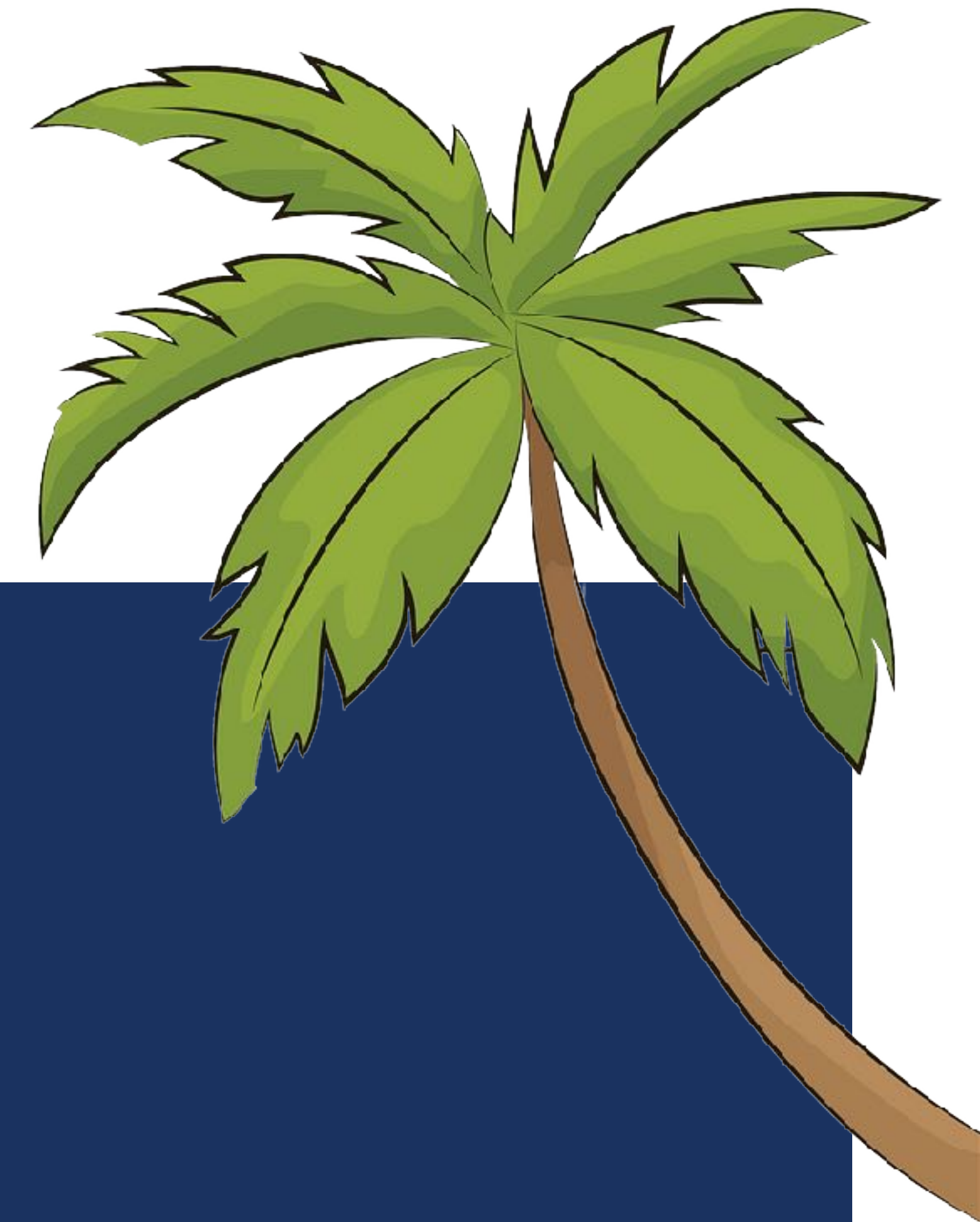


---

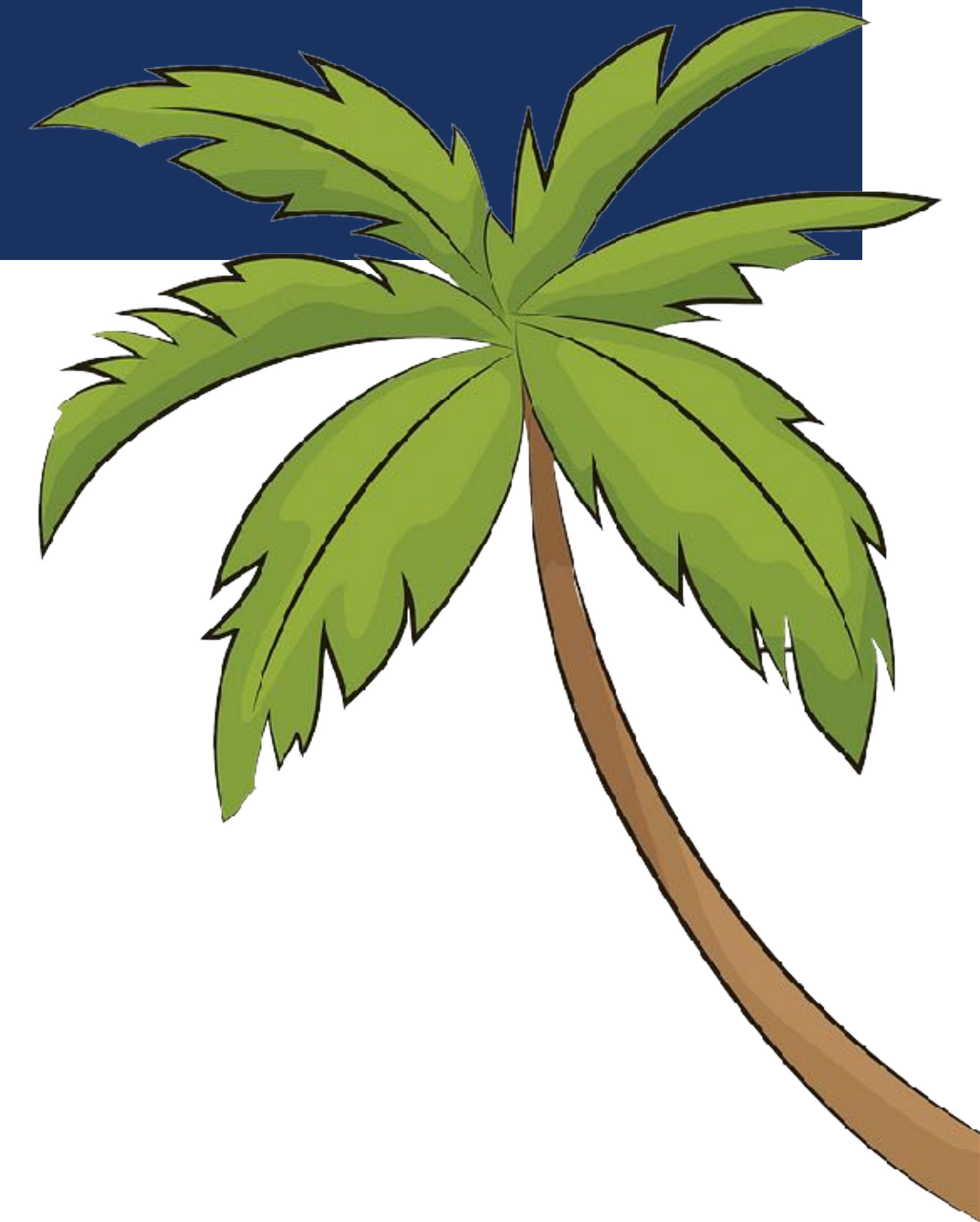
# LET'S TALK ABOUT PALM LEAVES FROM MINIMAL DATA TO TEXT UNDERSTANDING

MAGNUS BENDER<sup>1</sup>, MARCEL GEHRKE<sup>1</sup>, TANYA BRAUN<sup>2</sup>



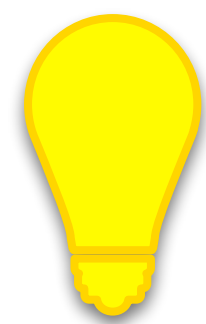
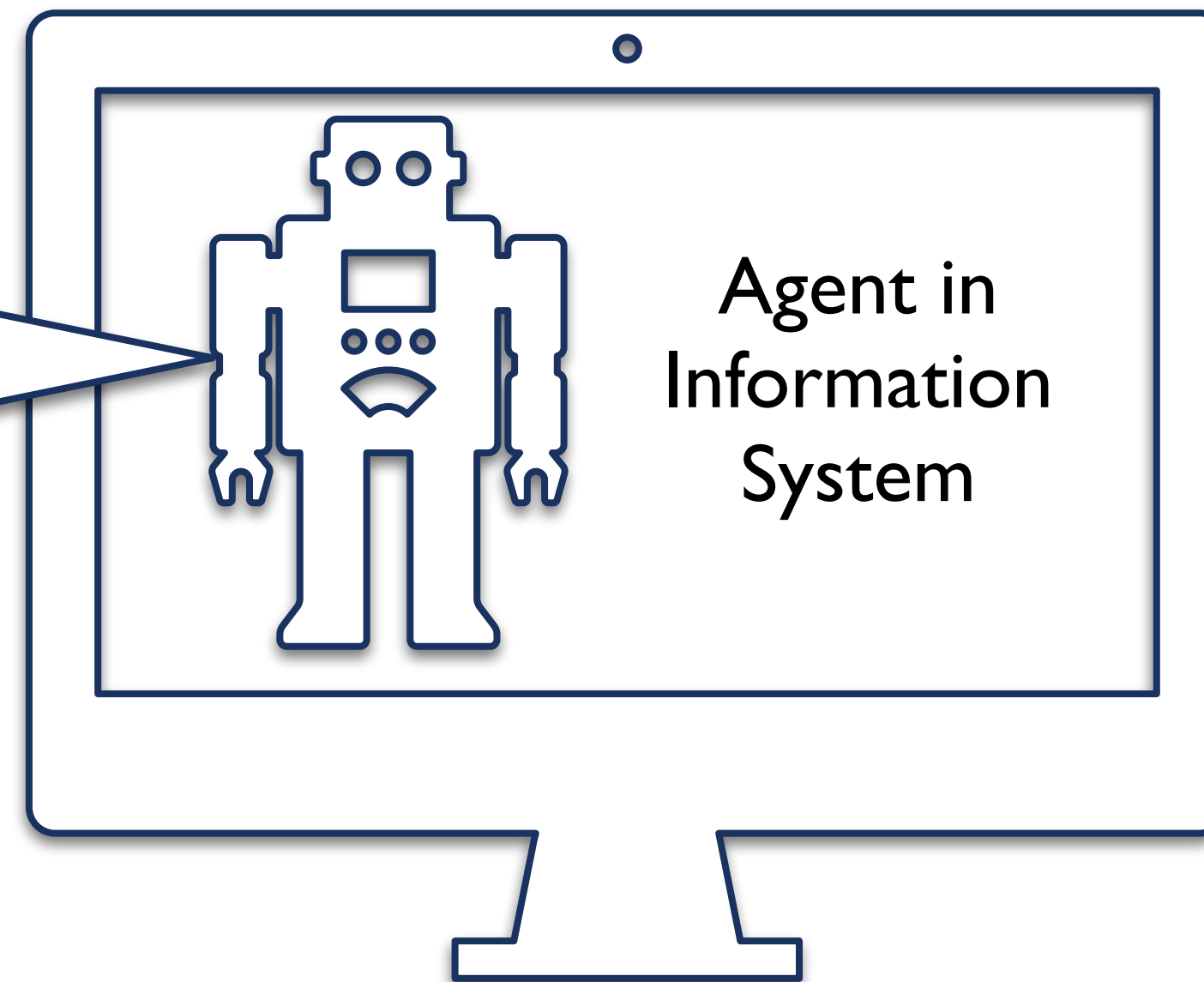
# AGENDA

1. Introduction to Semantic Systems [Tanya]
2. Supervised Learning [Marcel]
3. Unsupervised and Relational Learning [Magnus]
  - Unsupervised Estimation of SCDs
  - Continuous Improvement by Feedback
  - Labelling of SCDs
  - Inter- and Intra-SCD Relations
4. Summary [Tanya]



# SCENARIO

- Tasks, e.g.,
  - Information Retrieval
  - Corpus Enrichment
- Techniques
- SCDs

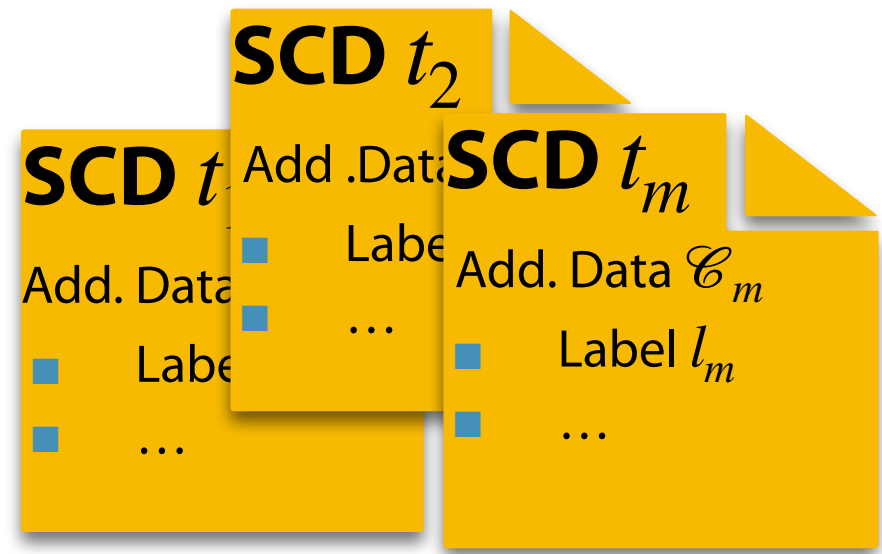


Any corpus brought e.g. by human.



# OVERVIEW

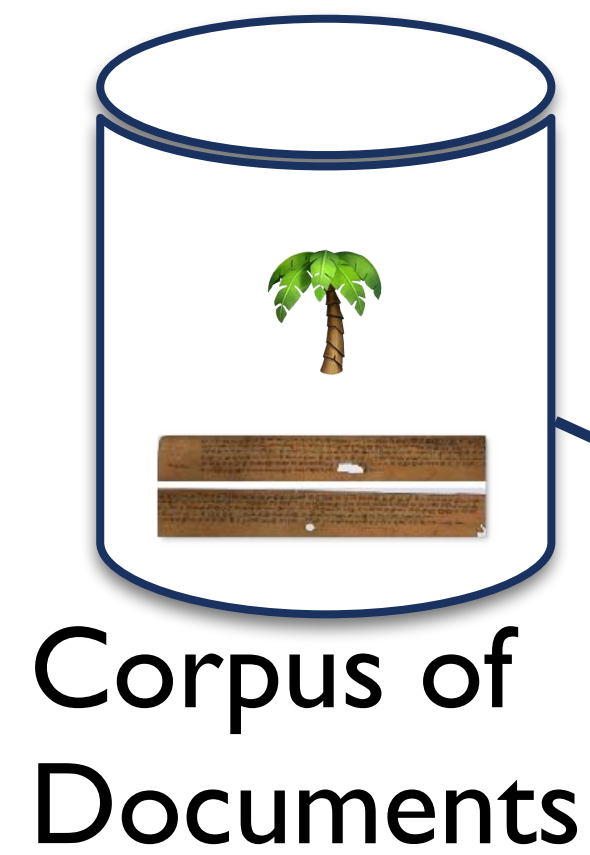
No initial SCDs



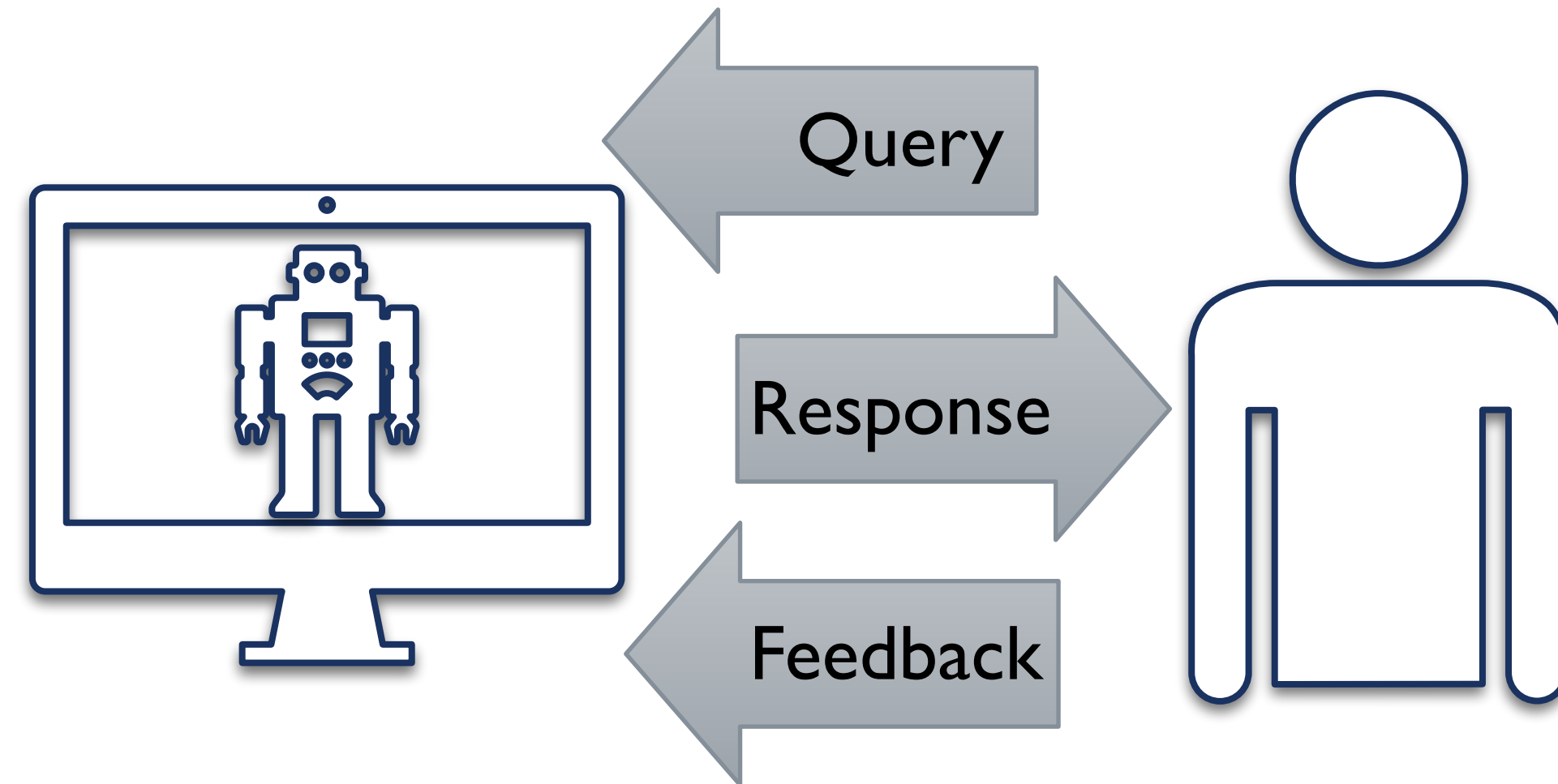
SCDs  $g(\mathcal{D})$

$$\begin{matrix} & w_1 & w_2 & \dots & w_n \\ \begin{matrix} t_1 \\ t_2 \\ \vdots \\ t_m \end{matrix} & \left( \begin{matrix} v_{1,1} & v_{1,2} & \dots & v_{1,n} \\ v_{2,1} & v_{2,2} & \dots & v_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ v_{m,1} & v_{m,2} & \dots & v_{m,n} \end{matrix} \right)
 \end{matrix}$$

How to incorporate Feedback?



