

# TO EXTEND OR NOT TO EXTEND? COMPLEMENTARY DOCUMENTS

MAGNUS BENDER<sup>1</sup>, FELIX KUHR<sup>1</sup>, TANYA BRAUN<sup>2</sup>



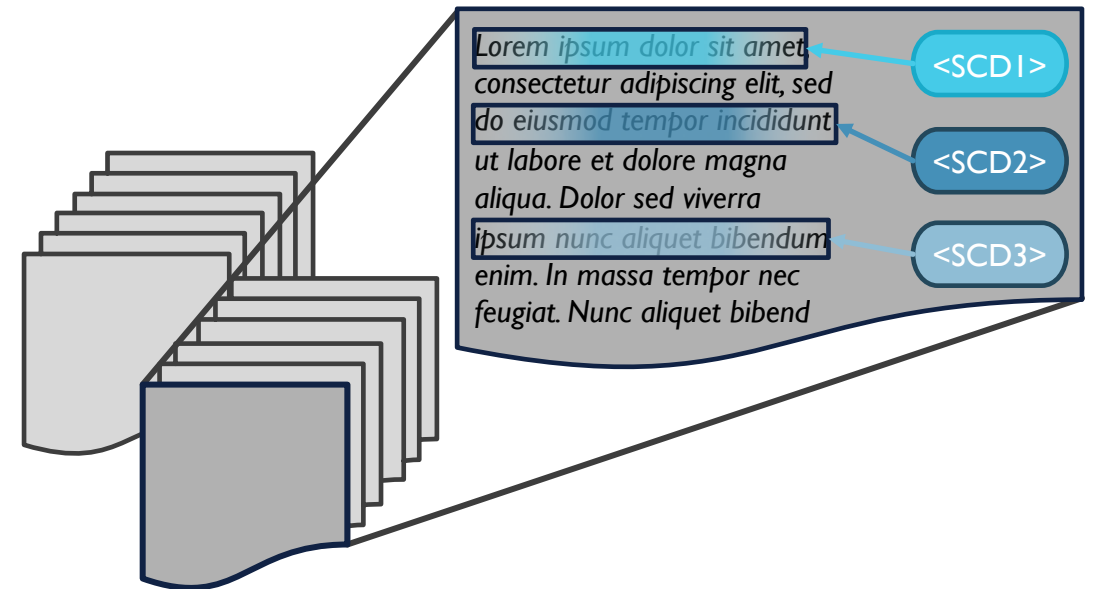
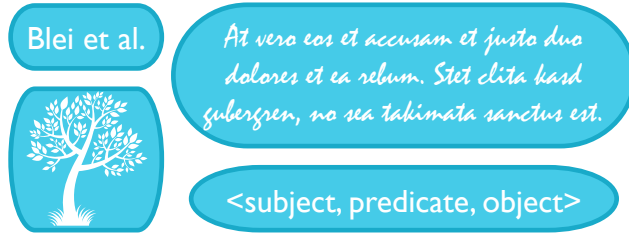
UNIVERSITÄT ZU LÜBECK

<sup>1</sup>Institute of Information Systems, University of Lübeck  
<sup>2</sup>Computer Science Department, University of Münster



# THE SETTING: A CORPUS OF DOCUMENTS AND ANNOTATIONS

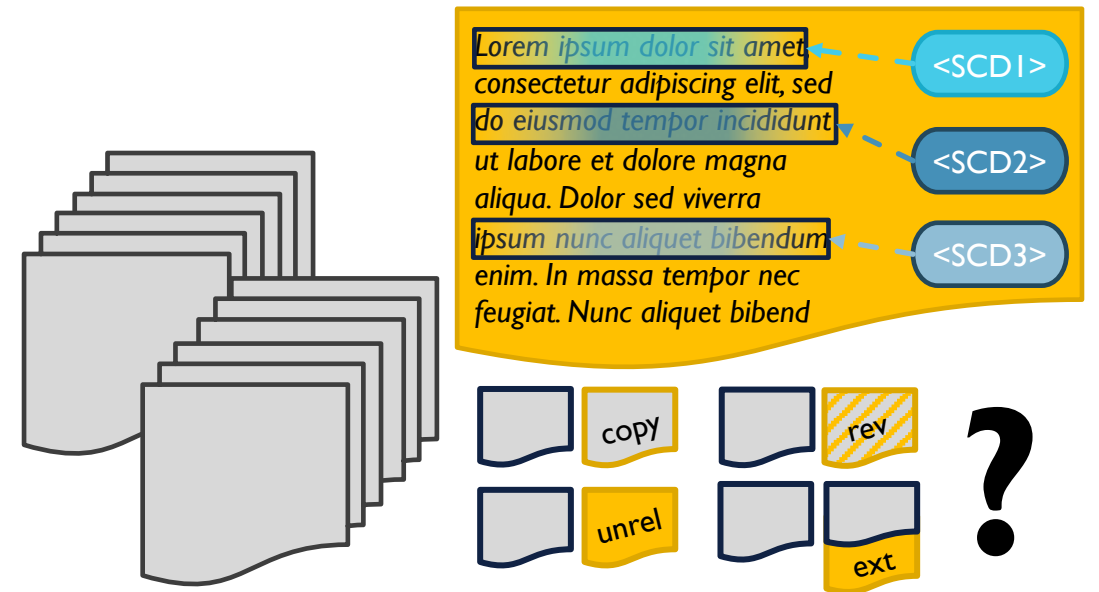
- Corpus = set of documents  $\mathcal{D}$
- Each document  $d$  has a set of annotations  $g(d)$ 
  - Annotation  $\triangleq$  *subjective content description* (SCD)
  - Reflect the *context* of the purpose of the corpus
- Types of SCDs can be manifold
  - Figures, notes, references, ...
- Each SCD associated with words at specific location
  - Assumption: Words closer to location, influence higher



