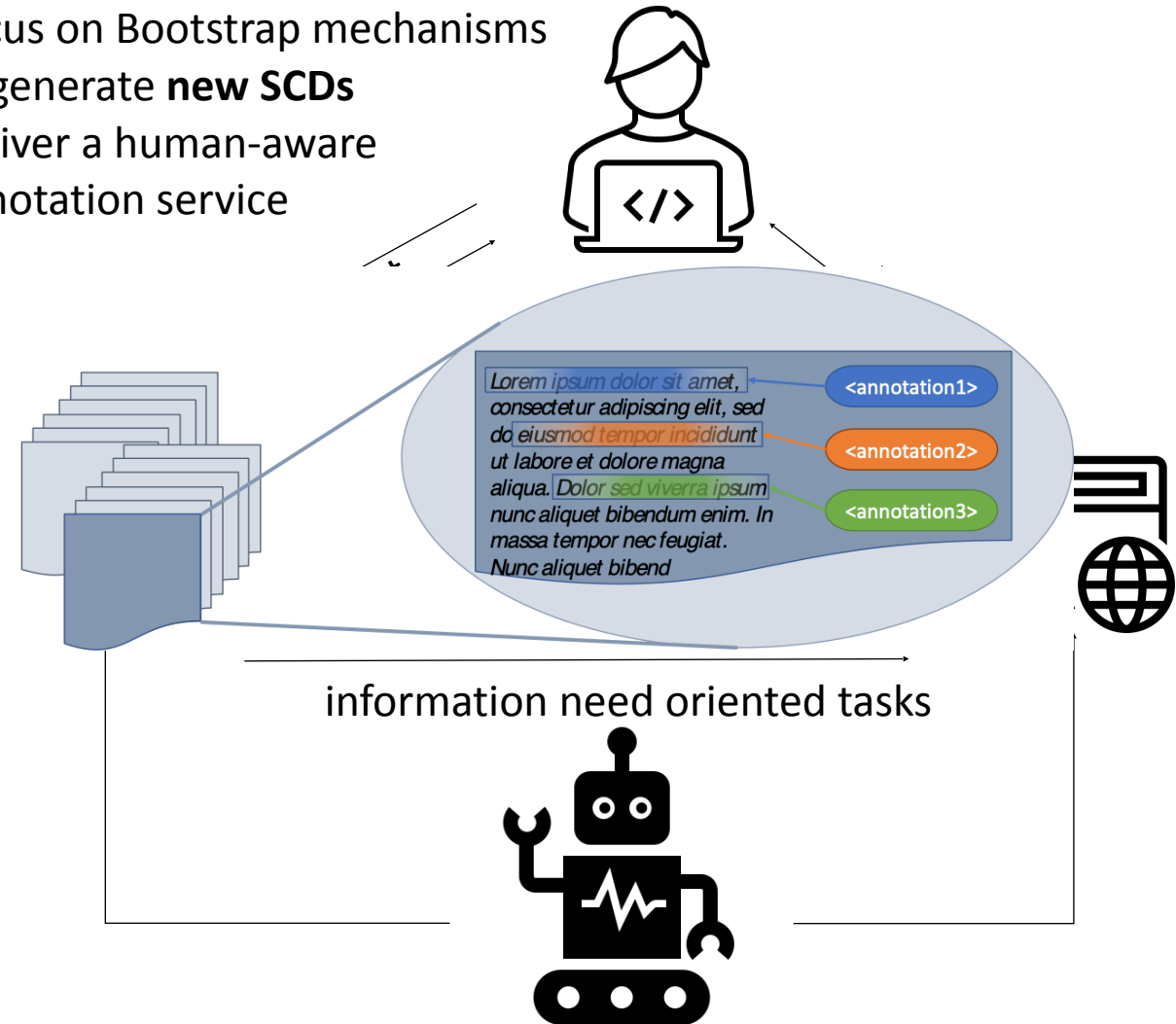
# Conclusion

- Human-aware information retrieval considering <u>not only</u> content of documents and queries

- Fully automated annotation approach considering the human information need represented by a corpus and SCDs

- Approach for the bootstrap problem considering inline-SCDs

<u>Focus on human-aware AI approaches:</u>
→ Data linking **services** in a fashion that takes into aware **human expectations**

**Future Work:**
- Focus on Bootstrap mechanisms to generate **new SCDs**
- Deliver a human-aware annotation service



Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Dolor sed viverra ipsum nunc aliquet bibendum enim. In massa tempor nec feugiat. Nunc aliquet bibend

<annotation1>
<annotation2>
<annotation3>

information need oriented tasks

# Referenzen

[1] Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." *Journal of machine Learning research* 3.Jan (2003): 993-1022.

[2] Felix Kuhr, Tanya Braun, Magnus Bender, Ralf Möller: To Extend or not to Extend? Context-specific Corpus Enrichment. *Proceedings of AI 2019: Advances in Artificial Intelligence, 2019, Springer, p.357-368*

[3] Felix Kuhr, Tanya Braun, Ralf Möller: Augmenting and Automating Corpus Enrichment. *Proceedings of the 14th IEEE International Conference on Semantic Computing (ICSC-20), 2020, Best Student Paper Award*

[4] Felix Kuhr, Bjarne Witten, Ralf Möller: Corpus-Driven Annotation Enrichment. *Proceedings of the 13th IEEE International Conference on Semantic Computing (ICSC-19), 2019, Jan, p.138-141*

# Referenzen

[5] Magnus Bender, Tanya Braun, Marcel Gehrke, Felix Kuhr, Ralf Möller, Simon Schiff: Identifying Subjective Content Descriptions Among Texts. *15th IEEE International Conference on Semantic Computing, (ICSC 2021), Laguna Hills, CA, USA, January 27-29, 2021, IEEE, p.9-16*


[6] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Context-specific Adaptation of Subjective Content Descriptions. *15th IEEE International Conference on Semantic Computing, (ICSC 2021), Laguna Hills, CA, USA, January 27-29, 2021, IEEE, p.134-139*


[7] Felix Kuhr, Magnus Bender, Tanya Braun, Ralf Möller: Maintaining Topic Models for Growing Corpora. *Proceedings of the 14th IEEE International Conference on Semantic Computing (ICSC-20), 2020*