Used to Respond to Queries

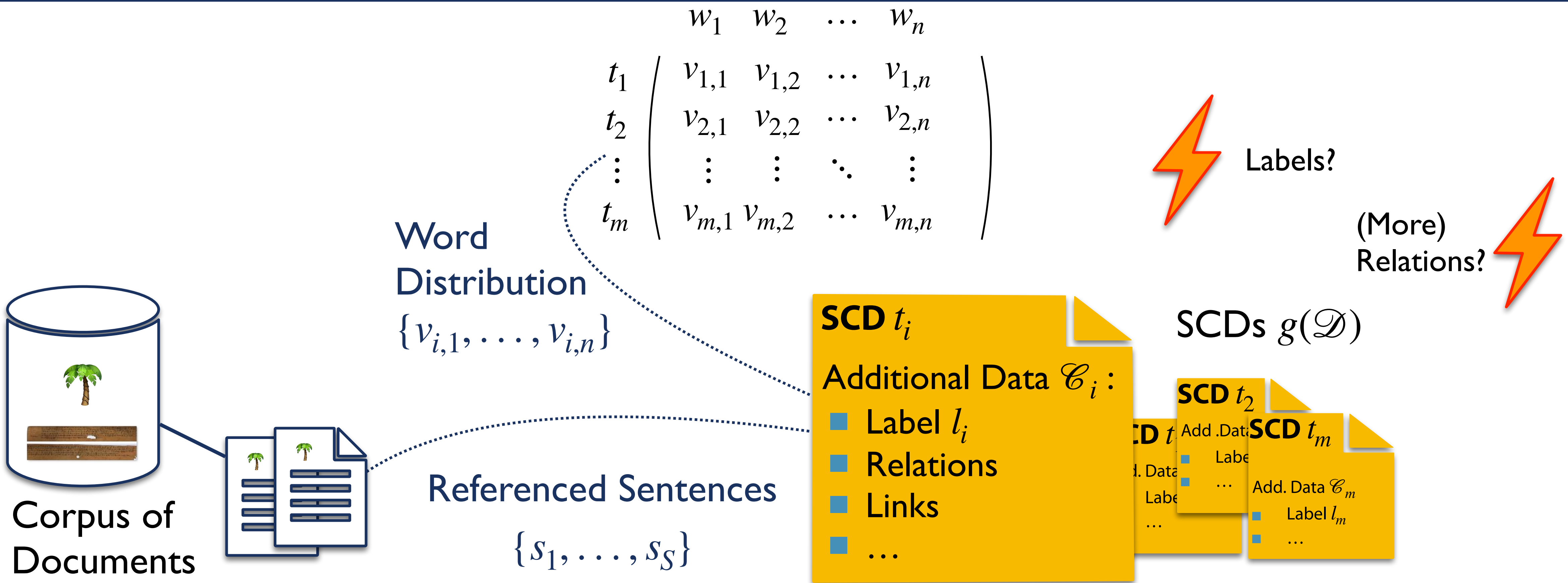


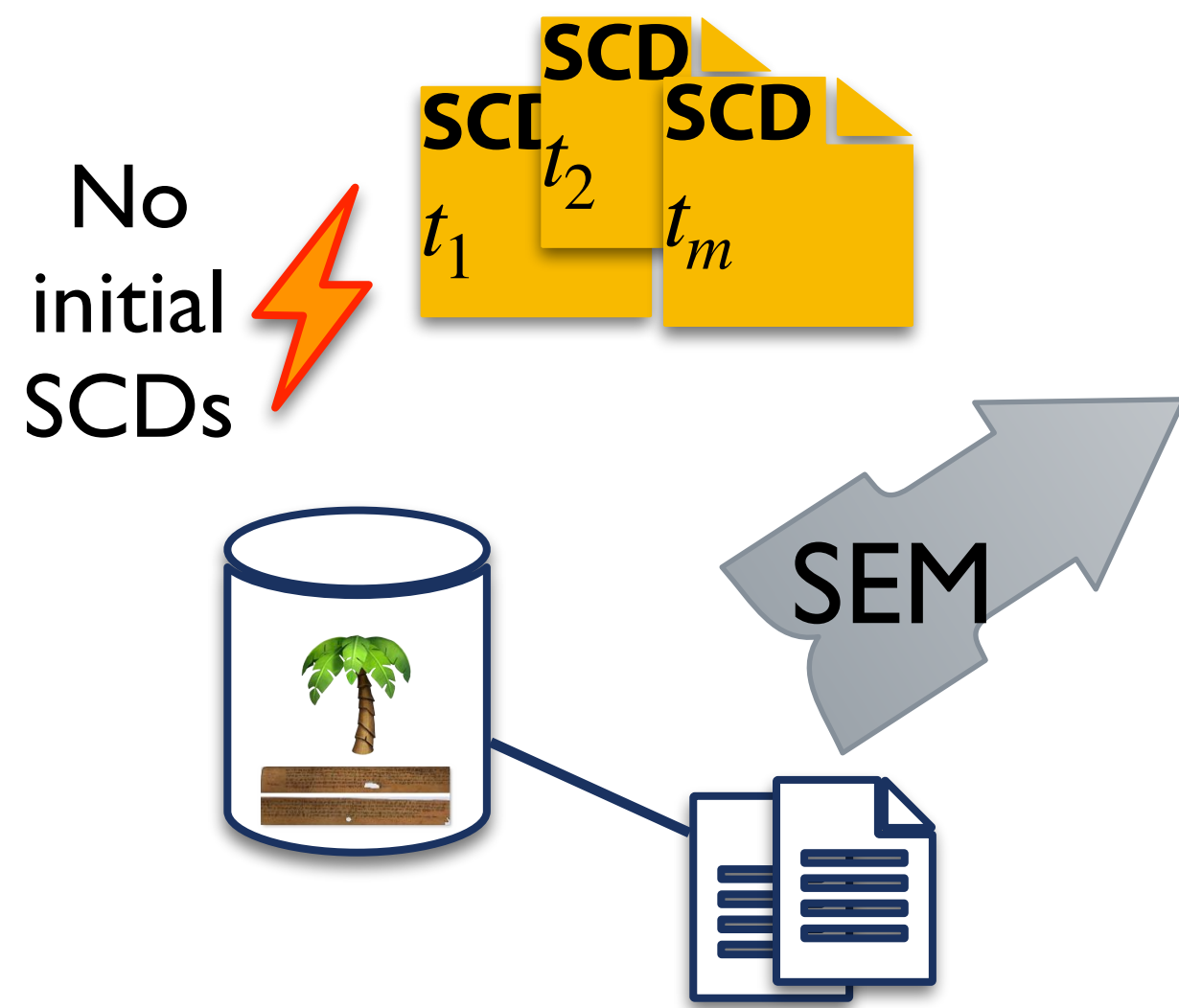
FROM MINIMAL DATA TO TEXT UNDERSTANDING

Unsupervised and Relational Learning

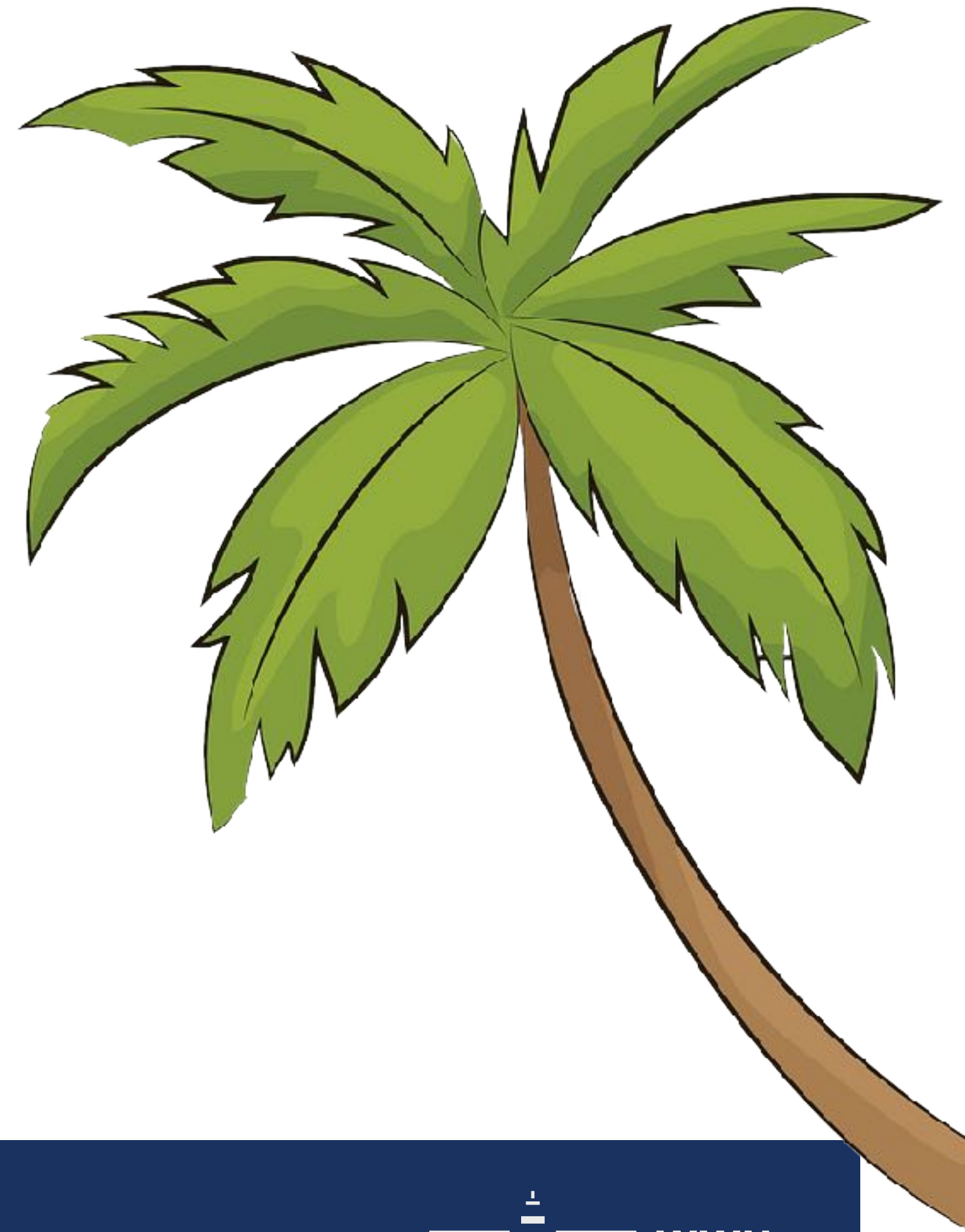


# SCD IN DETAIL





$$\begin{matrix}
 & w_1 & w_2 & \dots & w_n \\
 t_1 & \left( \begin{matrix} v_{1,1} & v_{1,2} & \dots & v_{1,n} \\
 t_2 & v_{2,1} & v_{2,2} & \dots & v_{2,n} \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 t_m & v_{m,1} & v_{m,2} & \dots & v_{m,n} \end{matrix} \right)
 \end{matrix}$$



# UNSUPERVISED ESTIMATION OF SCDS

USEM – UNSUPERVISED ESTIMATION OF SCD MATRICES



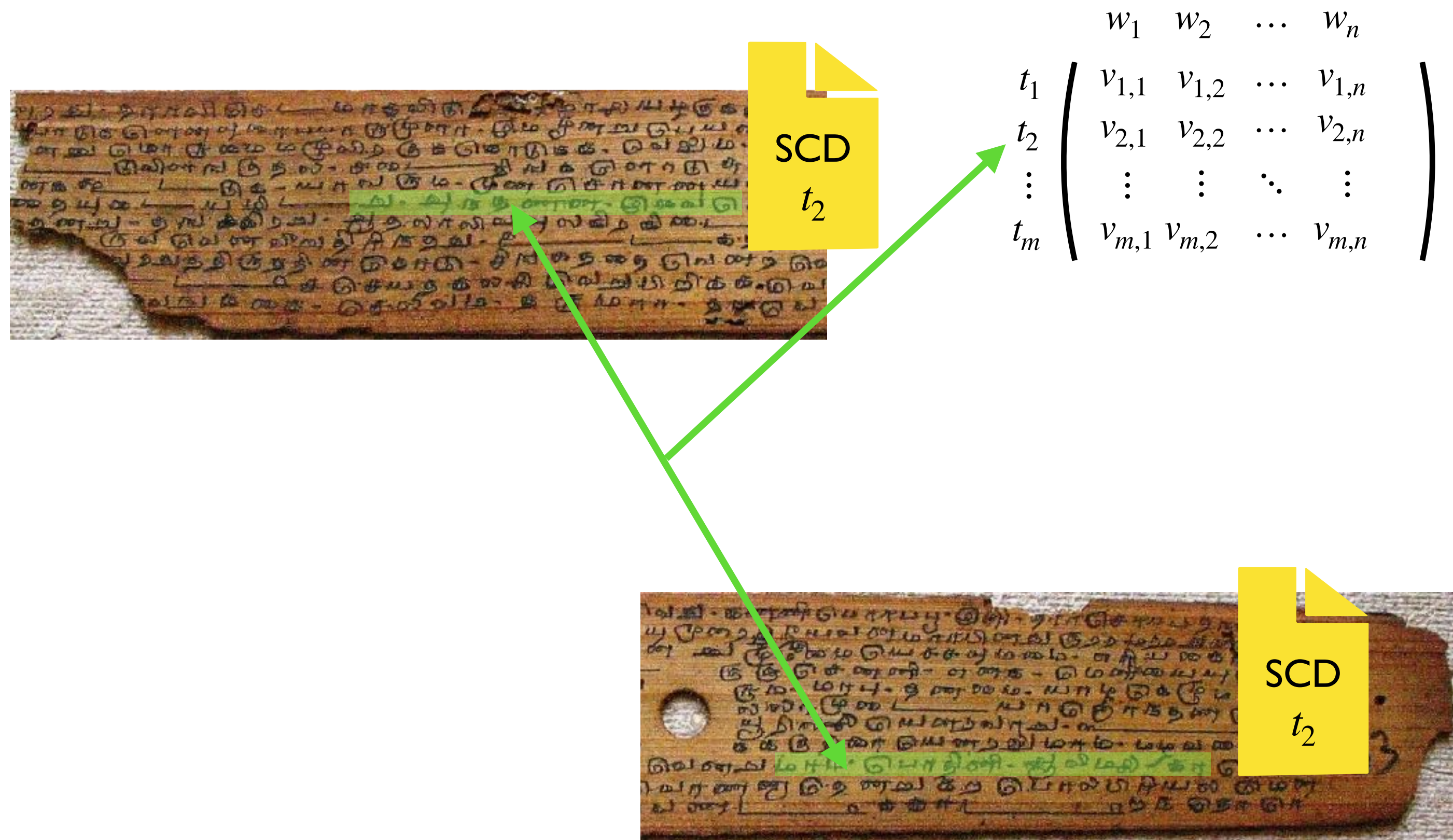
FROM MINIMAL DATA TO TEXT UNDERSTANDING

Unsupervised and Relational Learning



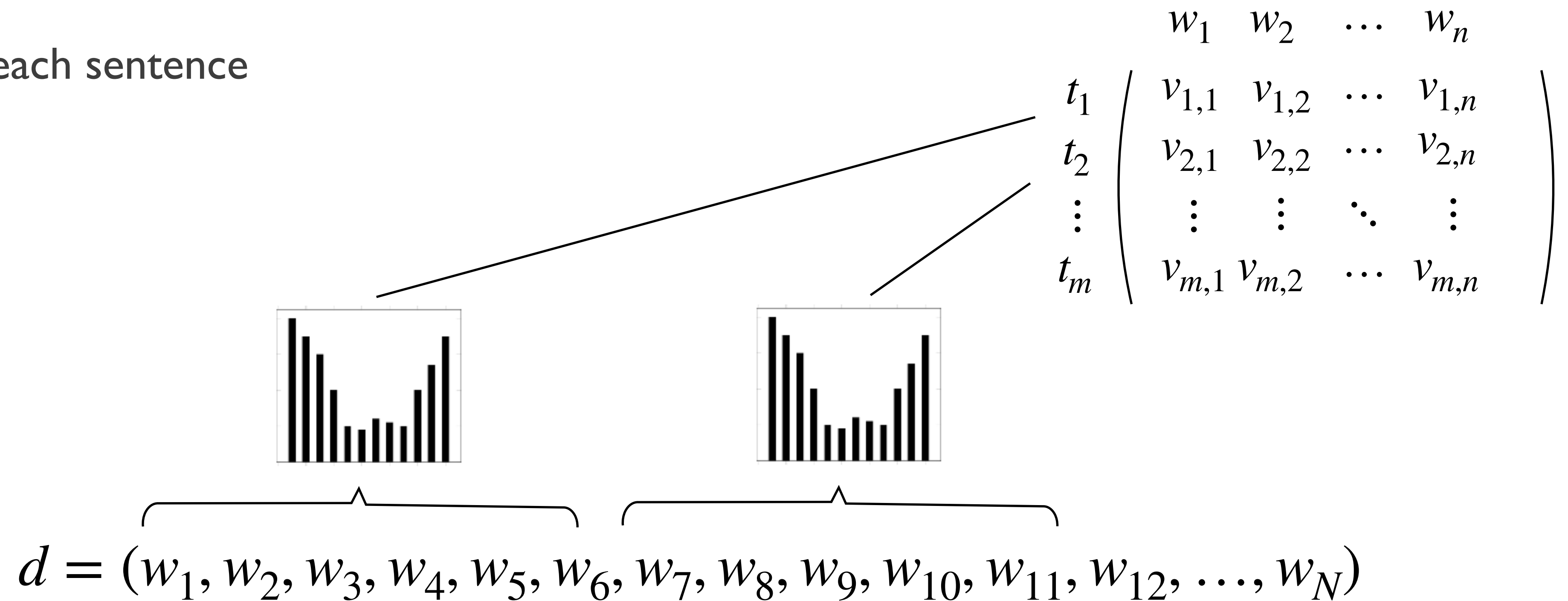
# UNSUPERVISED ESTIMATION OF SCD MATRICES

- Estimate SCDs in an unsupervised manner
- Focus on identifying similar sentences
- Estimate an SCD matrix
- Select the *best* from multiple matrices



# IDEA: USEM

I. Initially, one SCD for each sentence





# IDEA: USEM

1. Initially, one SCD for each sentence

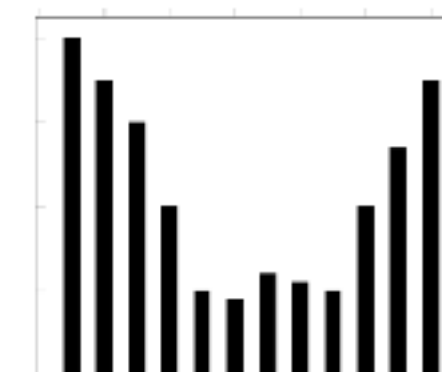
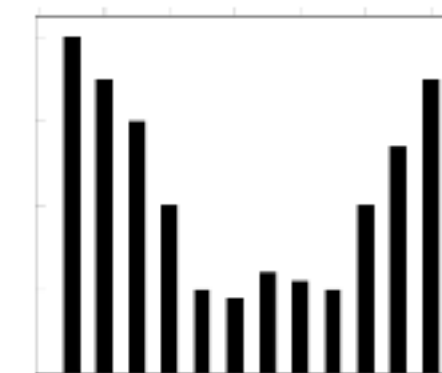
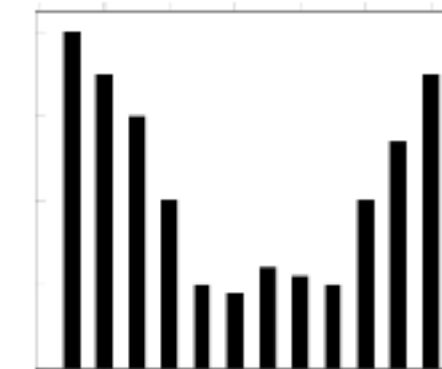
2. Identify similar distributions (sentences)

- Greedy

- K-Means

- DBSCAN

$$\begin{matrix}
 & w_1 & w_2 & \dots & w_n \\
 t_1 & \left( \begin{matrix} v_{1,1} & v_{1,2} & \dots & v_{1,n} \\
 t_2 & \begin{matrix} v_{2,1} & v_{2,2} & \dots & v_{2,n} \\
 \vdots & \begin{matrix} \vdots & \vdots & \ddots & \vdots \\
 t_m & \begin{matrix} v_{m,1} & v_{m,2} & \dots & v_{m,n} \end{matrix} \end{matrix} \right)
 \end{matrix}$$



Only K rows, multiple sentences represented in each SCD

Similar?

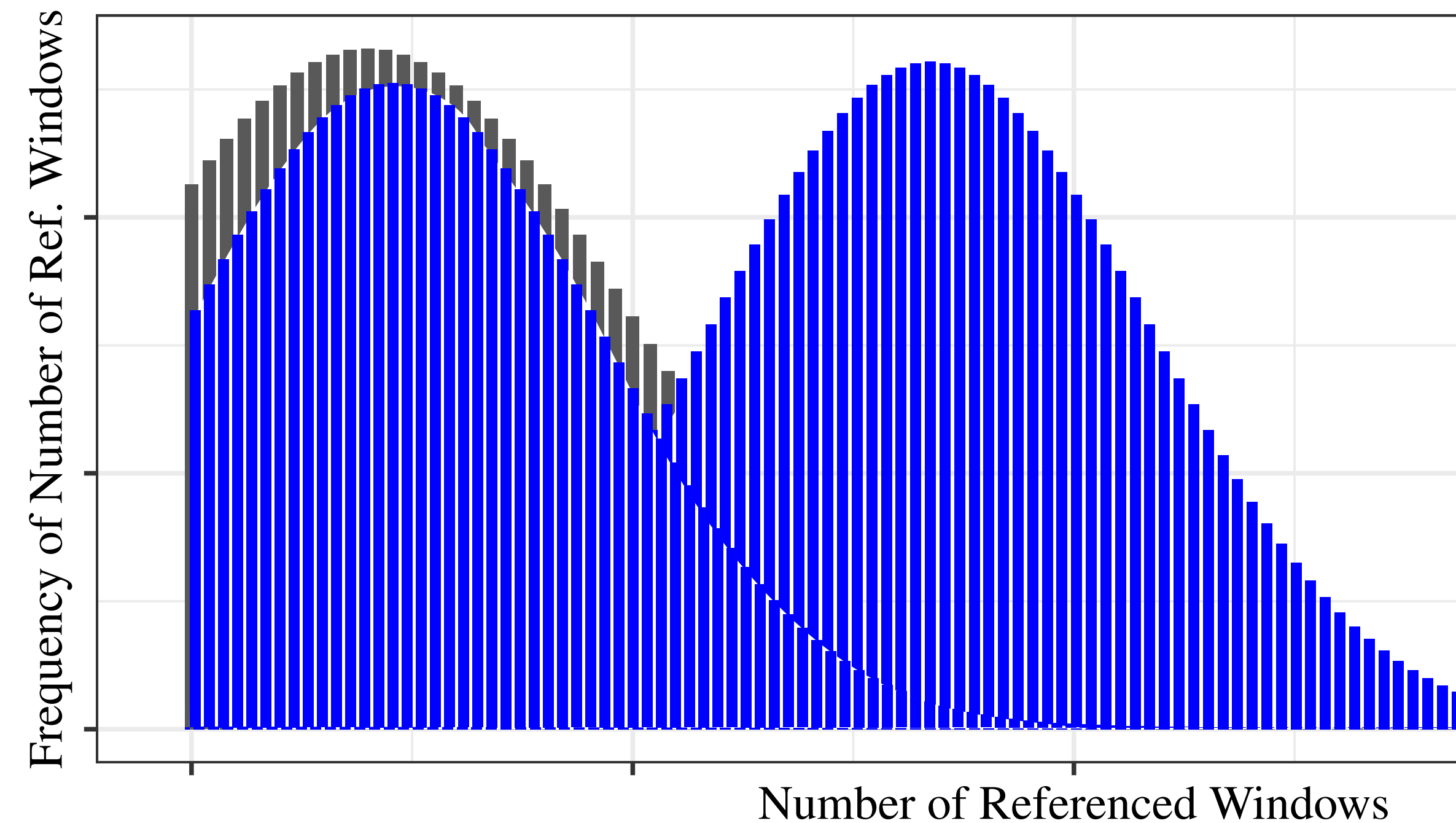
Similar?

$$\begin{matrix}
 & w_1 & w_2 & \dots & w_n \\
 t_1 & \left( \begin{matrix} v_{1,1} & v_{1,2} & \dots & v_{1,n} \\
 t_2 & \begin{matrix} v_{2,1} & v_{2,2} & \dots & v_{2,n} \\
 \vdots & \begin{matrix} \vdots & \vdots & \ddots & \vdots \\
 t_K & \begin{matrix} v_{K,1} & v_{K,2} & \dots & v_{K,n} \end{matrix} \end{matrix} \end{matrix} \right)
 \end{matrix}$$

3. Merge the similar sentences (incrementally)

# SCD MATRIX MODEL SELECTION

- **Problem:** Three Methods → multiple matrices
- **Goal:** Identify best hyperparameters for USEM
  - Method (one of DBSCAN, K-Means, Greedy) and
  - Hyperparameters for method
- **Idea:** Run USEM multiple times and choose best resulting matrix
- **Quality Score:** Similarity to optimal histogram depicting the different numbers of windows referenced in an SCD matrix →



Gray assumed optimal, left blue best matrix, right blue weaker matrix.

**Algorithm** UnSupervised Estimator for SCD Matrices  $\delta(\mathcal{D})$ 

```

1: function USEM( $\mathcal{D}$ ,  $m$ , [ $\theta$ ,] [ $K$ ,] [ $\varepsilon$ ,  $ms$ ])
2:   Input: Corpus  $\mathcal{D}$ ; Method with hyperparameters, i.e.,
3:      $m = \text{Greedy}$  and threshold  $\theta$ ,
4:      $m = \text{K-Means}$  and number of SCDs  $K$ , or
5:      $m = \text{DBSCAN}$ , distance  $\varepsilon$ , and threshold  $ms$ 
6:   Output: SCD-word distribution matrix  $\delta(\mathcal{D})$ 
7:   Initialize an  $(\sum_{d \in \mathcal{D}} M^d) \times L$  matrix  $\delta(\mathcal{D})$  with zeros
8:    $l \leftarrow 0$ 
9:                                      $\triangleright$  Build initial SCD matrix
10:  for each document  $d \in \mathcal{D}$  do
11:    for each sentence  $s^d \in d$  do
12:      for each word  $w_i \in s^d$  do
13:         $\delta(\mathcal{D})[l][w_i] += I(w_i, s^d)$ 
14:       $l \leftarrow l + 1$ 
15:                                      $\triangleright$  Use method  $m$  to merge rows
16:  if  $m = \text{Greedy}$  then
17:    repeat                                      $\triangleright$  Detect similar rows and merge
18:       $(r_i, r_j) \leftarrow \text{MOSTSIMILARROWS}(\delta(\mathcal{D}))$ 
19:       $\delta(\mathcal{D})[r_i] \leftarrow \delta(\mathcal{D})[r_i] + \delta(\mathcal{D})[r_j]$             $\triangleright$  Sum rows
20:       $\delta(\mathcal{D})[r_j] \leftarrow \text{Nil}$                                       $\triangleright$  Delete row
21:    until  $\text{SIMILARITY}(r_i, r_j) < \theta$ 
22:  else                                      $\triangleright$  Create clusters of similar rows
23:    if  $m = \text{K-Means}$  then
24:       $clusters \leftarrow \text{KMEANS}(\delta(\mathcal{D}), K)$ 
25:    else
26:       $clusters \leftarrow \text{DBSCAN}(\delta(\mathcal{D}), \varepsilon, ms)$ 
27:    for each cluster  $c \in clusters$  do
28:                                      $\triangleright$  Create sum of all cluster's rows in first row
29:       $r_i \leftarrow \text{FIRSTROW}(c)$ 
30:       $\delta(\mathcal{D})[r_i] \leftarrow \sum_{r_j \in c} \delta(\mathcal{D})[r_j]$ 
31:      for each row  $r_j \in c$  do
32:        if  $r_i \neq r_j$  then                                      $\triangleright$  Delete all non-first rows
33:           $\delta(\mathcal{D})[r_j] \leftarrow \text{Nil}$ 
34:  return  $\delta(\mathcal{D})$ 

```

**Algorithm** SCD Matrix Model Selection

```

1: function ESTIMATEBESTMATRIX( $\mathcal{D}$ )
2:   Input: Corpus  $\mathcal{D}$ 
3:   Output: Best SCD-word distribution matrix  $\delta(\mathcal{D})$ 
4:    $sim_{best} \leftarrow 0$ 
5:    $\delta_{best} \leftarrow \text{Nil}$                                       $\triangleright$  Iterate all methods
6:   for each method  $m \in \{\text{Greedy}, \text{K-Means}, \text{DBSCAN}\}$  do
7:                                      $\triangleright$  Take a set of hyperparameters depending on method
8:     if  $m = \text{Greedy}$  then
9:        $H \leftarrow (0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1)$             $\triangleright$  Values of  $\theta$ 
10:    if  $m = \text{K-Means}$  then
11:       $M' \leftarrow \sum_{d \in \mathcal{D}} M^d$             $\triangleright$  Number of sentences in  $\mathcal{D}$  to calculate  $K$ s
12:       $H \leftarrow ([M' \cdot 0.8], [M' \cdot 0.6], [M' \cdot 0.4], [M' \cdot 0.3], [M' \cdot 0.2], [M' \cdot 0.1])$ 
13:    else
14:       $H \leftarrow ((0.3, 1), (0.5, 10), (0.5, 5), (0.5, 2), (0.7, 10))$     $\triangleright$  Tuples of  $\varepsilon, ms$ 
15:    for each hyperparameter  $h \in H$  do
16:       $\delta(\mathcal{D}) \leftarrow \text{USEM}(\mathcal{D}, m, h)$                                       $\triangleright$  Run USEM
17:       $\triangleright$  Calculate score using Hellinger distance to normal distribution
18:       $sim \leftarrow 1 - \text{HD}(\text{SCALE}([0, 100], \delta(\mathcal{D})), \mathcal{N}([0, 100], \mu = 10, \sigma^2 = 15))$ 
19:      if  $sim > sim_{best}$  then
20:         $sim_{best} \leftarrow sim$ 
21:         $\delta_{best} \leftarrow \delta(\mathcal{D})$ 
22:  return  $\delta_{best}$ 

```



# EVALUATION & EXAMPLE: DATASET

- Bürgerliches Gesetzbuch (BGB)
  - German civil code (German language)
- Why BGB?
  - Easily to process
  - Uniform style of writing
- Identify and present similar paragraphs
  - Compare USEM to LDA topic model
- Only example for a corpus

Bürgerliches  
Gesetzbuch

BGB

**Buch 1**  
*Allgemeiner Teil*  
§§ 1 - 240

2022

BGB

**Buch 2**  
*Recht der Schuldverhältnisse*  
§§ 241 - 853

2022

Bürgerliches  
Gesetzbuch

BGB

Bürgerliches  
Gesetzbuch

BGB

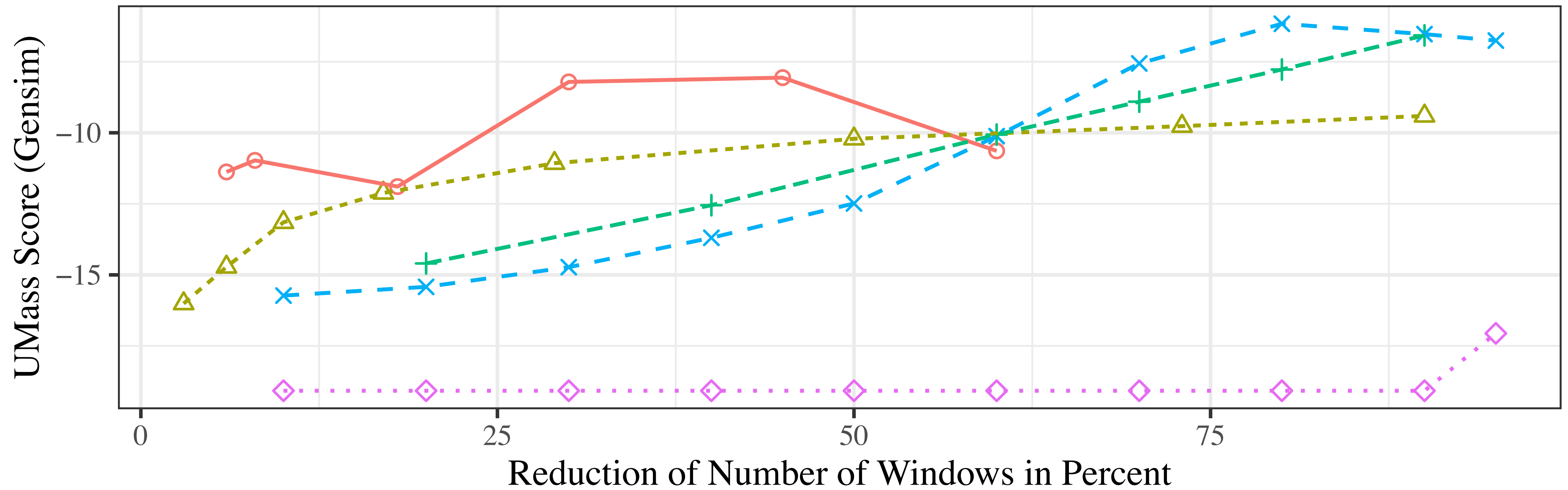
Bürgerliches  
Gesetzbuch

BGB

**Buch 5**  
*Erbrecht*  
§§ 1922 - 2385

2022

# USEM VS. LDA



—○— DBSCAN    -△- Greedy by Similarity    -+· K-Means    -x- LDA (Documents)    ·◇· LDA (Windows)

Well-know competitor

Provides SCD matrix

# USAGE EXAMPLE

„An association whose purpose is not directed towards a commercial business operation acquires legal capacity through entry in the register of associations at the competent local court.“

„An association whose purpose is to engage in commercial business shall acquire legal capacity, in the absence of special federal law, through state conferral. The grant is due to the state in whose territory the association has its registered office.“

„The seat of an association, unless otherwise provided, is the place where the administration is conducted.“

„The seat of a foundation, unless otherwise provided, is the place where the administration is conducted.“

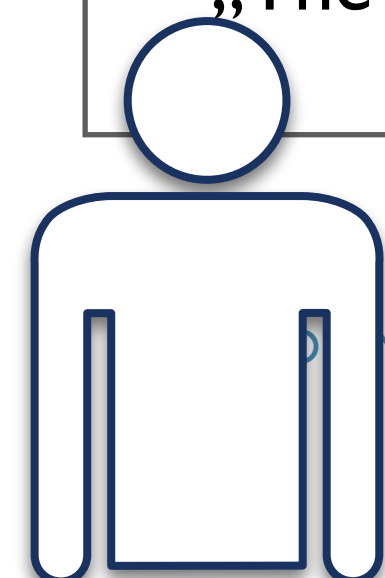
Find similar sentences in SCD matrix

	$w_1$	$w_2$	$\dots$	$w_n$
$t_1$	$v_{1,1}$	$v_{1,2}$	$\dots$	$v_{1,n}$
$t_2$	$v_{2,1}$	$v_{2,2}$	$\dots$	$v_{2,n}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$t_m$	$v_{m,1}$	$v_{m,2}$	$\dots$	$v_{m,n}$

Find sentences referenced by same SCD

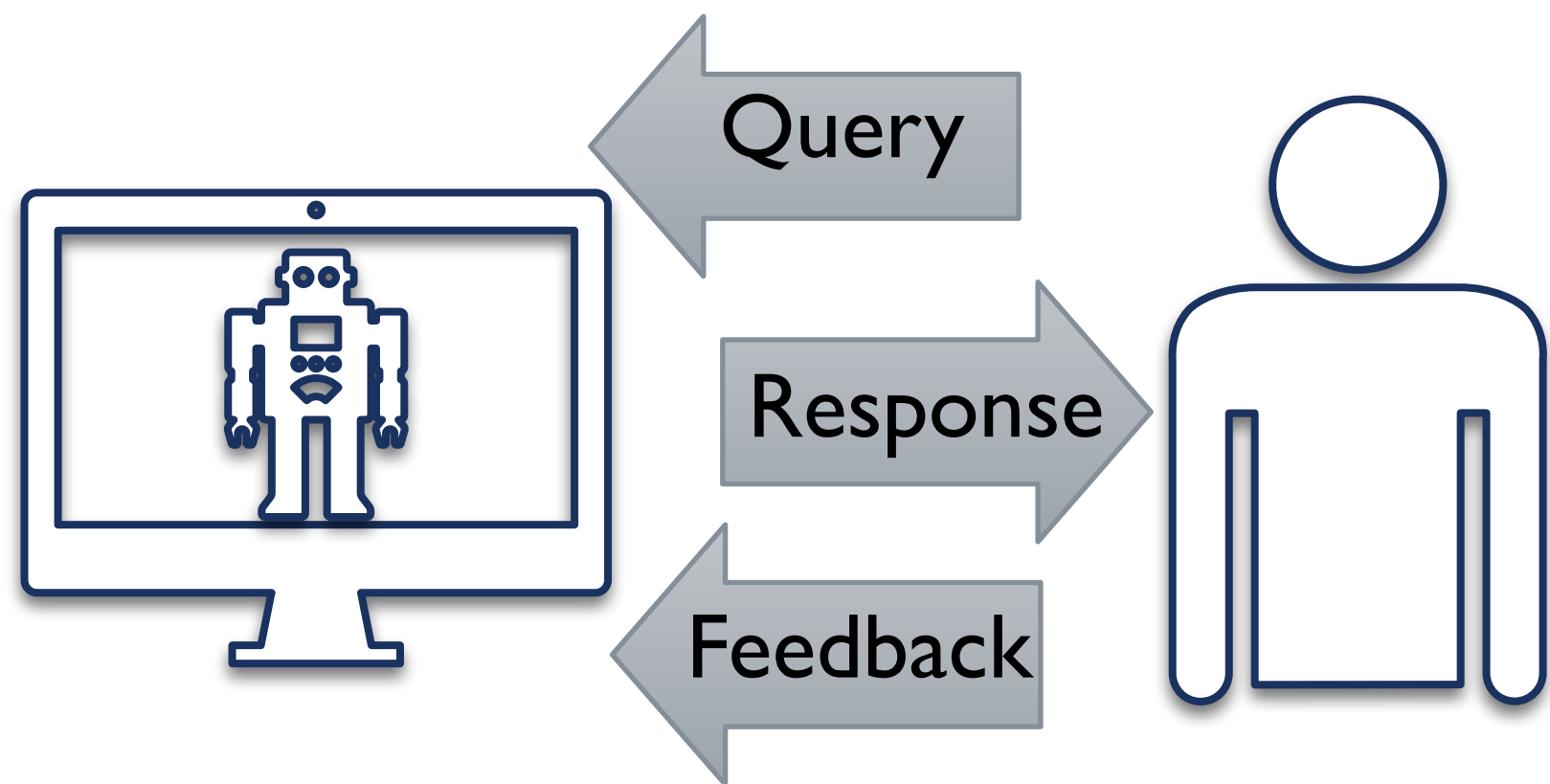
Unsupervised and Reinforcement Learning

USEM

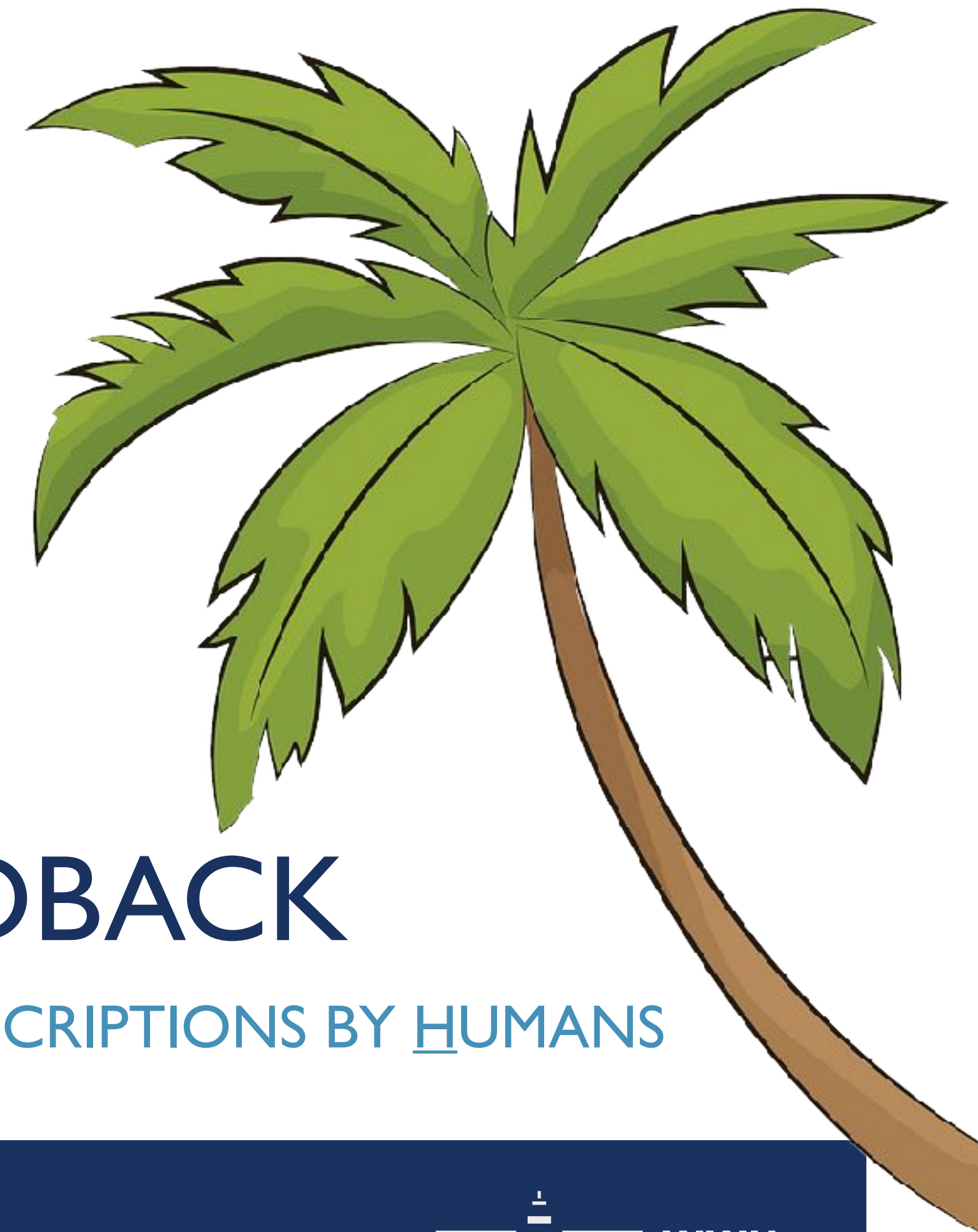


MINIMAL DATA TO TEXT UNDERSTANDING





How to incorporate Feedback?



# CONTINUOUS IMPROVEMENT BY FEEDBACK

FRESH – FEEDBACK-RELIANT ENHANCEMENT OF SUBJECTIVE CONTENT DESCRIPTIONS BY HUMANS



UNIVERSITÄT ZU LÜBECK



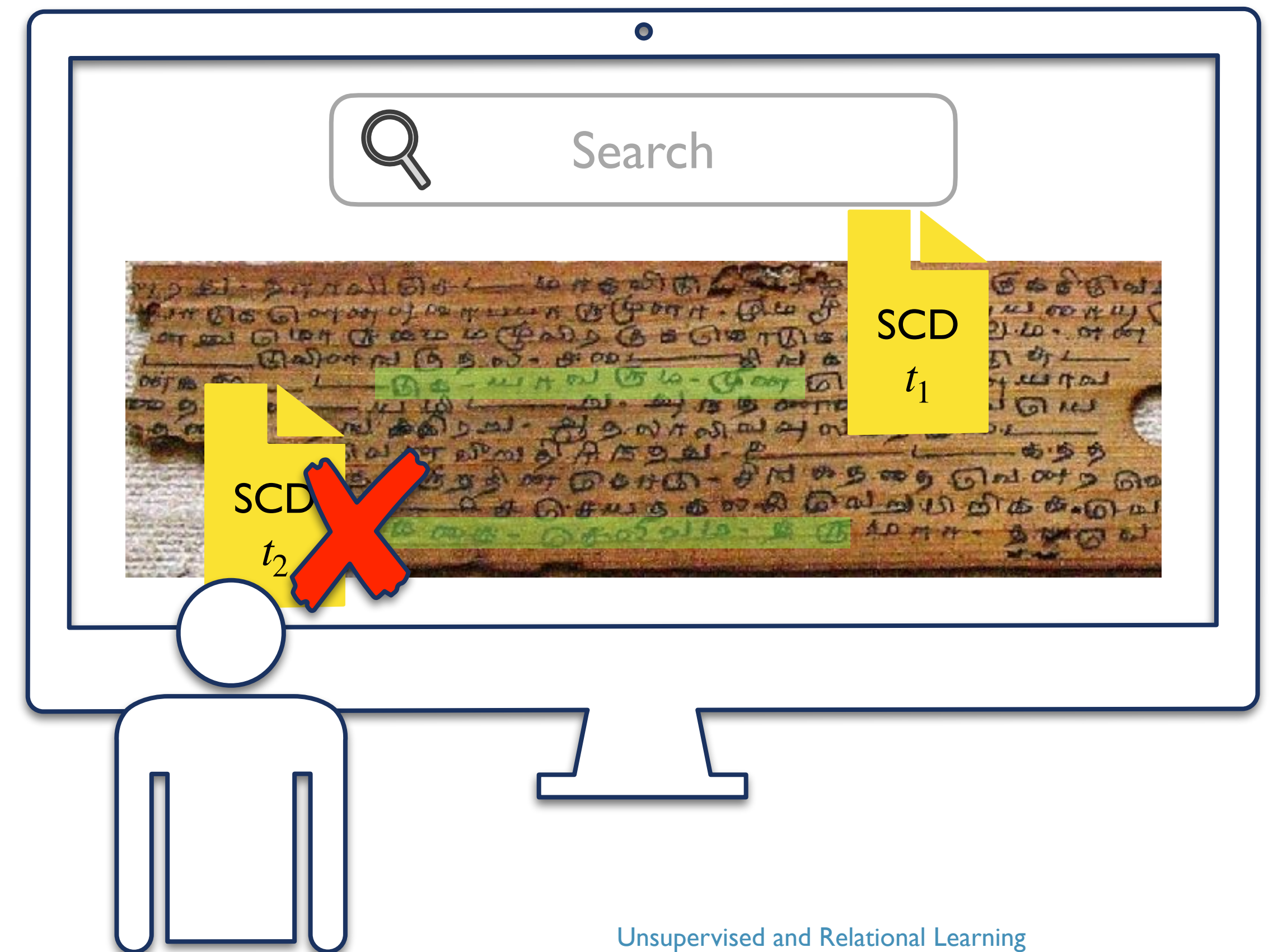
FROM MINIMAL DATA TO TEXT UNDERSTANDING

Unsupervised and Relational Learning

15

# FRESH – FEEDBACK-RELIANT ENHANCEMENT OF SUBJECTIVE CONTENT DESCRIPTIONS BY HUMANS

- Information retrieval service uses SCDs
  - Corpus of documents associated with SCD
  - SCD matrix for corpus
- **Problem:** Faulty SCDs, faulty content like *fake-news*, or privacy-protected content
  - Delete from corpus ✓
  - Retrain SCD matrix from scratch?
  - Update SCD matrix ✓





# UPDATE SCD MATRIX: DELETE SINGLE SENTENCE

- Update distribution (matrix row) of SCD
- Reverse SEM for sentences  $p$  and SCD
- Cases
  - **C1:** Sentence and SCD known
  - **C2:** SCDs not known  
→ MPS<sup>2</sup>CD
  - **C3:** Distribution instead of frequencies in matrix  
→ Assume factor
- C2+C3 may be combined

---

## Algorithm Feedback-reliant Enhancement of SCDs

---

```
1: function FRESH(SCD Matrix  $\delta(\mathcal{D})$ , Set of faulty Sentences  $p$ )
2:   for each  $(s, t) \in p$  do
3:     if  $t = nil$  then
4:        $t = \text{MPS}^2\text{CD}(\delta(\mathcal{D}), s)$ 
5:     if DISTRIBUTIONMATRIX( $\delta$ ) then
6:        $m = \min_{j=1, \dots, n; \delta(\mathcal{D})[t][j] > 0} \delta(\mathcal{D})[t][j]$ 
7:     else
8:        $m = 1$ 
9:     for each word  $w_i \in s$  do
10:       $\delta'(\mathcal{D})[t][w_i] \ominus I(w_i, s) \cdot m$ 
11:   return  $\delta(\mathcal{D})$ 
```

---



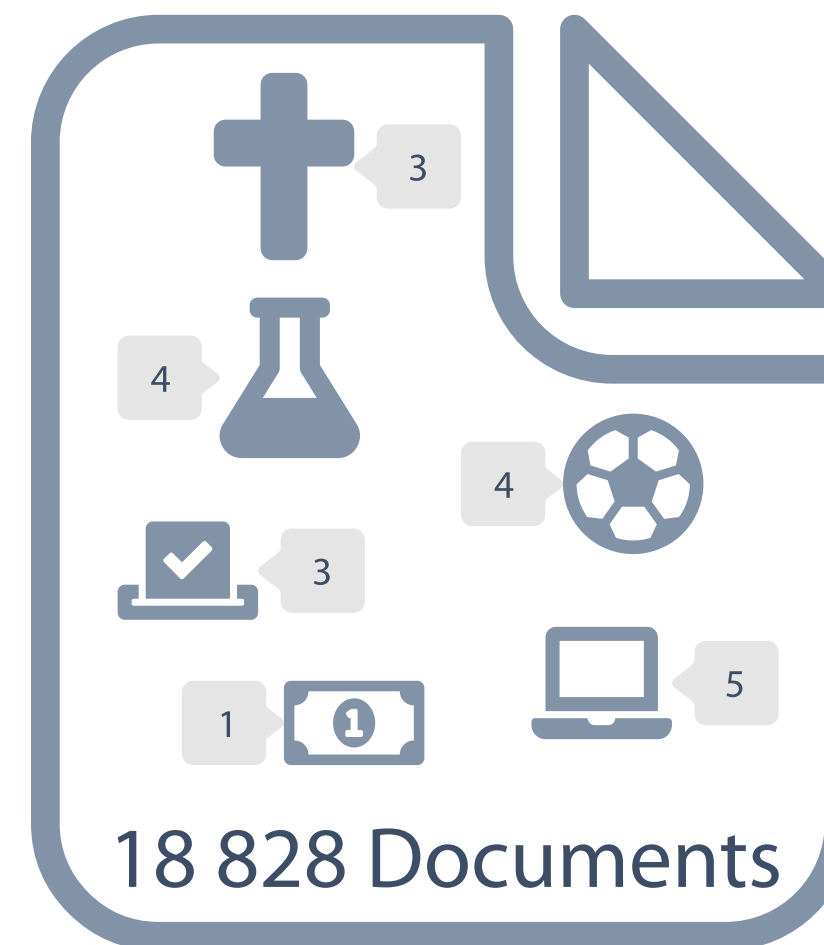
# EVALUATION

## ■ Corpora

- Assumed faulty  $\mathcal{D}_s$
- Assumed correct  $\mathcal{D}_k$
- Full corpus  $\mathcal{D}_f = \mathcal{D}_s \cup \mathcal{D}_k$

## ■ Workflow

1. SCD matrices  $\delta(\mathcal{D}_f)$  and  $\delta(\mathcal{D}_k)$
2. Run update  $\delta' = \text{FrESH}(\delta(\mathcal{D}_f), \mathcal{D}_s)$
3. Evaluate distance between  $\delta(\mathcal{D}_k)$  and  $\delta'$

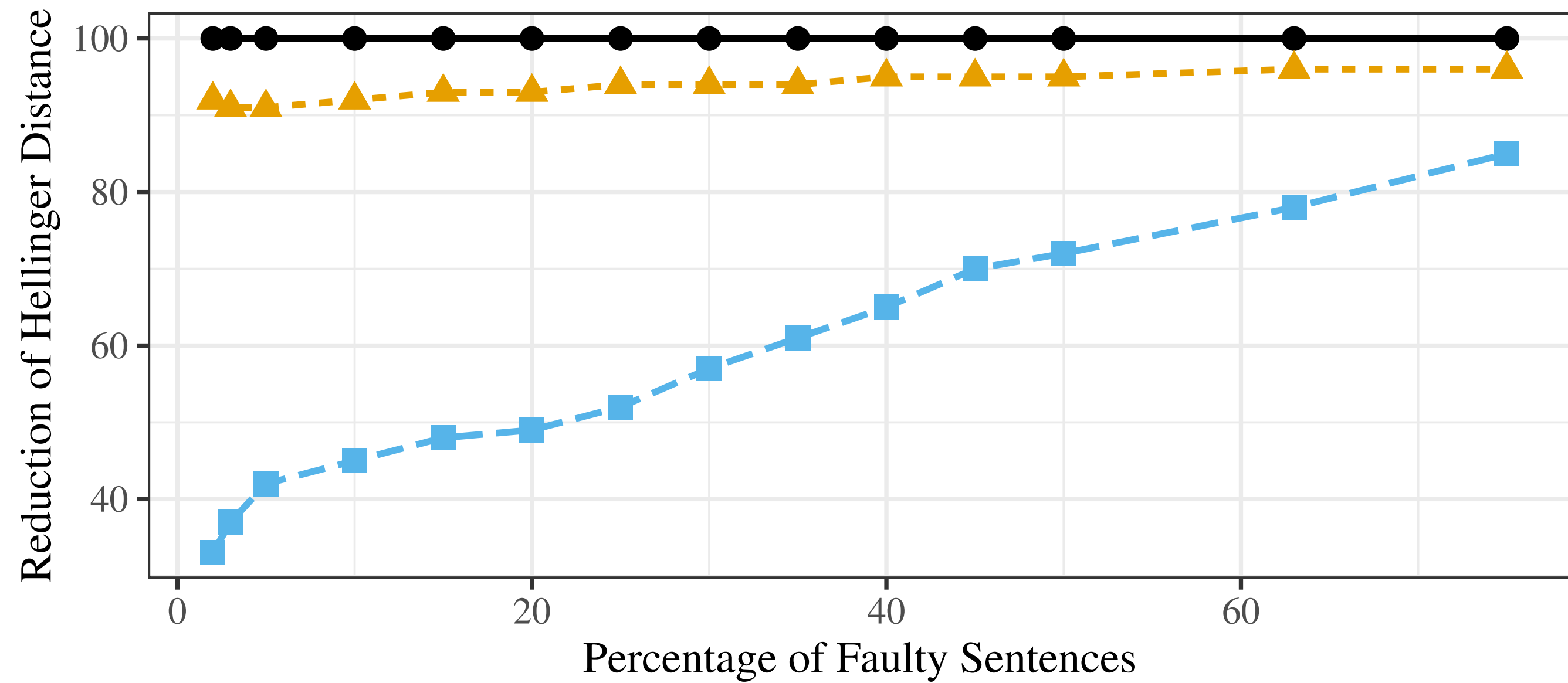
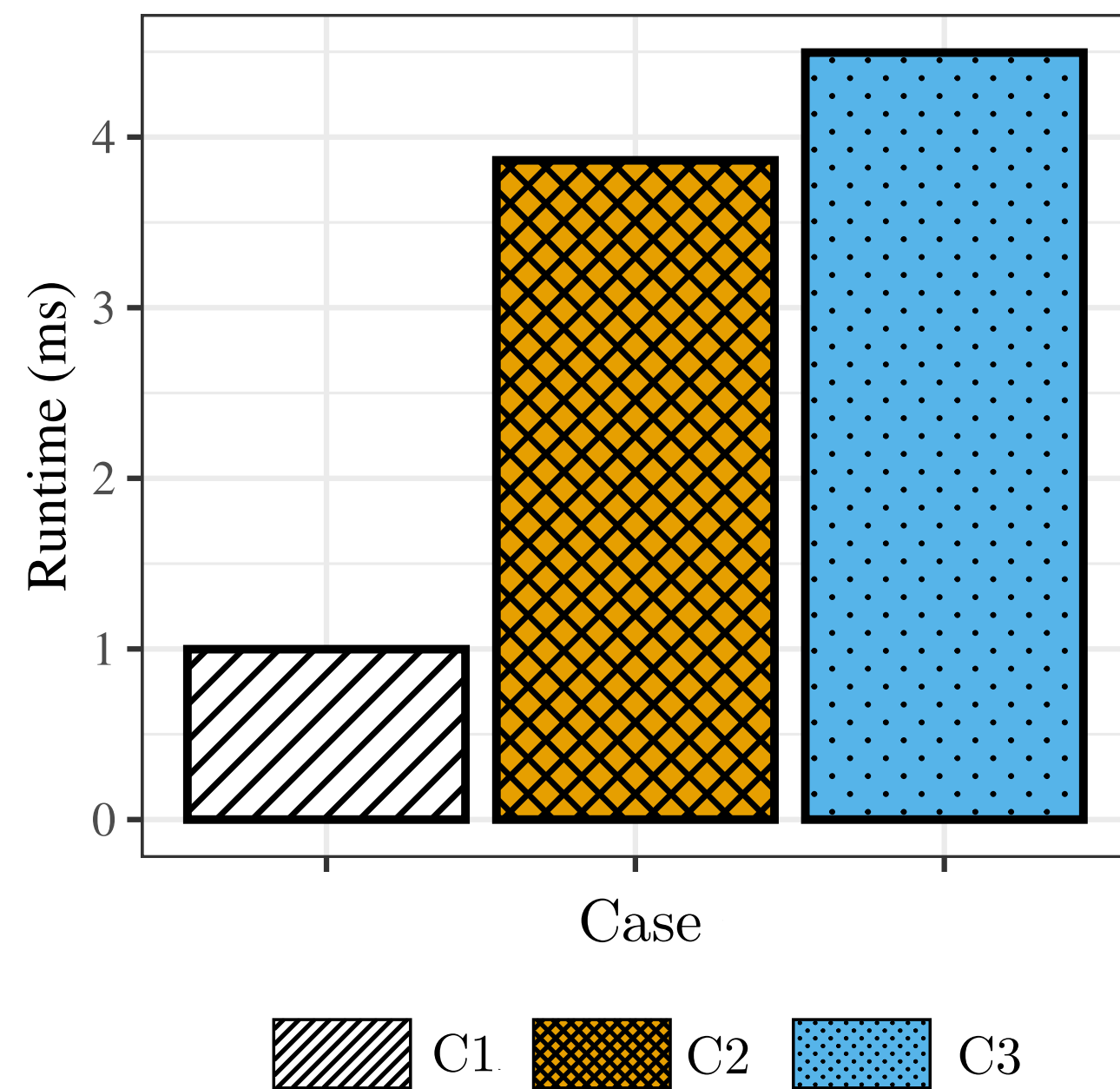


## ■ Dataset

- 20 newsgroups
- SCD using SEM and Open IE

$$HD_t(\delta', \delta(\mathcal{D}_k)) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^n \left( \sqrt{\delta'[t][i]} - \sqrt{\delta(\mathcal{D}_k)[t][i]} \right)^2}$$

# RESULTS: RUNTIME & DELETION ACCURACY



Case ● C1 ▲ C2 ■ C3

# USAGE EXAMPLE

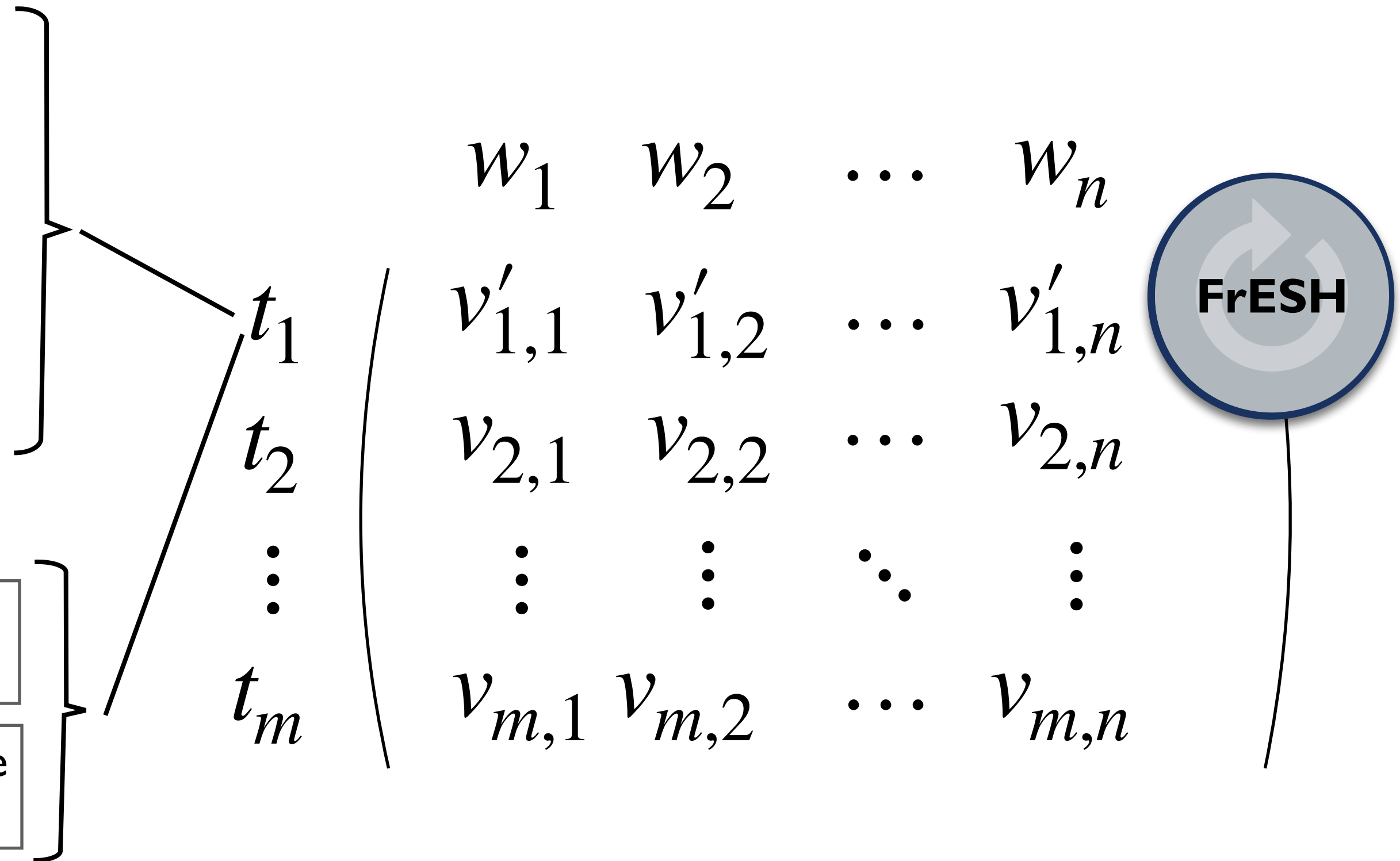
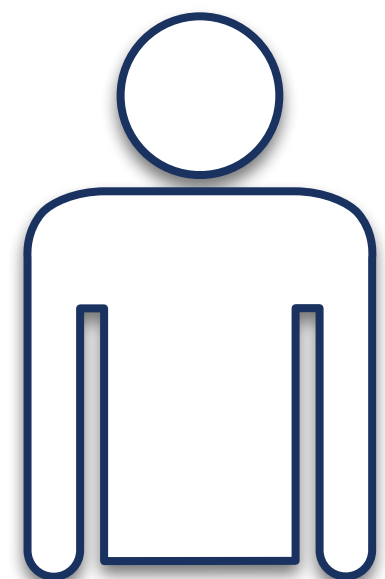
„The seat of an association, unless otherwise provided, is the place where the administration is conducted.“

„The seat of a foundation, unless otherwise provided, is the place where the administration is conducted.“

„A minor child shares the parents' residence.“

„The seat of an association, unless otherwise provided, is the place where the administration is conducted.“

„The seat of a foundation, unless otherwise provided, is the place where the administration is conducted.“







Referenced Sentences

$\{s_1, \dots, s_S\}$

**SCD**  $t_i$

Additional Data  $\mathcal{C}_i$ :

- Label  $l_i$
- Relations
- Links
- ...

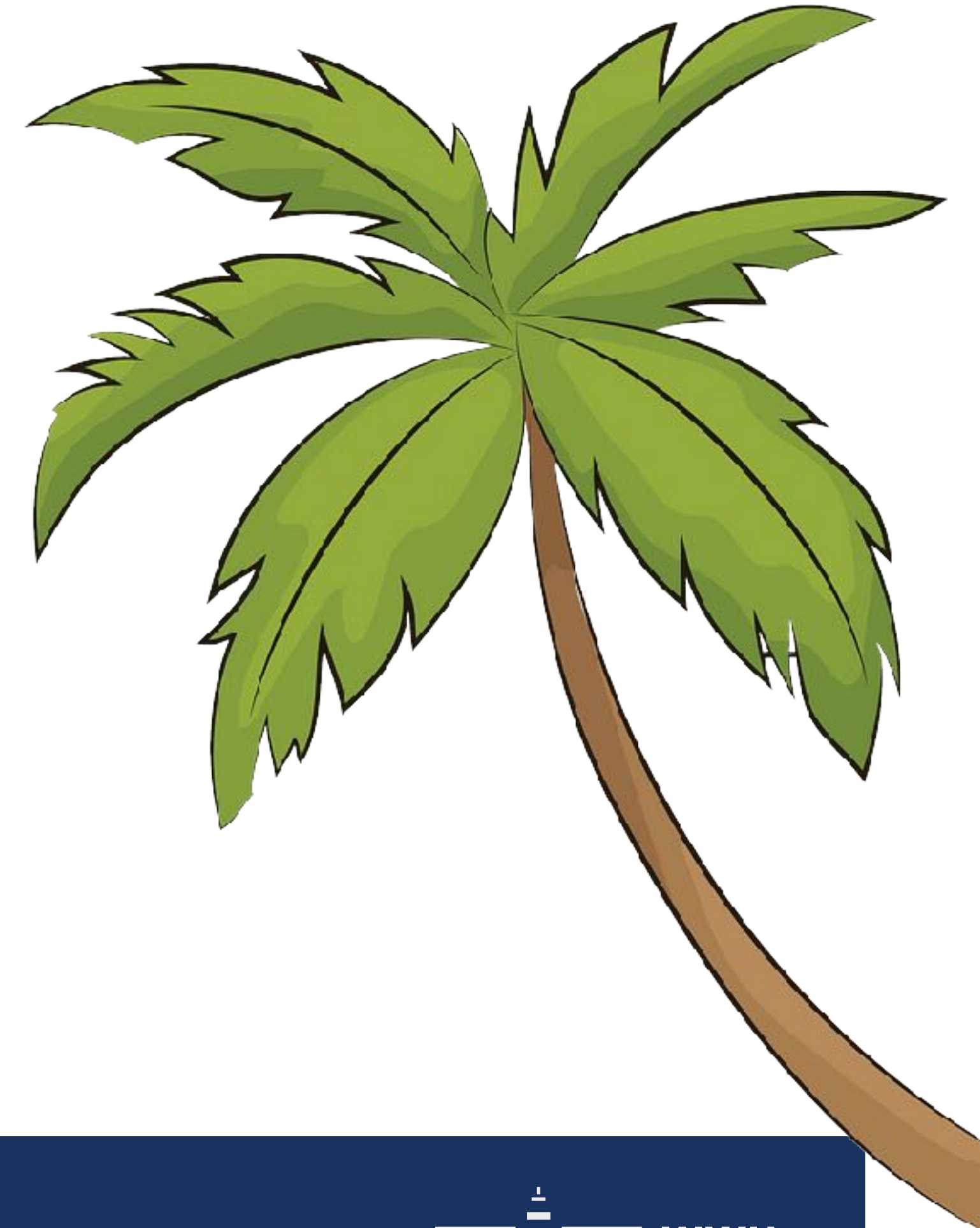


Labels?

$$\begin{matrix}
 & w_1 & w_2 & \dots & w_n \\
 t_1 & v_{1,1} & v_{1,2} & \dots & v_{1,n} \\
 t_2 & v_{2,1} & v_{2,2} & \dots & v_{2,n} \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 t_m & v_{m,1} & v_{m,2} & \dots & v_{m,n}
 \end{matrix}$$

Word Distribution

$\{v_{i,1}, \dots, v_{i,n}\}$



# LABELLING OF SCDS

LESS – LEAN COMPUTING FOR SELECTIVE SUMMARIES

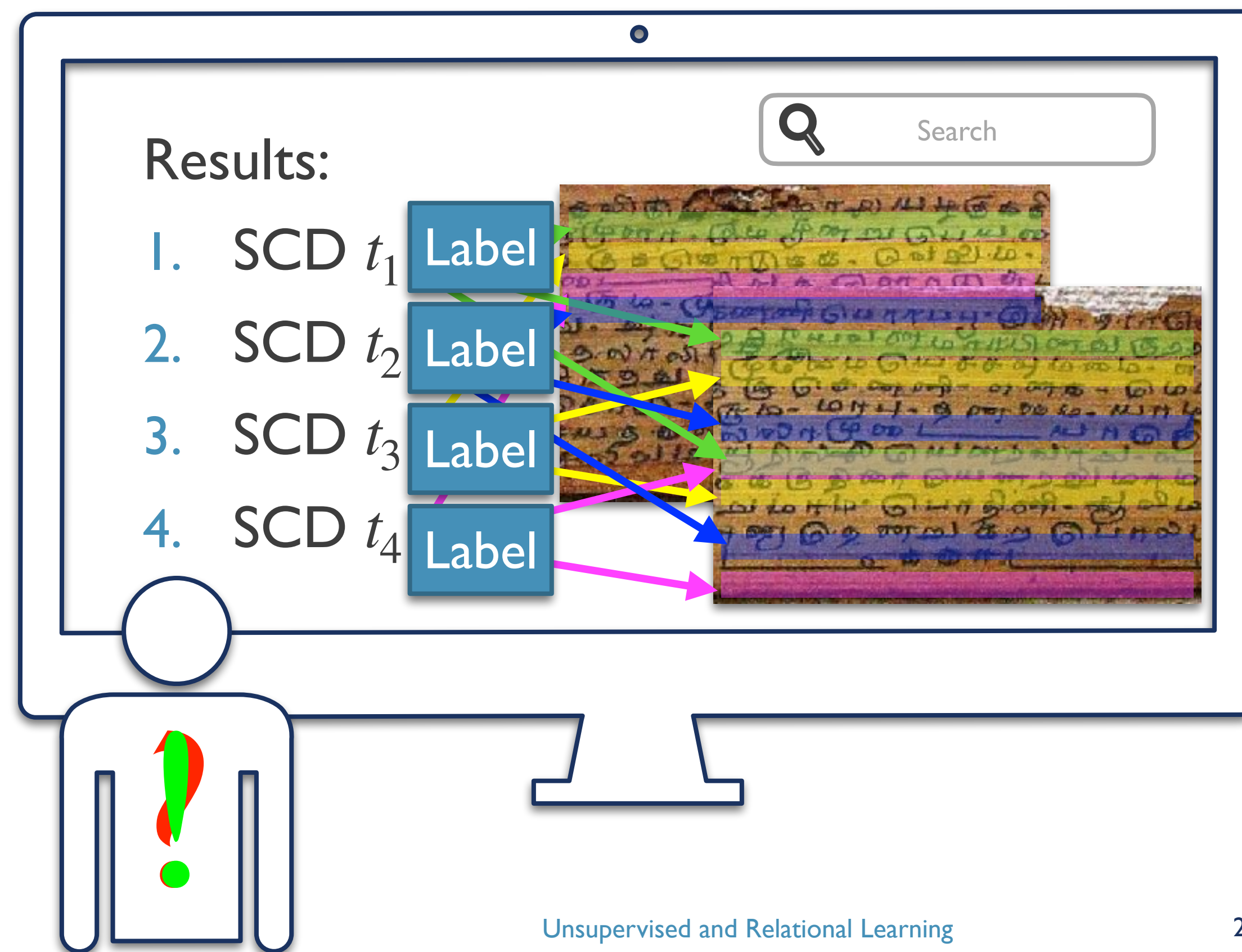


UNIVERSITÄT ZU LÜBECK



# LABELS AS DESCRIPTIONS FOR SCDs

- User browses corpus with SCDs
- SCDs represent concepts mentioned in corpus
- SCDs contain references to sentence
- **Problem:** System needs to describe SCDs to user
- ➔ **Solution:** Label for SCDs



# INFORMATION SOURCES

- Available per SCD

- ✓ References sentences  $\{s_1, \dots, s_S\}$
- ✓ Word distribution  $\{v_{i,1}, \dots, v_{i,n}\}$
- ✗ Label
- ✗ Other data like *relations*

- Formalised problem

$$l_i = \underset{l_j \in \text{all possible labels}}{\operatorname{argmax}} \operatorname{Utility}\left(l_j, t_i = \left((v_{i,1}, \dots, v_{i,n}), \{s_1, \dots, s_S\}\right)\right)$$

SCD  $t_i$

Additional Data  $\mathcal{C}_i$ :

- Label  $l_i$
- Relations
- Links
- ...

Referenced Sentences

$\{s_1, \dots, s_S\}$



Word Distribution  $\{v_{i,1}, \dots, v_{i,n}\}$

	$w_1$	$w_2$	$\dots$	$w_n$
$t_1$	$v_{1,1}$	$v_{1,2}$	$\dots$	$v_{1,n}$
$t_2$	$v_{2,1}$	$v_{2,2}$	$\dots$	$v_{2,n}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$t_m$	$v_{m,1}$	$v_{m,2}$	$\dots$	$v_{m,n}$



# LABEL CANDIDATES & UTILITY OF LABELS

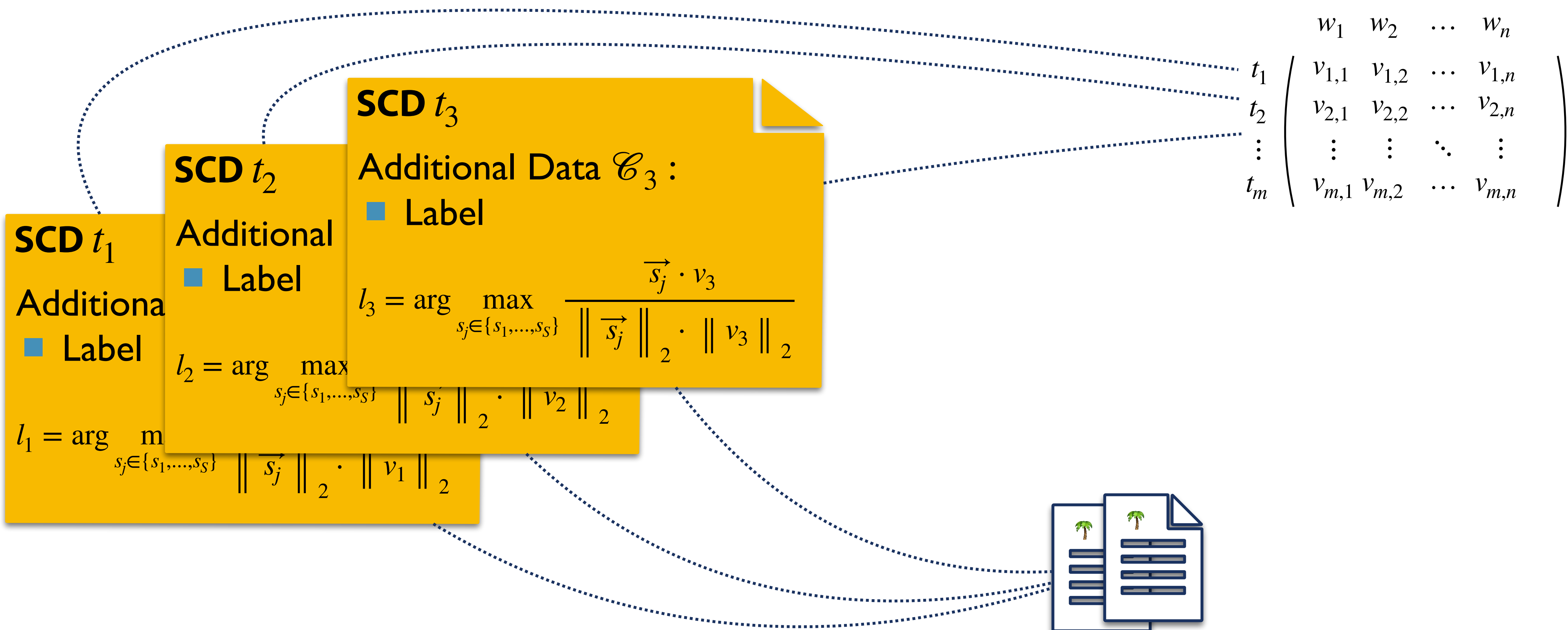
- What should a label look like?
  - Sequence of words like a short description
  - Summary of SCDs
- **Candidates:**
  - Referenced sentences  $\{s_1, \dots, s_S\}$
  - Reformulate problem
- What is a *good* label?
  - Similar to references sentences of SCDs
  - Word distributions generates sentences
- **Utility:** Cosine similarity
  - Use word distribution  $\{v_{i,1}, \dots, v_{i,n}\}$
  - Reformulate problem

$$l_i = \arg \max_{s_j \in \{s_1, \dots, s_S\}} Utility(s_j, (v_{i,1}, \dots, v_{i,n}))$$

$$l_i = \arg \max_{s_j \in \{s_1, \dots, s_S\}} \frac{\vec{s}_j \cdot v_i}{\|\vec{s}_j\|_2 \cdot \|v_i\|_2}$$



# APPROACH – LESS



## ALGORITHM

**Algorithm** LEan computing for Selective Summaries

---

```

1: function LESS( $\mathcal{D}$ )
2:   Input: Corpus  $\mathcal{D}$ 
3:   Output: SCD matrix  $\delta(\mathcal{D})$ ; SCD set  $g(\mathcal{D})$  containing labels  $l_i$  for SCDs  $t_i$ 
4:    $\delta(\mathcal{D}) \leftarrow \text{USEM}(\mathcal{D})$  ▷ Run USEM
5:    $g(\mathcal{D}) \leftarrow \{\}$  ▷ Initialize empty SCD set  $g(\mathcal{D})$ 
6:   for each row of matrix  $i = 1, \dots, K$  do
7:      $v_i \leftarrow g(\mathcal{D})[i]$  ▷ Extract SCD-word distribution
8:      $\{s_1, \dots, s_S\} \leftarrow \text{REFERENCEDSENTENCES}(i)$  ▷ Get referenced sentences
9:      $l_i \leftarrow \arg \max_{s_j \in \{s_1, \dots, s_S\}} \frac{\vec{s}_j \cdot v_i}{\|\vec{s}_j\|_2 \cdot \|v_i\|_2}$  ▷ Compute label
10:     $t_i \leftarrow (\{l_i\}, \{s_1, \dots, s_S\})$  ▷ Compose associated SCD with computed label
11:     $g(\mathcal{D}) \cup \{t_i\}$  ▷ Add to SCD set
12:   return  $\delta(\mathcal{D}), g(\mathcal{D})$ 

```

---

# EVALUATION & DATASET: AGAIN BGB

- Bürgerliches Gesetzbuch (BGB)
  - German civil code (German language)
- First run USEM
- Second add labels with LESS
  - Compare to BERT-based approach

Bürgerliches  
Gesetzbuch

BGB

**Buch 1**

*Allgemeiner Teil*  
§§ 1 - 240

2022

BGB

**Buch 2**

*Recht der Schuldverhältnisse*  
§§ 241 - 853

2022

Bürgerliches  
Gesetzbuch

BGB

Bürgerliches  
Gesetzbuch

BGB

Bürgerliches  
Gesetzbuch

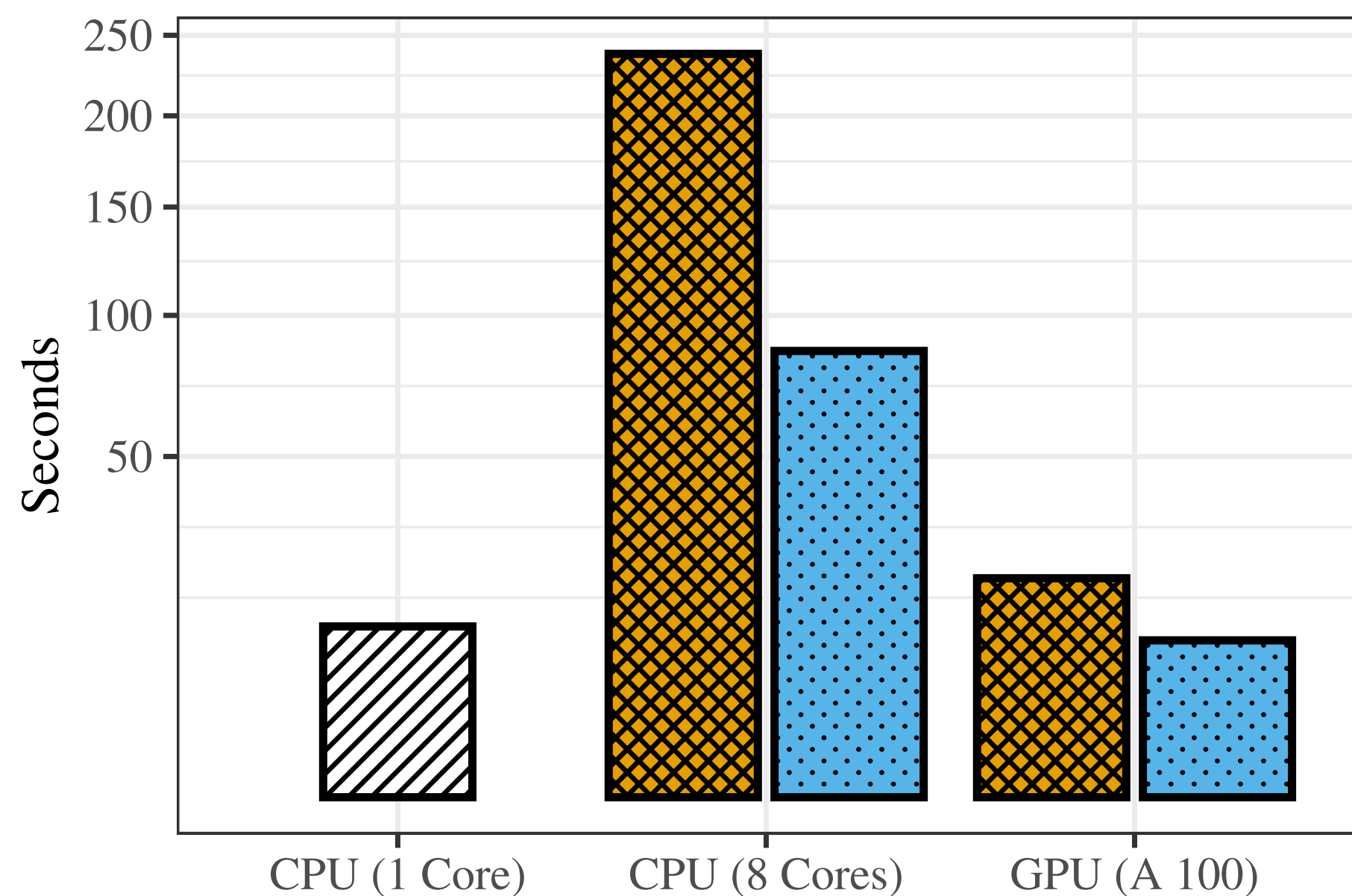
BGB

**Buch 5**

*Erbrecht*  
§§ 1922 - 2385

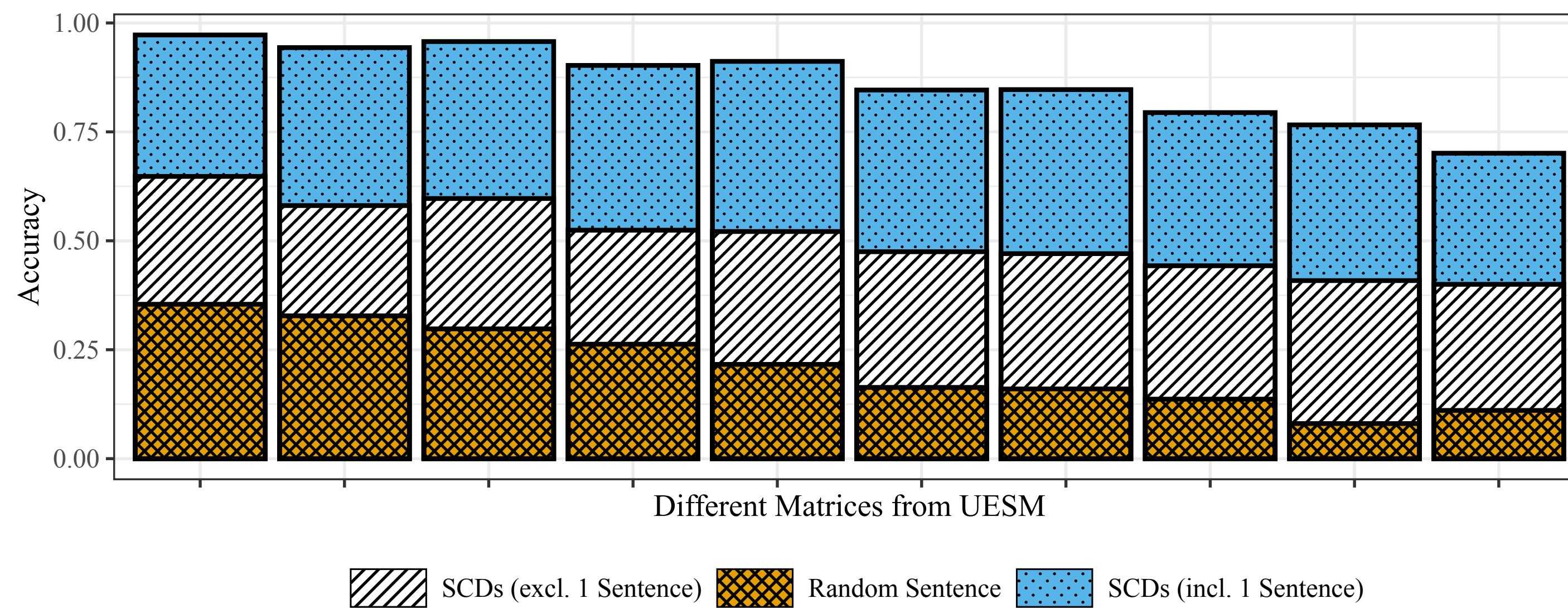
2022

# RUNTIME AND ACCURACY



LESS
  BERT Q&A
  BERT Vectors

Two techniques using BERT; 1 or 8 Intel CPU cores and single NVIDIA A100 40GB GPU



Random sentence: Theoretical accuracy a random approach would result in.



# USAGE EXAMPLE

## Association Commercial Business Operation

„An association whose purpose is not directed towards a commercial business operation acquires legal capacity through entry in the register of associations at the competent local court.“

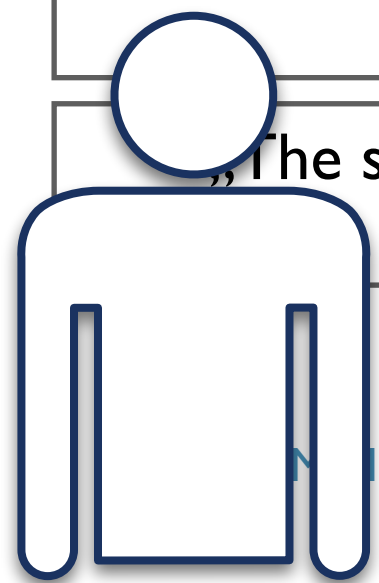
„An association whose purpose is to engage in commercial business shall acquire legal capacity, in the absence of special federal law, through state conferral. The grant is due to the state in whose territory the association has its registered office.“

## Seat Foundation Place Administration

„The seat of an association, unless otherwise provided, is the place where the administration is conducted.“

„The seat of a foundation, unless otherwise provided, is the place where the administration is conducted.“

	$w_1$	$w_2$	$\dots$	$w_n$
$t_1$	$v_{1,1}$	$v_{1,2}$	$\dots$	$v_{1,n}$
$t_2$	$v_{2,1}$	$v_{2,2}$	$\dots$	$v_{2,n}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$t_m$	$v_{m,1}$	$v_{m,2}$	$\dots$	$v_{m,n}$



MINIMAL DATA TO TEXT UNDERSTANDING

Unsupervised and Reinforcement Learning

USEM  
+ LESS



SCD  $t_i$

Additional Data  $\mathcal{C}_i$ :

- Label  $l_i$
- Relations
- Links
- ...

Word Distribution  
 $\{v_{i,1}, \dots, v_{i,n}\}$

	$w_1$	$w_2$	$\dots$	$w_n$
$t_1$	$v_{1,1}$	$v_{1,2}$	$\dots$	$v_{1,n}$
$t_2$	$v_{2,1}$	$v_{2,2}$	$\dots$	$v_{2,n}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$t_m$	$v_{m,1}$	$v_{m,2}$	$\dots$	$v_{m,n}$



Relations?



# INTER- AND INTRA-SCD RELATIONS

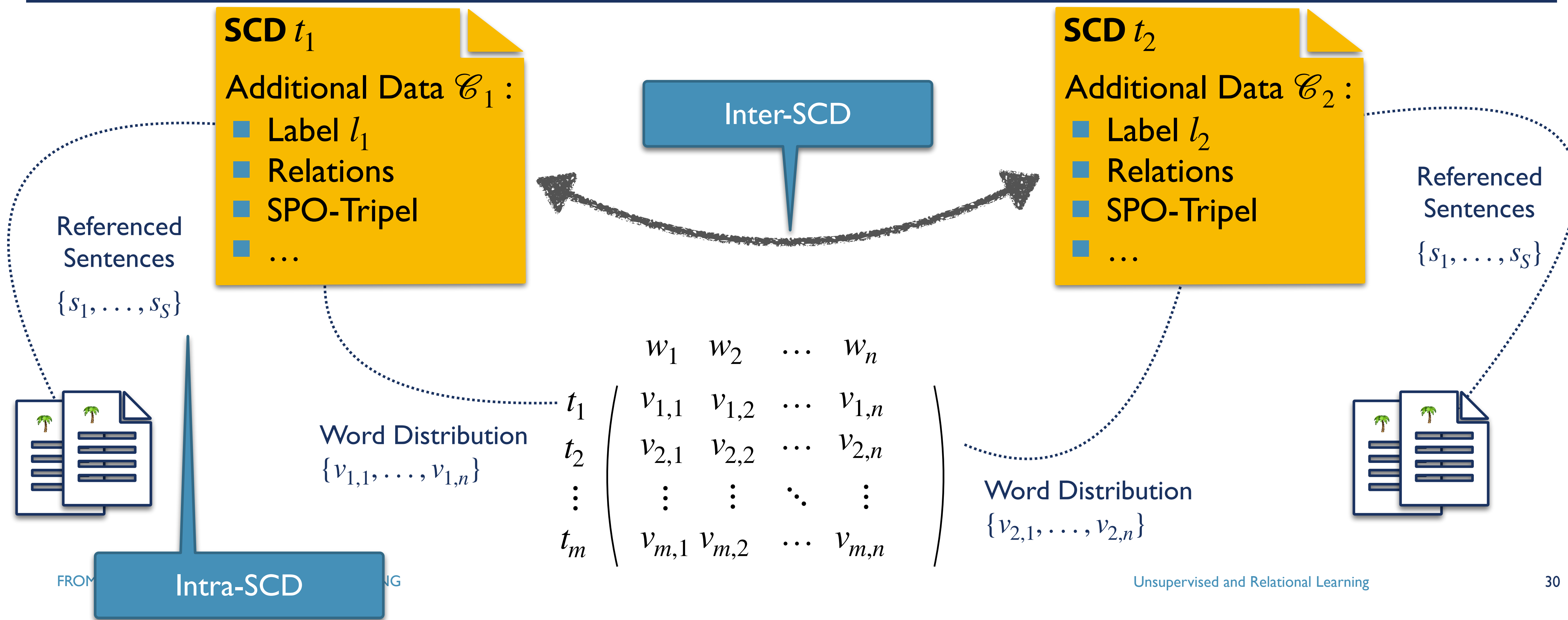
ENRICHING A CORPUS WITH DOCUMENTS USING THE INTER-SCD RELATION COMPLEMENT



UNIVERSITÄT ZU LÜBECK



# RELATIONS AMONG SCDS



# EXAMPLE INTER-SCD RELATION: COMPLEMENT

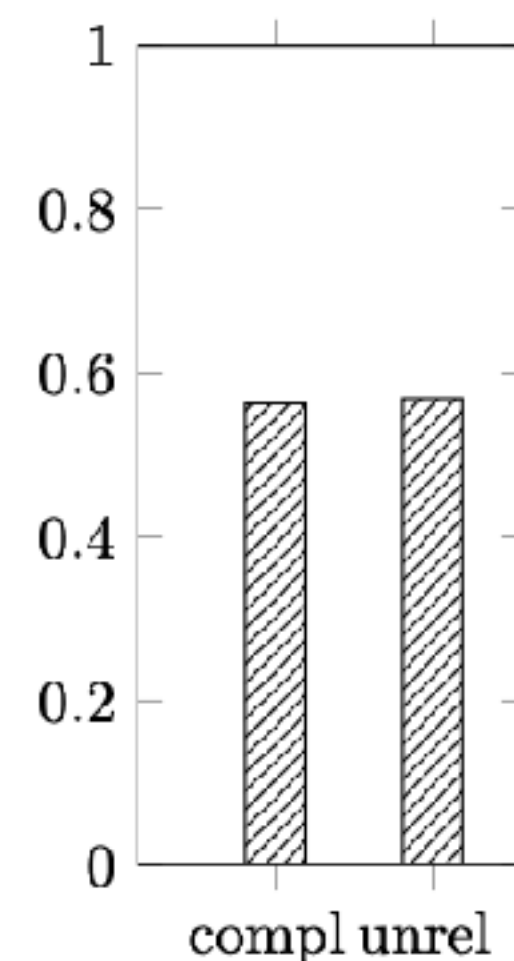
- **Goal:** Identify documents that are complementary to a corpus/ a document in a corpus
  - Binary classification problem:  
*Complement = true or Complement = false*
- Solution approach to corpus enrichment uses cosine **similarity** at its core
  - Sequence of similarity values between vector representations of SCDs and the words in the new document
- Also applies to many document retrieval approaches: return documents similar in some regard
  - Topic distribution similar, entities match (equality), etc.



- Corpus  $\mathcal{D}_r$  on sporting events
  - Olympics 2020, UEFA Euro 2020



- Corpus  $\mathcal{D}_c$  with *complementary* documents
  - Covid-19 spread in cities



Similarity values of complements and unrelated documents for corpus enrichment.

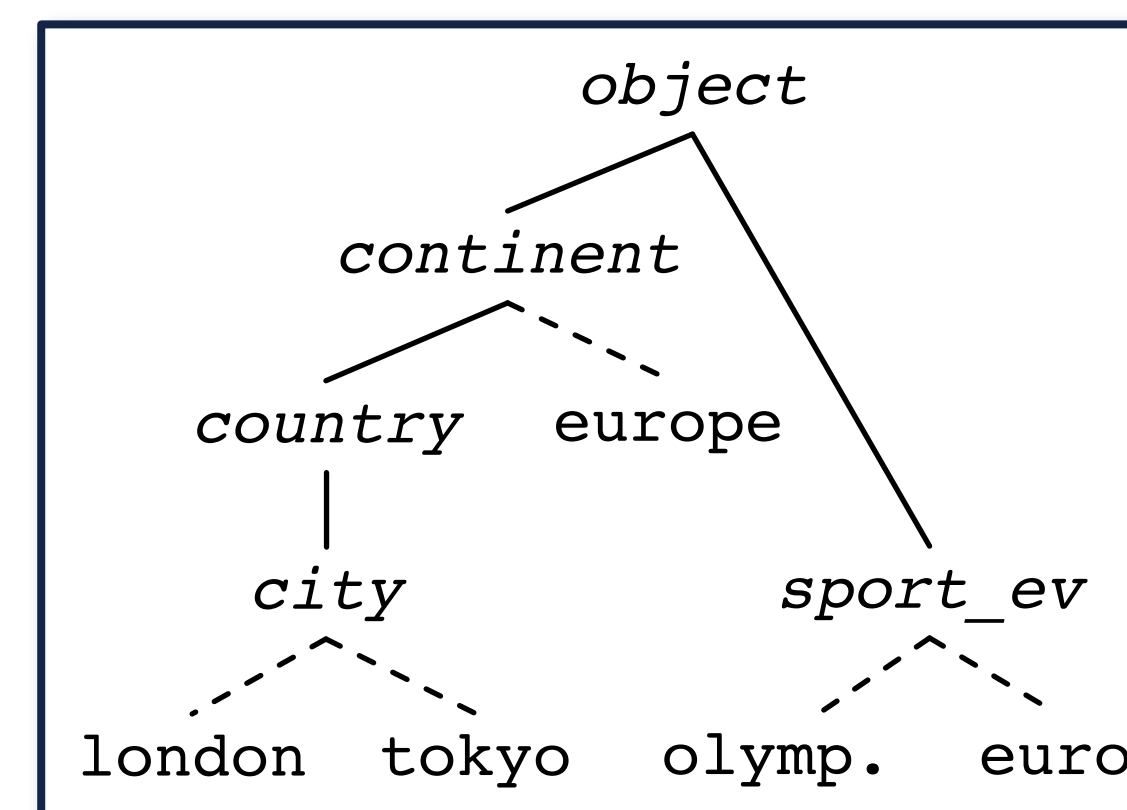
**Problem:** How do we *formally* define complementarity accounting for semantics?

**Problem:** Similarity-based approaches might only provide more of the same.



# HOW TO GRASP COMPLEMENTARITY ON A FORMAL LEVEL?

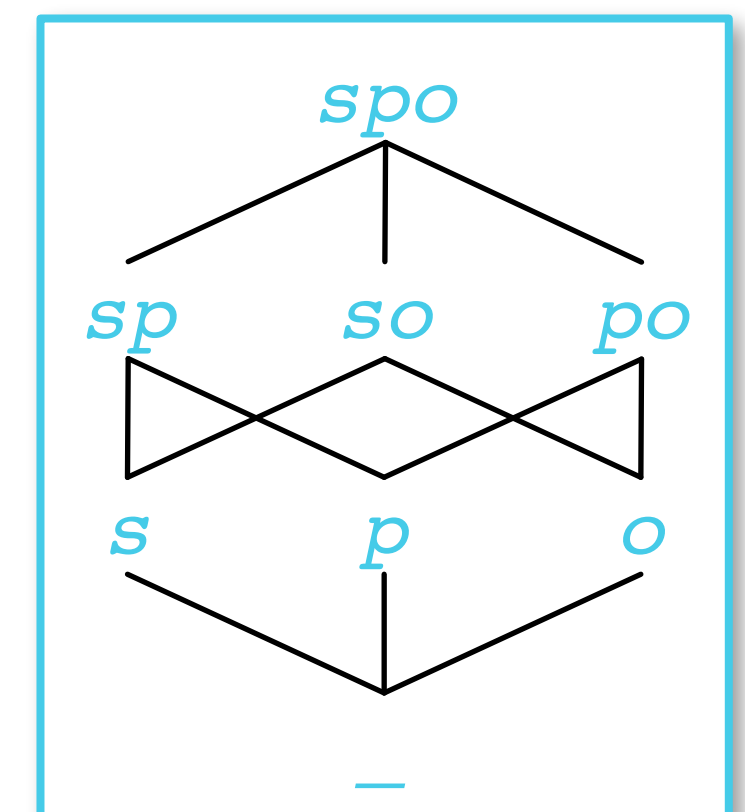
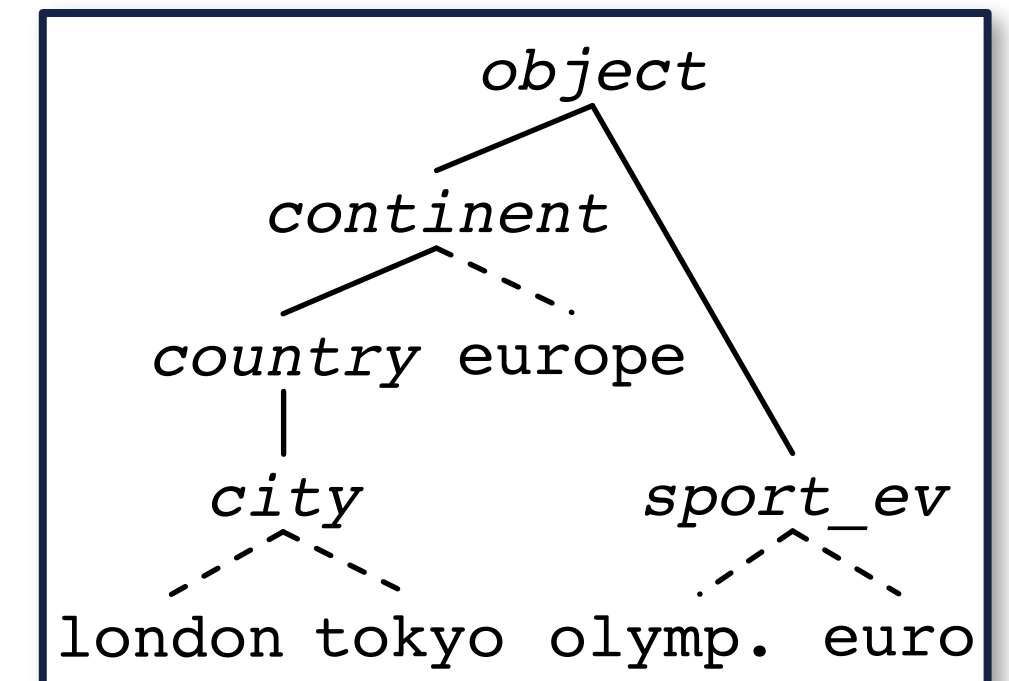
- Use SCDs specifically in the form of **SPO-Triples**
  - SPO-Triple: <subject, predicate, object>
  - Extract for any document, e.g., with OpenIE tools
- Words of *complement* very different compared to corpus
  - Different (topic / SCD) distributions
  - Likely to be classified as unrelated
- Together with a taxonomy
  - Hierarchy of concepts
  - Dictionary of synonyms
- Allow for grasping complementarity on a semantic level by
  - Looking at shared concepts in the SPO triples
  - While also accounting for hierarchy and synonyms
- $t_1$ : <Olympics '21, in, Tokyo>
- $t_2$ : <UEFA euro '20, in, Europe>
- $t_3$ : <Covid-19, in, Tokyo>
- $t_4$ : <Covid-19, in, London>



# A FORMAL DEFINITION: COMPLEMENTARY SCDS

- Let  $x^\uparrow$  refer to the concept or meaning of  $x$
- Seven types of complementarity between SCDSs  $t_i, t_j$ 
  1.  $s$       $t_i = \langle s^\uparrow, p_i, o_i \rangle, t_j = \langle s^\uparrow, p_j, o_j \rangle$
  2.  $p$       $t_i = \langle s_i, p^\uparrow, o_i \rangle, t_j = \langle s_j, p^\uparrow, o_j \rangle$
  3.  $o$       $t_i = \langle s_i, p_i, o^\uparrow \rangle, t_j = \langle s_j, p_j, o^\uparrow \rangle$
  4.  $sp$      $t_i = \langle s^\uparrow, p^\uparrow, o_i \rangle, t_j = \langle s^\uparrow, p^\uparrow, o_j \rangle$
  5.  $so$      $t_i = \langle s^\uparrow, p_i, o^\uparrow \rangle, t_j = \langle s^\uparrow, p_j, o^\uparrow \rangle$
  6.  $po$      $t_i = \langle s_i, p^\uparrow, o^\uparrow \rangle, t_j = \langle s_j, p^\uparrow, o^\uparrow \rangle$
  7.  $spo$     $t_i = \langle s^\uparrow, p^\uparrow, o^\uparrow \rangle, t_j = \langle s^\uparrow, p^\uparrow, o^\uparrow \rangle$
- Types gets more strict  $\rightarrow$  Order in lattice

- $t_1: \langle \text{Olympics '21, in, Tokyo} \rangle$
- $t_2: \langle \text{UEFA euro '20, in, Europe} \rangle$
- $t_3: \langle \text{Covid-19, in, Tokyo} \rangle$
- $t_4: \langle \text{Covid-19, in, London} \rangle$
- $t_1, t_3$   $o$ -complementary
  - $s_1, s_3$  share *object*;  $p_1 = p_3; o_1 = o_3$
  - And  $p, po$  complementary
  - Same holds for  $t_1, t_4$
- $spo$ -complementary
  - All three items share same concept or are identical
- $s$ -complementary
  - $s$  shares same concept, other two different



# CORPUS ENRICHMENT: COMPLEMENTARY DOCUMENTS

- Let  $\mathfrak{C}_x(t_i, t_j), x \in \mathcal{X} = \{s, p, o, sp, so, op, spo\}$  be an **indicator function**

- Returns 1 if  $t_i, t_j$   $x$ -complementary; otherwise 0

- $\mathfrak{C}_x$  is symmetric, i.e.,  $\mathfrak{C}_x(t_i, t_j) = \mathfrak{C}_x(t_j, t_i)$

- Complementarity value between documents  $d', d$ :

$$c(d', d) = \sum_{t_i \in g(d')} \sum_{t_j \in g(d)} \sum_{x \in \mathcal{X}} w_x \mathfrak{C}_x(t_i, t_j)$$

- Sum over all pairs of SCDs  $t_i \in g(d'), t_j \in g(d)$ , indicating if  $t_i, t_j$  are  $x$ -complementary

- $c$  is symmetric, i.e.,  $c(d', d) = c(d, d')$

- Assign **weights**  $w_x, \sum_{w \in \mathcal{X}} w_x = 1$  to complementarity types  $x$  to encode which complementarity interested in

Corpus  
 $\mathcal{D}$



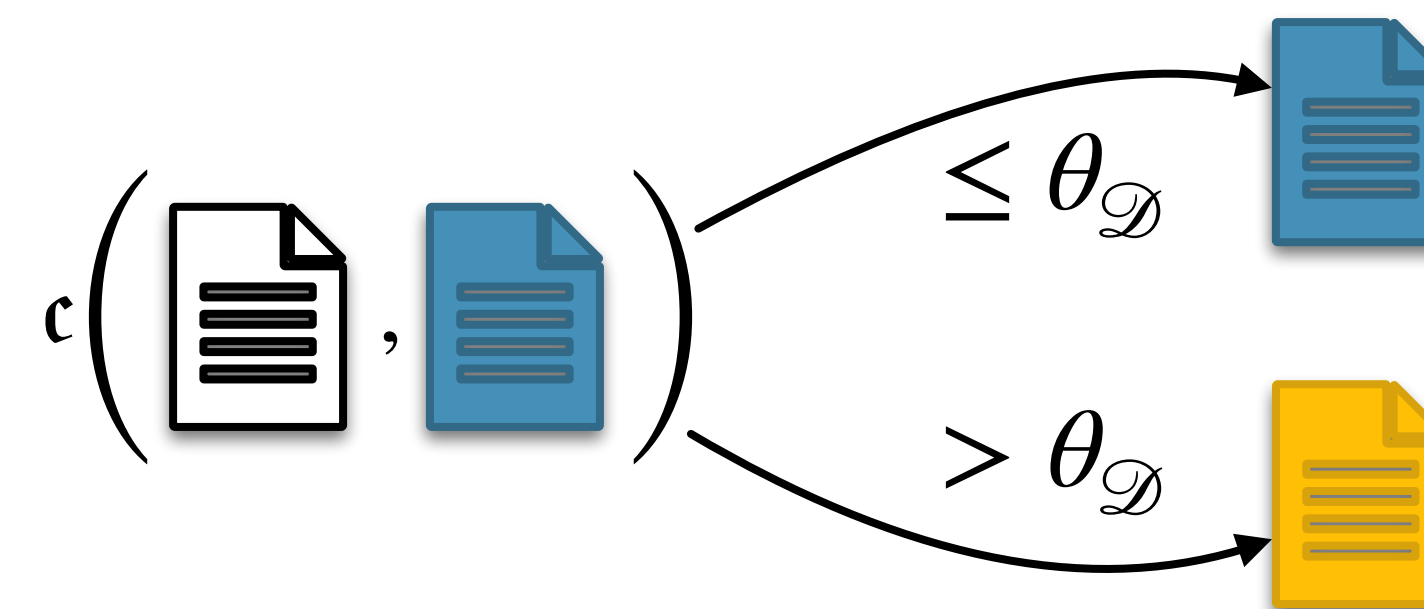
Complementary  
Documents



Related  
Documents

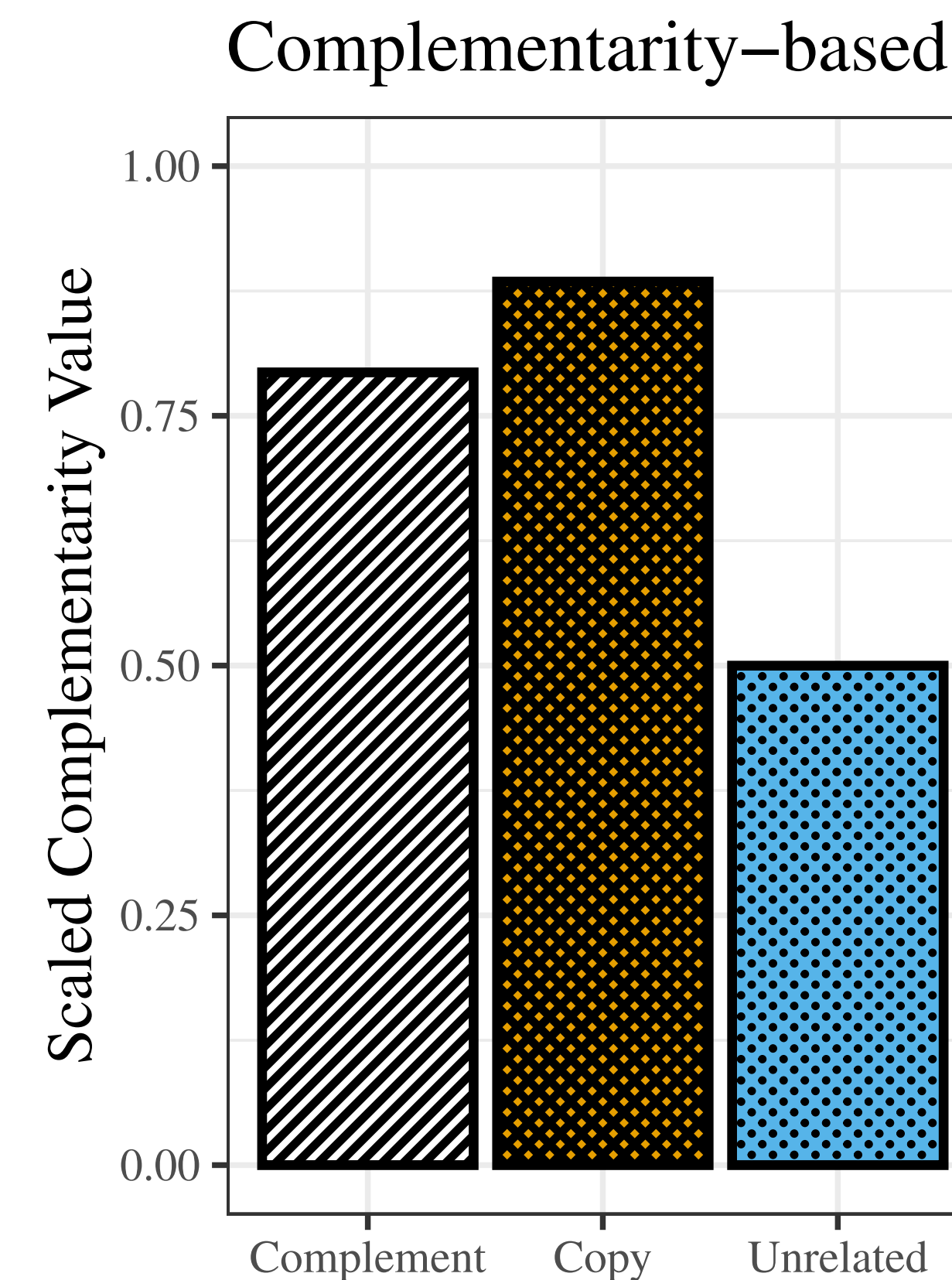
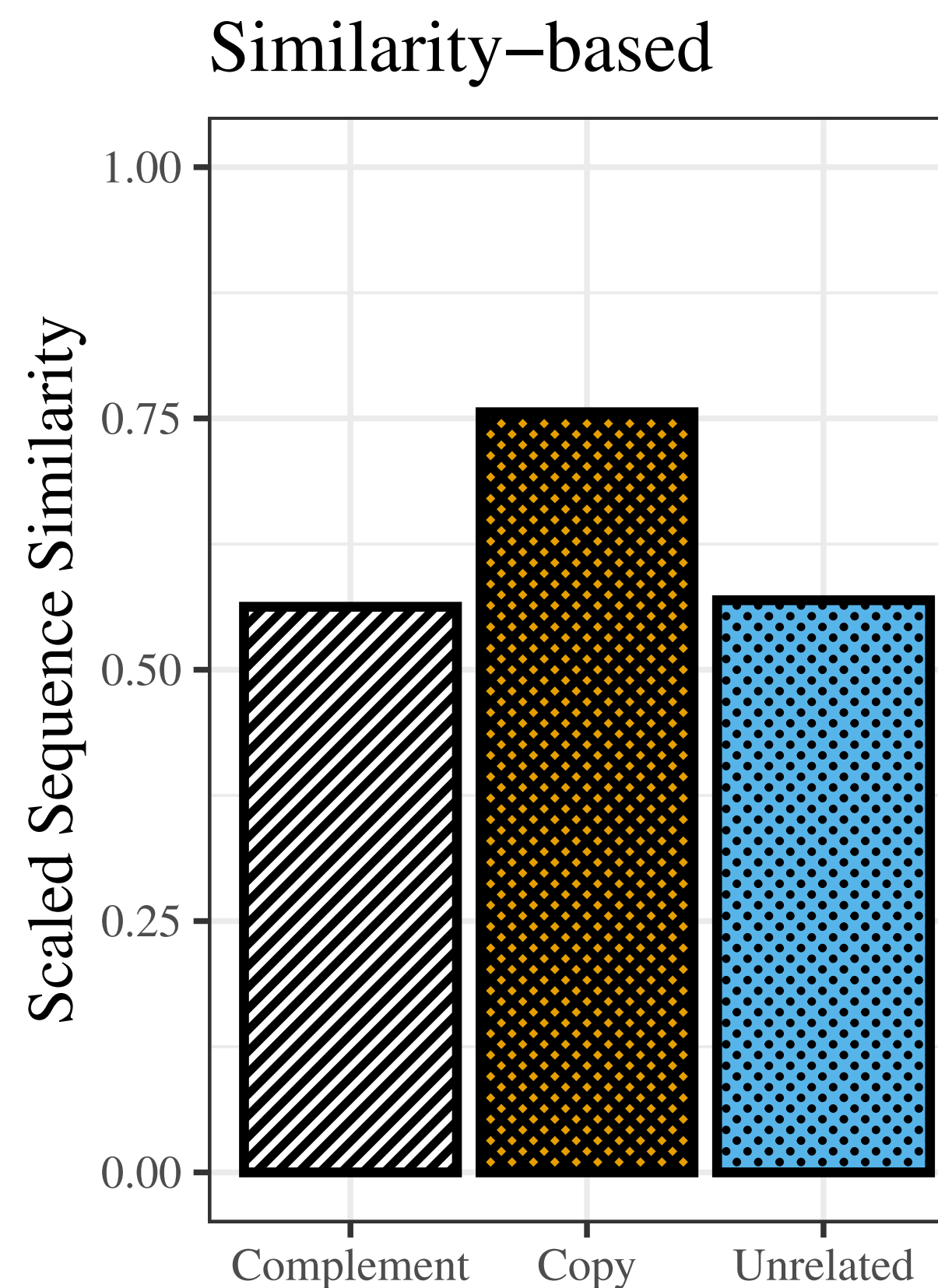


Complement?  
Add to corpus?



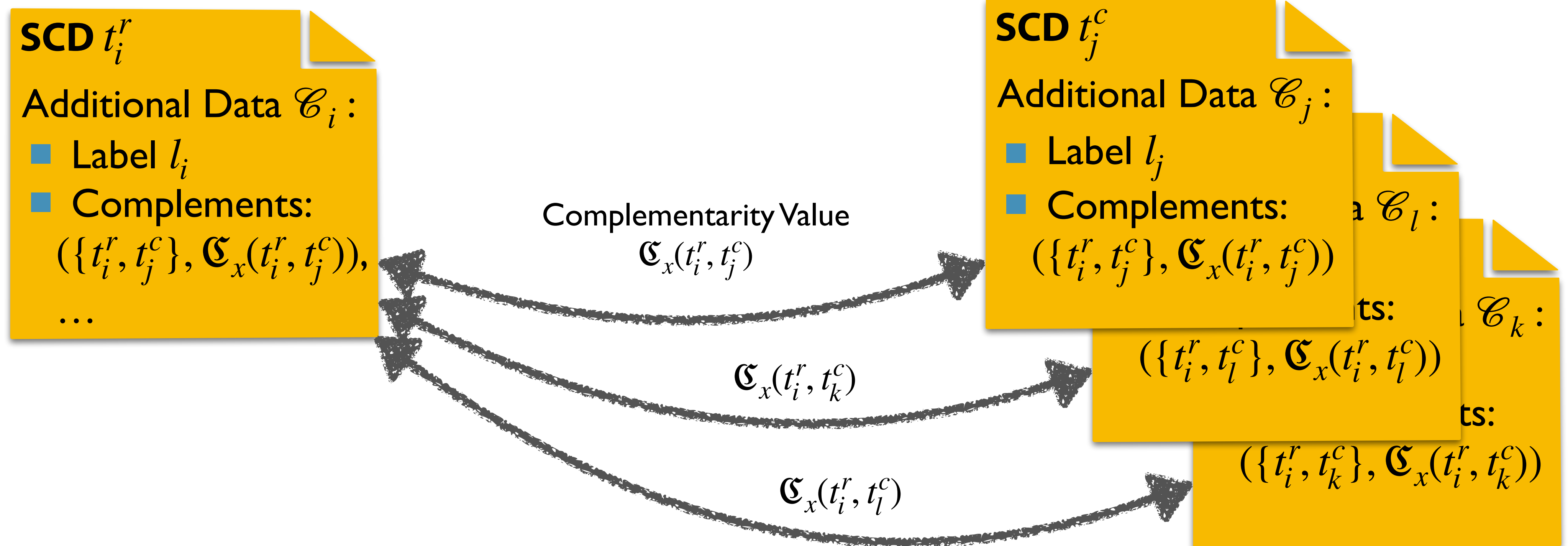
# COMPLEMENT DETECTION

- Similarity-based technique does not distinguish between *complement* and *unrelated*
- Complementarity-based technique uses  $c_x(d', d)$
- ➔ Resulting values differ for *complement* and *unrelated*



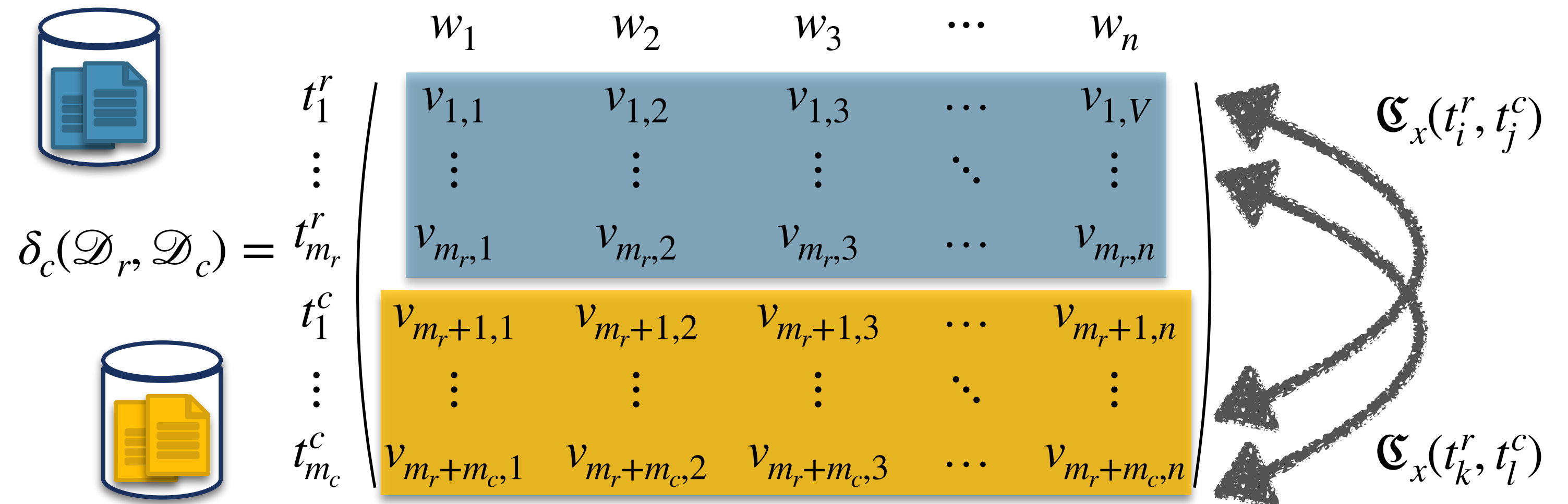


# BACK TO RELATIONS: COMPLEMENTARITY BETWEEN SCDS



# RELATIONS IN SCD MATRIX: COMBINED SCD MATRIX

- Combine SCDs of two different corpora in one SCD matrix
  - Corpus  $\mathcal{D}_r$  related documents
  - Corpus  $\mathcal{D}_c$  with complementary documents
- Model the relations among the SCDs in the matrix
  - Filter matrix to keep only SCDs from  $\mathcal{D}_c$  which are complementary
- Adapted of MPS<sup>2</sup>CD yields negative similarity value for complementary SCDs



May be generalised to any type of corpora and relations among them.

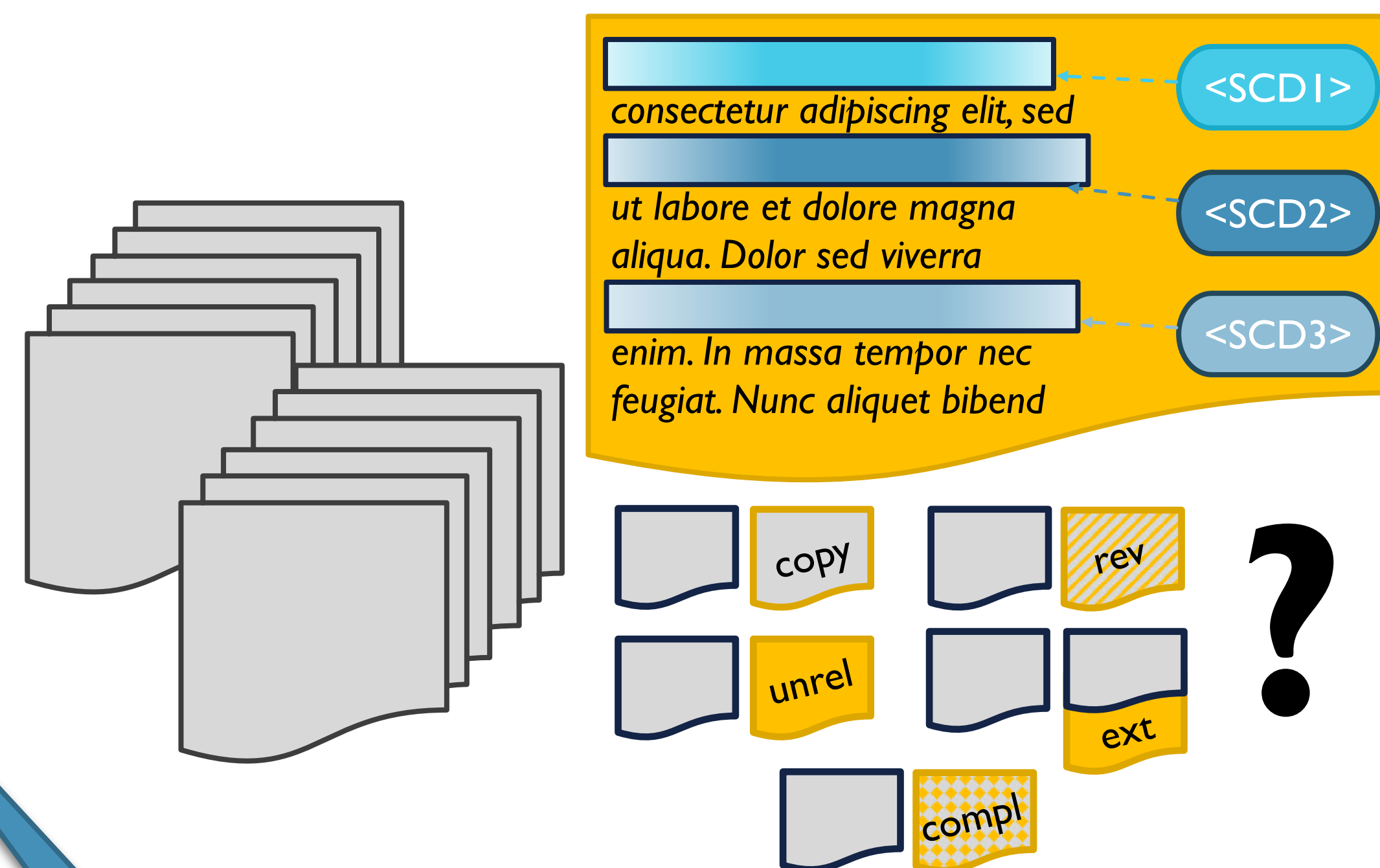
# CORPUS ENRICHMENT INCL. COMPLEMENTS

- Same problem as earlier

*Classify a new document before adding to corpus.*

- Now five types  $\mathcal{Y} = \{copy, ext, rev, unrel, compl\}$
- Find most probable type  $\arg \max_{y \in \mathcal{Y}} P(\text{Type} = y \mid d', \mathcal{D})$

1. Build combined SCD matrix (needs corpus of related and complementary documents, use  $c(d', d)$ )
2. Filter matrix by removing complementary documents with no relation to related document
3. Train an HMM on MPS<sup>2</sup>CD similarity values for classification
4. Run MPS<sup>2</sup>CD on new documents and use HMM

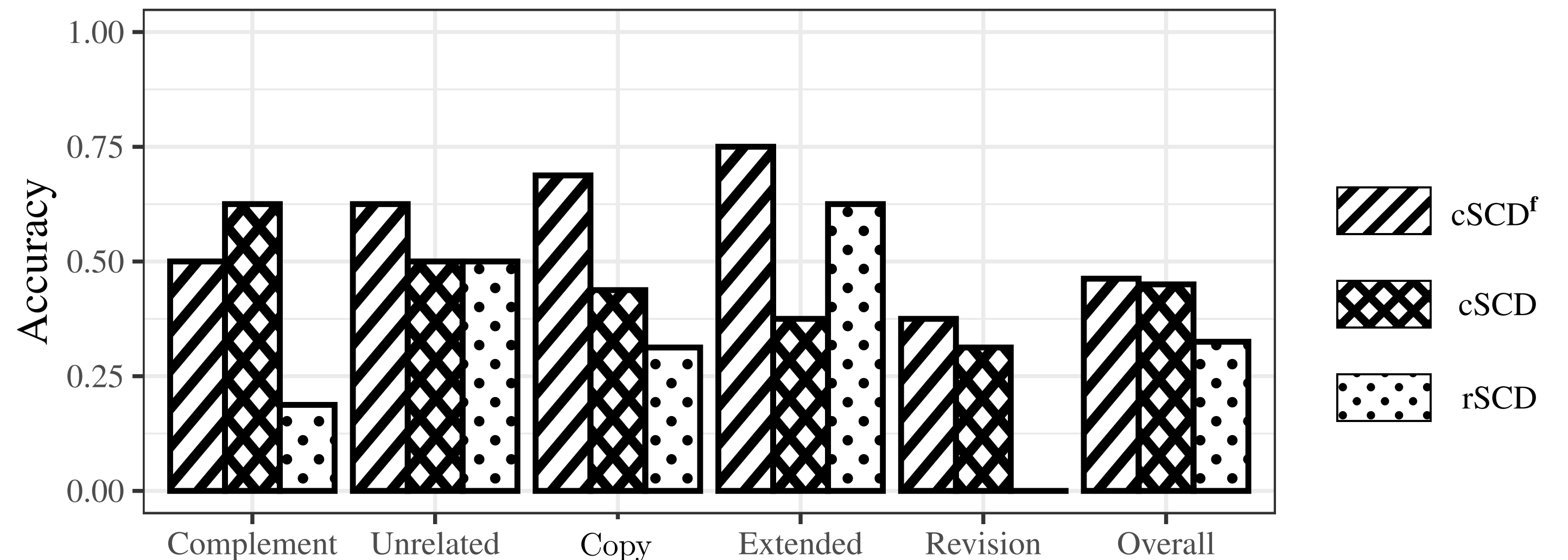


$c(d', d)$  only in 1. needed 😊  
Querying taxonomy quite costly.

# RESULTS: COMPLEMENT DETECTION BY COMBINED SCD MATRIX

## ■ Document classification accuracy using

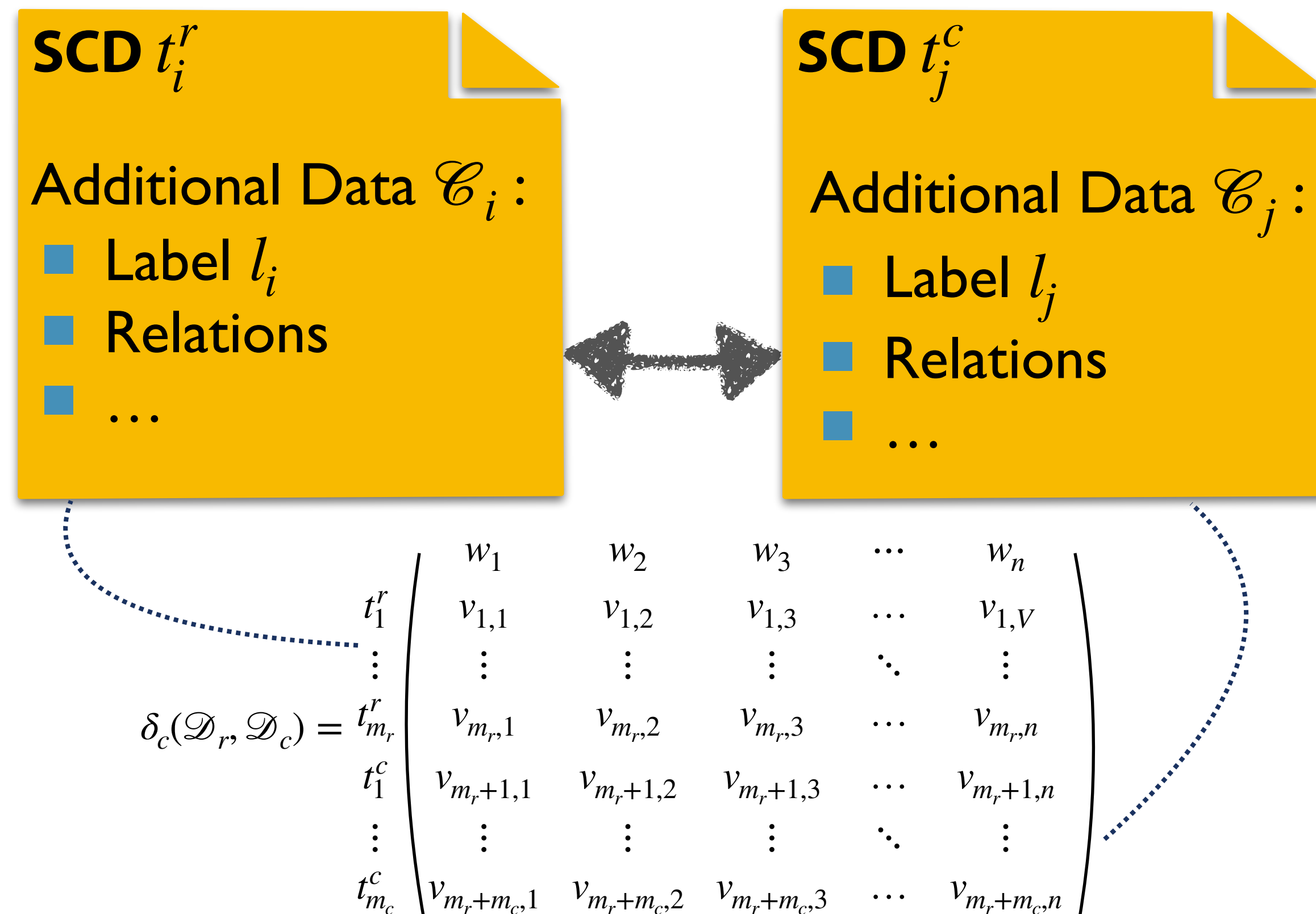
- Combined SCD matrix  
→ cSCD
- Filtered combined SCD matrix  
→ cSCD<sup>f</sup>
- Related (normal) SCD matrix  
→ rSCD





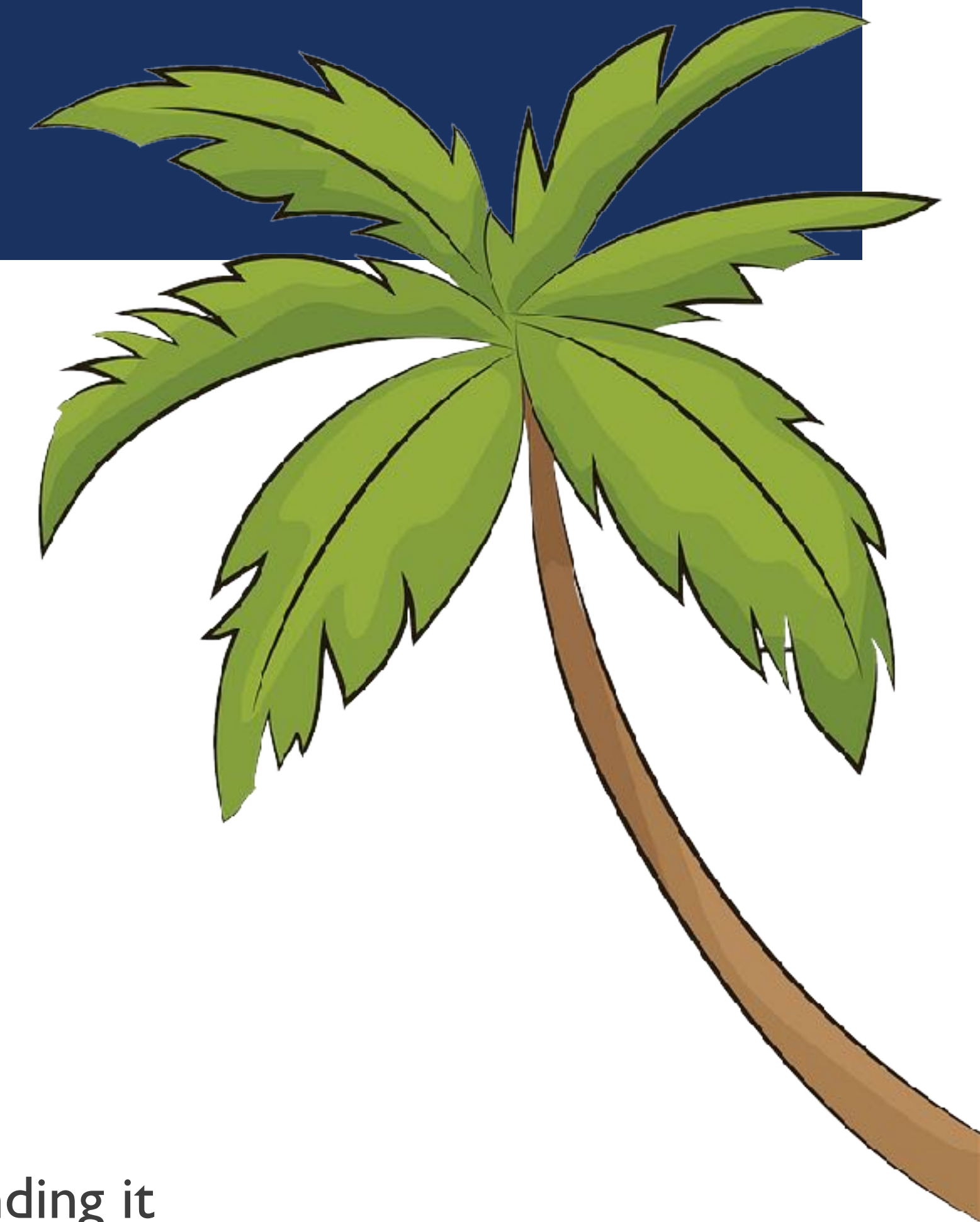
# GENERALISE RELATIONS AND COMBINED MATRIX

- Inter-SCD relations
  - Stored as links in additional data
  - Represented by combined SCD Matrix
    - Adapted MPS<sup>2</sup>CD yields adjusted similarity value
    - Apply techniques originally for related corpora
  - ➔ Example type complement used for corpus enrichment and document classification
  
- Intra-SCD relations
  - Referenced Sentences
  - Word-Distribution in matrix



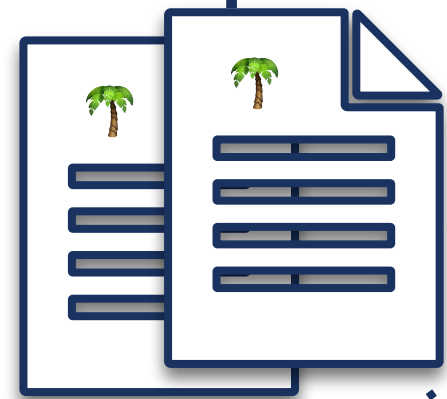
# INTERIM SUMMARY

1. Unsupervised Estimation of SCDs
    - SCDs (an SCD matrix) for any corpus
  2. Continuous Improvement by Feedback
    - Feedback from users used to update and enhance SCD matrix
  3. Labelling of SCDs
    - SCDs get a human friendly label for display and description
  4. Intra- and Inter-SCD Relations
    - Intra: Each SCD references sentences, has word distribution, and data incl. label
    - Inter: SCDs have relations, e.g., complement, among each other
- ➔ Apply SCD on any corpus (e.g., small and without initial SCDs) to help understanding it



Considered in Part 3

Corpus of Documents



USEM + LESS

$$\begin{matrix}
 & w_1 & w_2 & \dots & w_n \\
 \begin{matrix} t_1 \\ t_2 \\ \vdots \\ t_m \end{matrix} & \begin{pmatrix} v_{1,1} & v_{1,2} & \dots & v_{1,n} \\ v_{2,1} & v_{2,2} & \dots & v_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ v_{m,1} & v_{m,2} & \dots & v_{m,n} \end{pmatrix}
 \end{matrix}$$

Relations, e.g. Complement

Word Distribution  $\{v_{i,1}, \dots, v_{i,n}\}$

Referenced Sentences

$\{s_1, \dots, s_S\}$

SCDs  $g(\mathcal{D})$

SCD  $t_1$  Add. Data  $\mathcal{C}_1$  Label  $l_1$

SCD  $t_2$  Add. Data  $\mathcal{C}_2$  Label  $l_2$

SCD  $t_m$  Add. Data  $\mathcal{C}_m$  Label  $l_m$

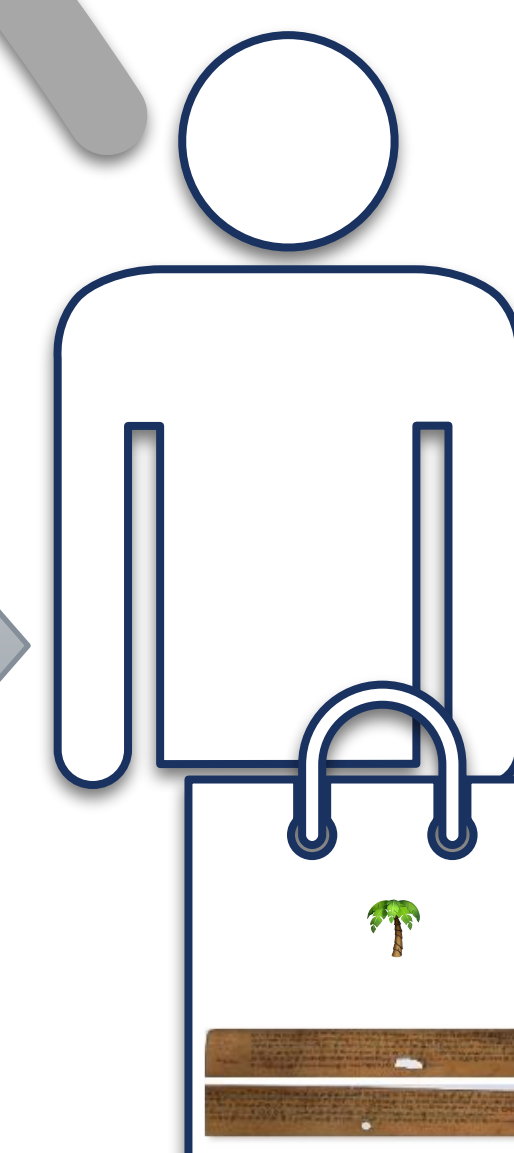
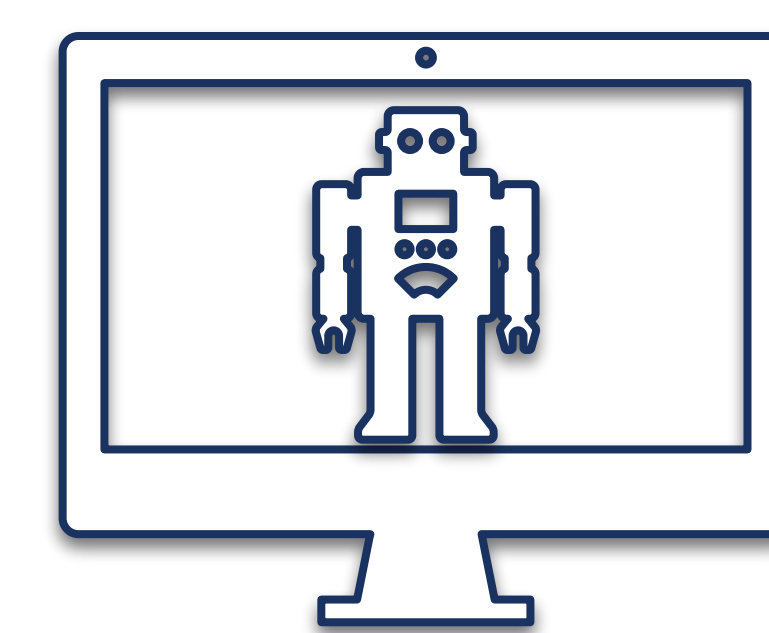
OVERVIEW IN  
DETAIL

Feedback (FRESH)

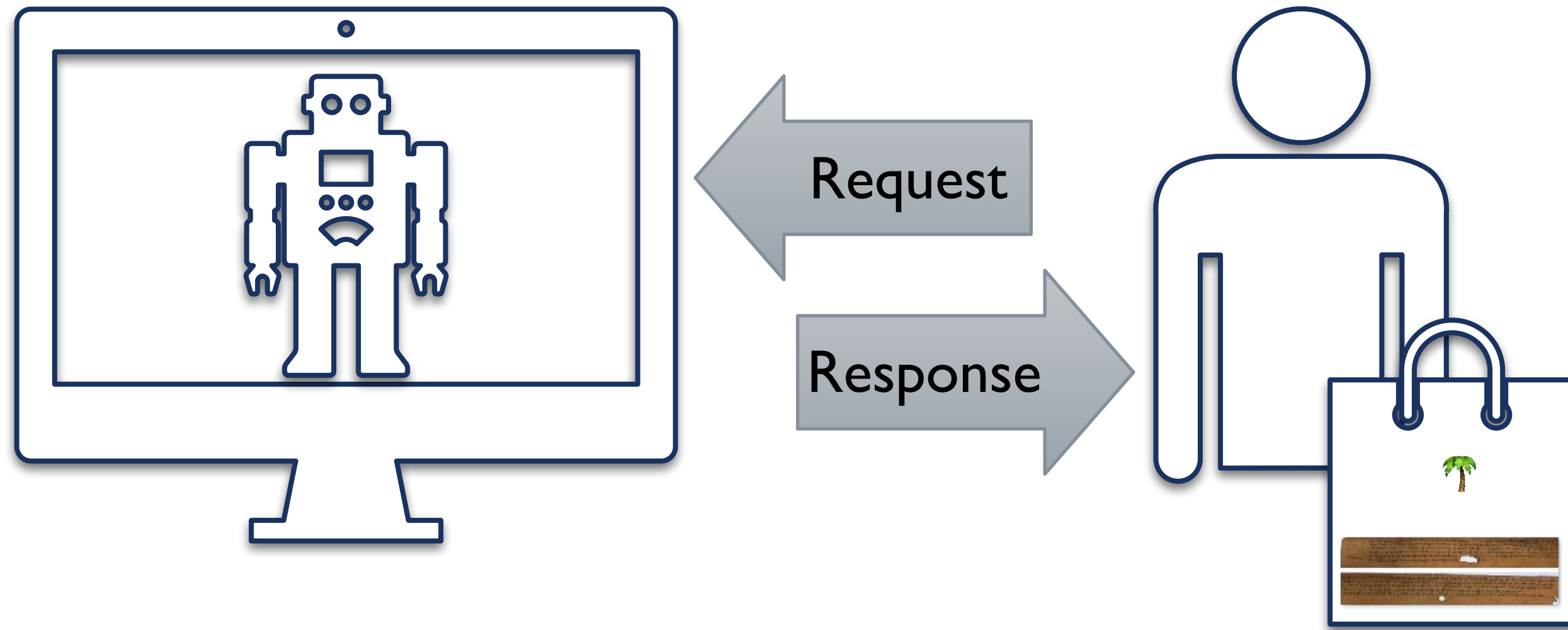
Used to  
Respond to  
Queries

Query

Response







# INFORMATION SYSTEM

BASED ON SCDS USING USEM, LESS, FRESH, ...



UNIVERSITÄT ZU LÜBECK



FROM MINIMAL DATA TO TEXT UNDERSTANDING

Unsupervised and Relational Learning

43



← Previous ↖ Parent → Next

## § 83

### *Foundation upon Death*

<sup>1</sup>If the foundation transaction consists of a disposition upon death, the probate court shall notify this to the competent authority for recognition, unless requested by the heir or the executor. <sup>2</sup>If the foundation transaction does not satisfy the requirements of §81 (1) 1 sentence 3, the foundation shall be given articles of association by the competent authority before recognition or incomplete articles of association shall be supplemented; the will of the founder shall be taken into account. <sup>3</sup>The seat of an foundation, unless otherwise provided, is the place where the administration is conducted. <sup>4</sup>In case of doubt, the founder's last domicile in Germany shall be deemed to be the registered office.

← Previous ↖ Parent → Next

Hover SCD Windows

1 Chapter 1 – Gen

2 Chapter 2 – Regi

2 Subtitle 2 – Foundatio

3 Subtitle 3 – Legal Pers

2 Section 2 – Things and Animals

3 Section 3 – Legal Transactions

1 Title 1 – Legal Capacity

2 Title 2 – Declaration of Inte

3 Title 3 – Contract

4 Title 4 – Conditions and Ti

5 Title 5 – Agency and Attorn

Search

[Book 1](#) / [Section 1](#) / [Title 2](#) / [Subtitle 1](#) / [Chapter 1](#) / § 24

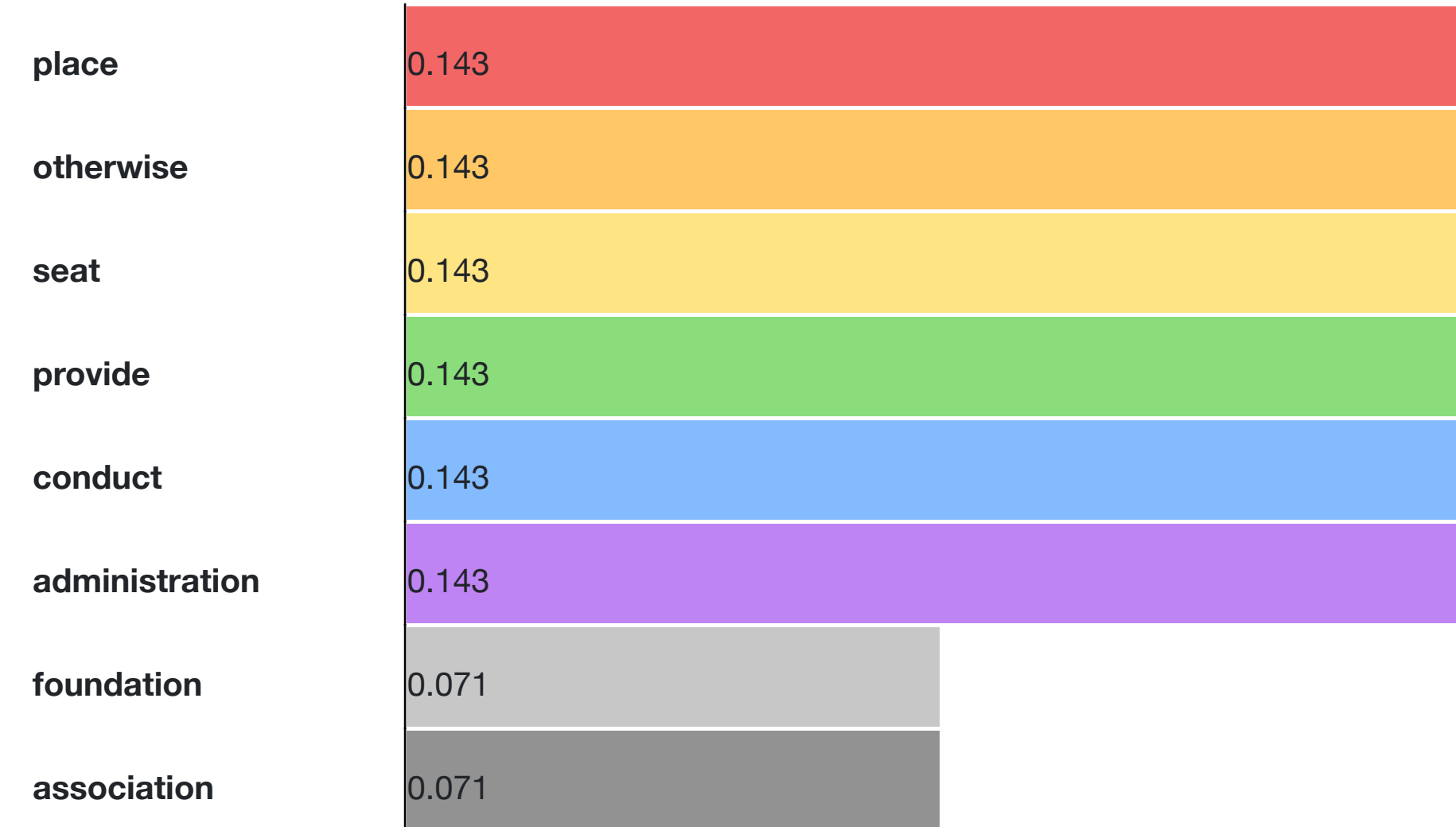
← Previous ↖ Parent → Next

### § 24 Seat

The seat of an association, unless otherwise provided, is the place where the administration is conducted.

← Previous ↖ Parent → Next

## Words of this SCD's Distribution



## Windows of this SCD 2

Paragraph   Window   Excerpt

> § 24  
 > 25  
 Seat  
 The seat of an association, unless otherwise provided, is the place where the administration is conducted.

> § 83  
 > 284  
 Foundation upon Death  
<sup>1</sup>If the foundation transaction consists of a disposition upon death, the probate court shall notify this to the competent authority for recognition, unless requested by the heir or the executor. <sup>2</sup>If the foundation transaction does not satisfy the requirements of § 81 (1) 1 sentence 3, the foundation shall be given articles of association by the competent authority before recognition or incomplete articles of association shall be supplemented; the will of the founder shall be taken into account. <sup>3</sup>The seat of an foundation, unless otherwise provided, is the place where the administration is conducted. <sup>4</sup>In case of doubt, the founder's last domicile in Germany shall be deemed to be the registered office.

# SCDs and Windows

[First](#)
[«](#)
[...](#)
[14](#)
[15](#)
[16](#)
[17](#)
[18](#)
[19](#)
[20](#)
[21](#)
[...](#)
[»](#)
[Last](#)

SCD	Top 10 Words	Windows	# Windows	
25	place otherwise seat provide conduct administration foundation association		<input type="text" value="↗ 25"/> <input type="text" value="↗ 284"/>	2
		<b>Label</b>		
28	legal representative status without		<input type="text" value="↗ 29"/>	2
		<b>Label</b>		
41	association seat district which		<input type="text" value="↗ 42"/> <input type="text" value="↗ 133"/>	2
		<b>Label</b>		



## Search Similar SCDs (using MPSCD)

The seat of an association shall be the place of its administration.

Search

## Results

> 25

0.7893522173763263

Label



Paragraph Window Excerpt

> § 24

> 25

Seat

The seat of an association, unless otherwise provided, is the place where the administration is conducted.

> § 83

> 284

Foundation upon Death

<sup>1</sup>If the foundation transaction consists of a disposition upon death, the probate court shall notify this to the competent authority for recognition, unless requested by the heir or the executor. <sup>2</sup>If the foundation transaction does not satisfy the requirements of § 81 (1) 1 sentence 3, the foundation shall be given articles of association by the competent authority before recognition or incomplete articles of association shall be supplemented; the will of the founder shall be taken into account. <sup>3</sup>The seat of an foundation, unless otherwise provided, is the place where the administration is conducted. <sup>4</sup>In case of doubt, the founder's last domicile in Germany shall be deemed to be the registered office.

> 41

0.49613893835683387



> 489

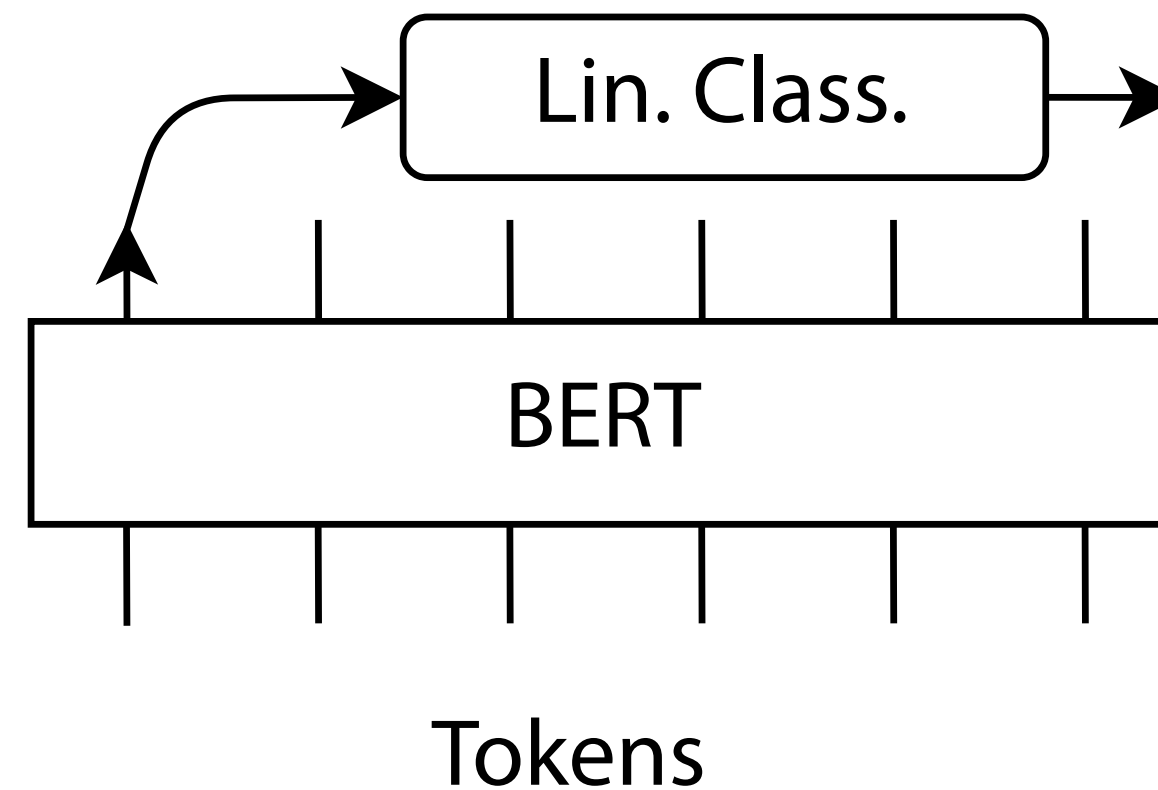
0.42499999999999993



① Identify SCDs among text – iSCD



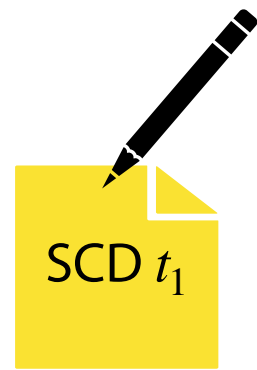
② Most probably suited SCD – MPS<sup>2</sup>CD



# LLMS IN REPLACEMENT FOR SCDS

ESTIMATING CONTEXT-SPECIFIC SCDS USING BERT

Backup



# TASK: APPLY BERT ON SCDS



## ① Identify SCDS among text – iSCD

- Given a text document  $d'$  where SCDS and content are interleaved
- Asked for set  $g(d)$  containing SCDS and the content of text document  $d \subseteq d'$

$d' = (\text{"We visited the bisons large animals in the zoo  
a place where non-domestic animals are exhibited."})$

$d = (\text{"We visited the bisons in the zoo."})$   
 $g(d) = \{ (\text{"large animals"}, 4),$   
 $(\text{"a place where non-domestic animals are exhibited"}, 7) \}$

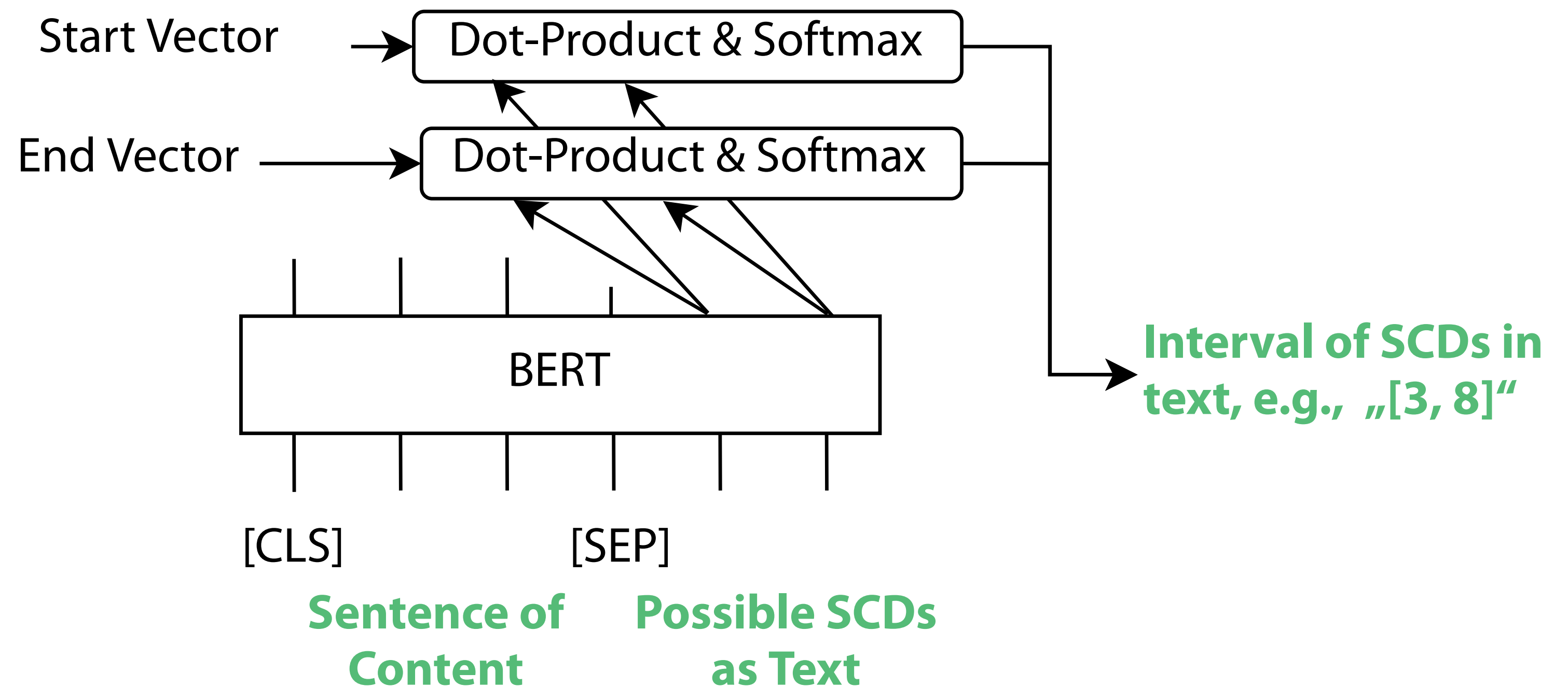
## ② Most probably suited SCD – MPS<sup>2</sup>CD

- Given a text document  $d$  without associated SCDS
- Asked for set  $g(d)$  containing best suited SCDS  $t$  for  $d$

$d = (\text{"We visited the bisons in the zoo."})$

# APPROACH: APPLYING BERT ON SCDS

- iSCD
- BERT Classify
- BERT Next
- MPS<sup>2</sup>CD
- BERT Choose
- BERT Highlight ←





# EVALUATION

## ■ Corpus

- 20 newsgroups
- Definitions from Wiktionary

## ■ Dataset

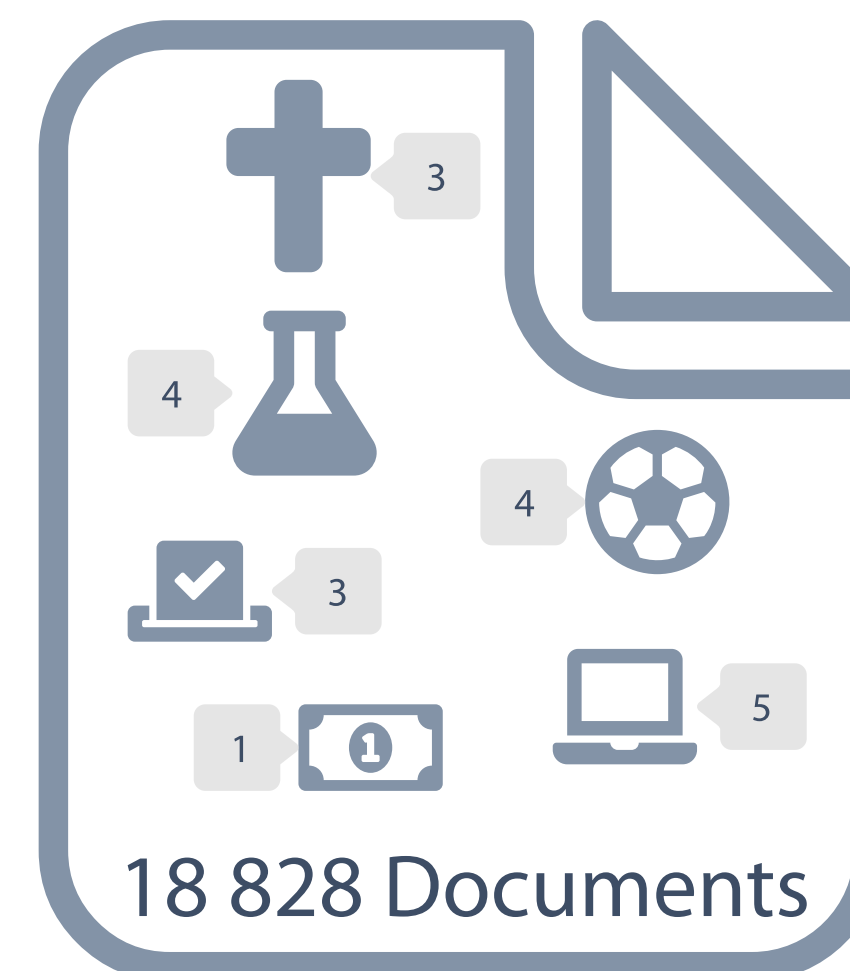
- 80 % training and 20 % testing
- Disjoint and same sets of definitions for SCDs

## ■ Hardware

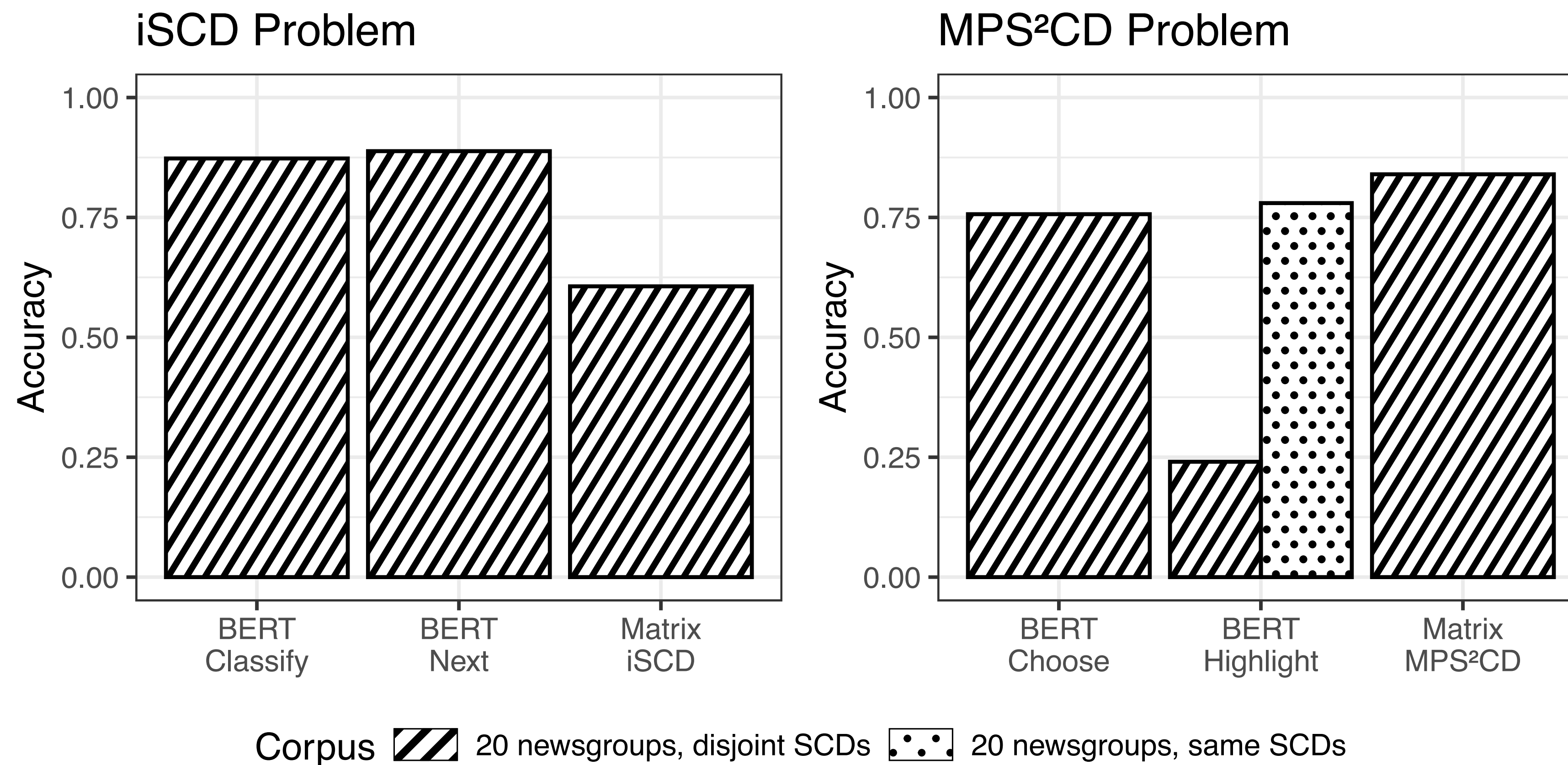
- NVIDIA DGX A100 320GB
- 8 Intel 6248 with 2.50GHz (3.90GHz), 16GB RAM

## ■ Model

- “Bert-Base-Uncased“

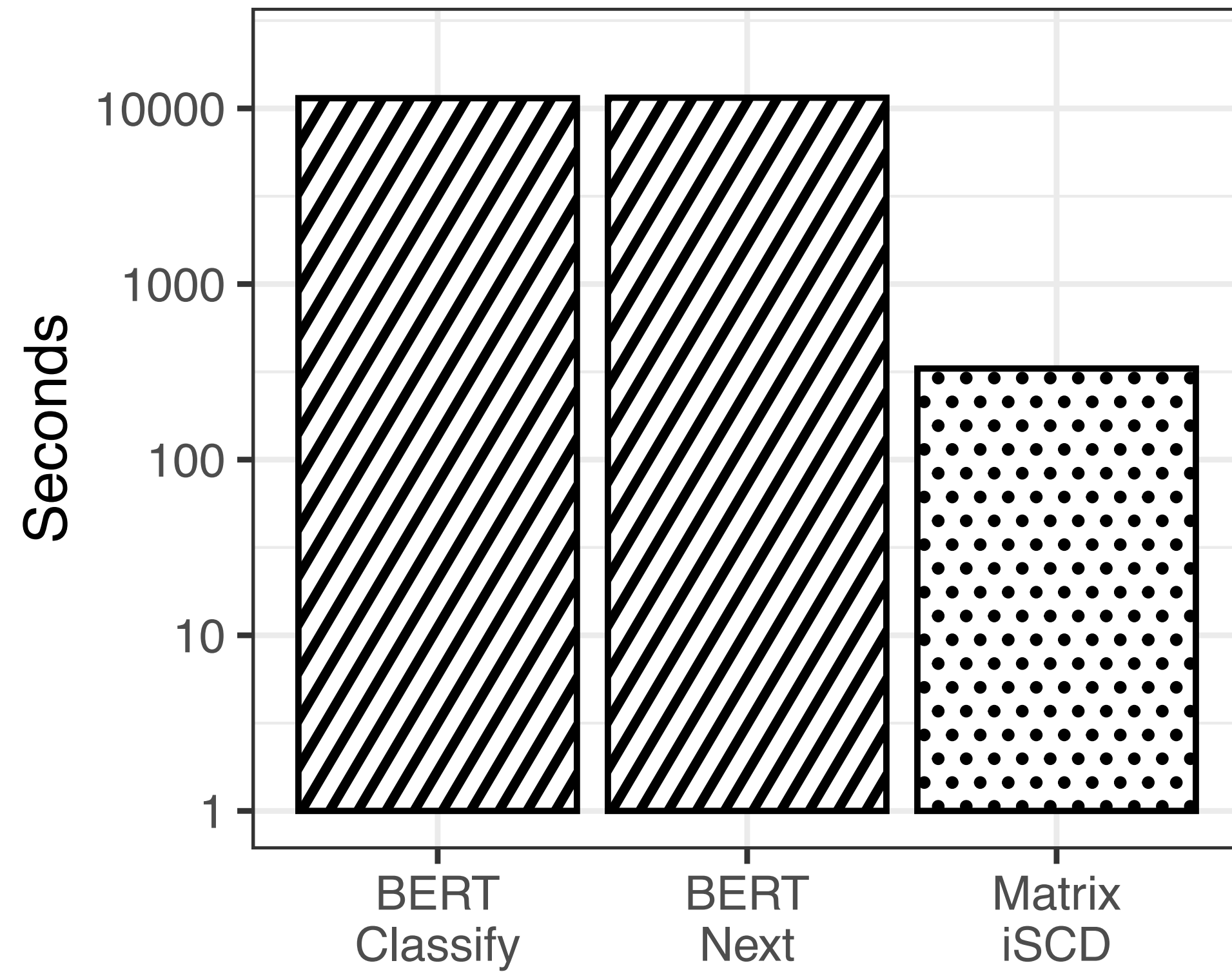


# RESULTS: ACCURACY

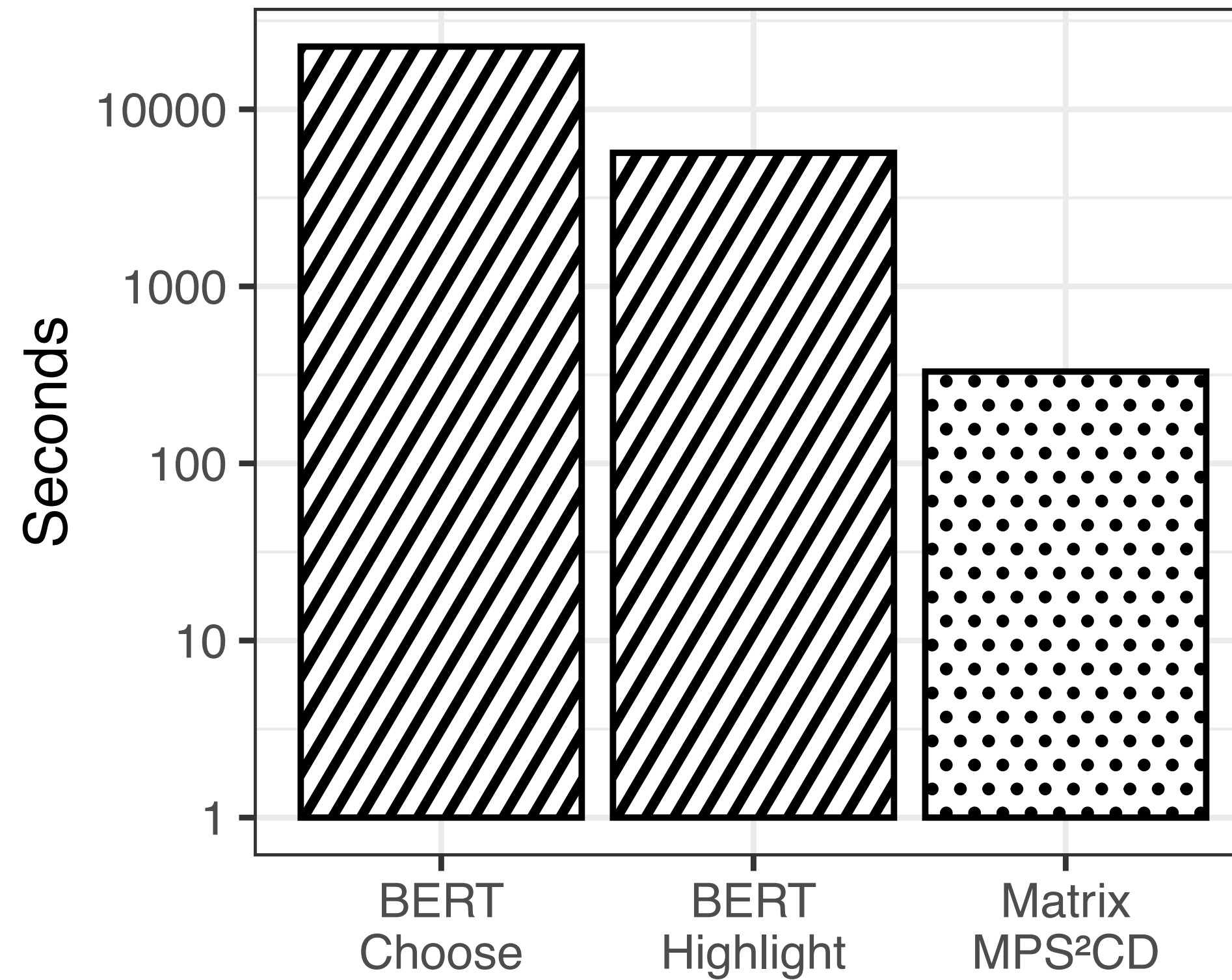


# RESULTS: RUNTIME

## iSCD Problem



## MPS<sup>2</sup>CD Problem



# CONCLUSION

- BERT and the SCD matrix solve the MPS<sup>2</sup>CD and iSCD problem well
- BERT needs much more time and computational resources in contrast to the SCD matrix

„We demonstrate that BERT is able to grasp the concept of SCDs, in a way that BERT can be trained to solve SCD-related tasks.“



# AGENDA

1. Introduction to Semantic Systems [Tanya]
2. Supervised Learning [Marcel]
3. Unsupervised and Relational Learning [Magnus]
  - Unsupervised Estimation of SCDs
  - Continuous Improvement by Feedback
  - Labelling of SCDs
  - Inter- and Intra-SCD Relations
4. Summary [Tanya]