**Proposition 1:**  
SCDs generate the words in a document

# TASKS: DOCUMENT RETRIEVAL & CORPUS ENRICHMENT

- Document retrieval user-driven
  - External task
  - Well-rounded corpus needed for high-quality retrieval
- Corpus enrichment to extend corpus with documents that provide *added* value in task context
  - Internal task
  - Classification problem
    - Input: new document  $d$ , corpus  $\mathcal{D}$
    - Classify  $d$  as: quasi-copy, revision, extension, unrelated
      - Documents with assumed added value

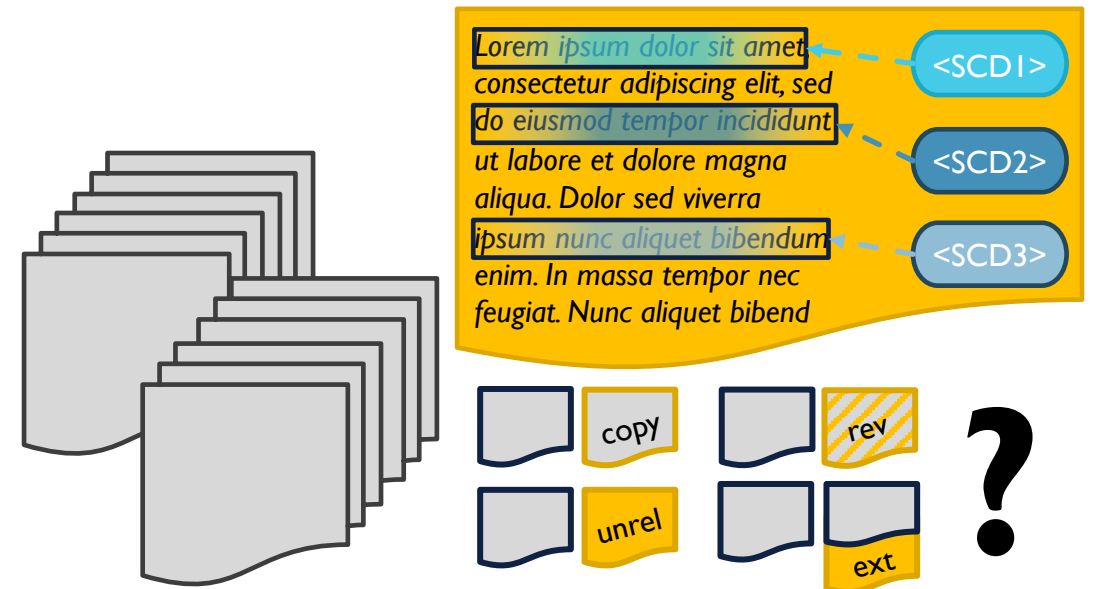


Based on Proposition 1:  
How much of  $d$  can SCDs of  $\mathcal{D}$  generate?

# PROBLEM: SIMILARITY AS A MEASURE

- 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.

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



# DREAM SOLUTION: COMPLEMENTARITY

- Goal: identify documents that are complementary to a corpus / a document in a corpus
  - Binary classification problem:  
*Complement = true or Complement = false*
- *BUT: Numbers-based representations make it hard to grasp complementarity on a semantic level*

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

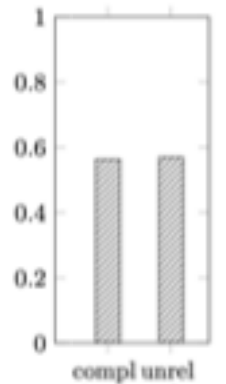


- Corpus on sporting events
  - Olympics 2020, UEFA Euro 2020



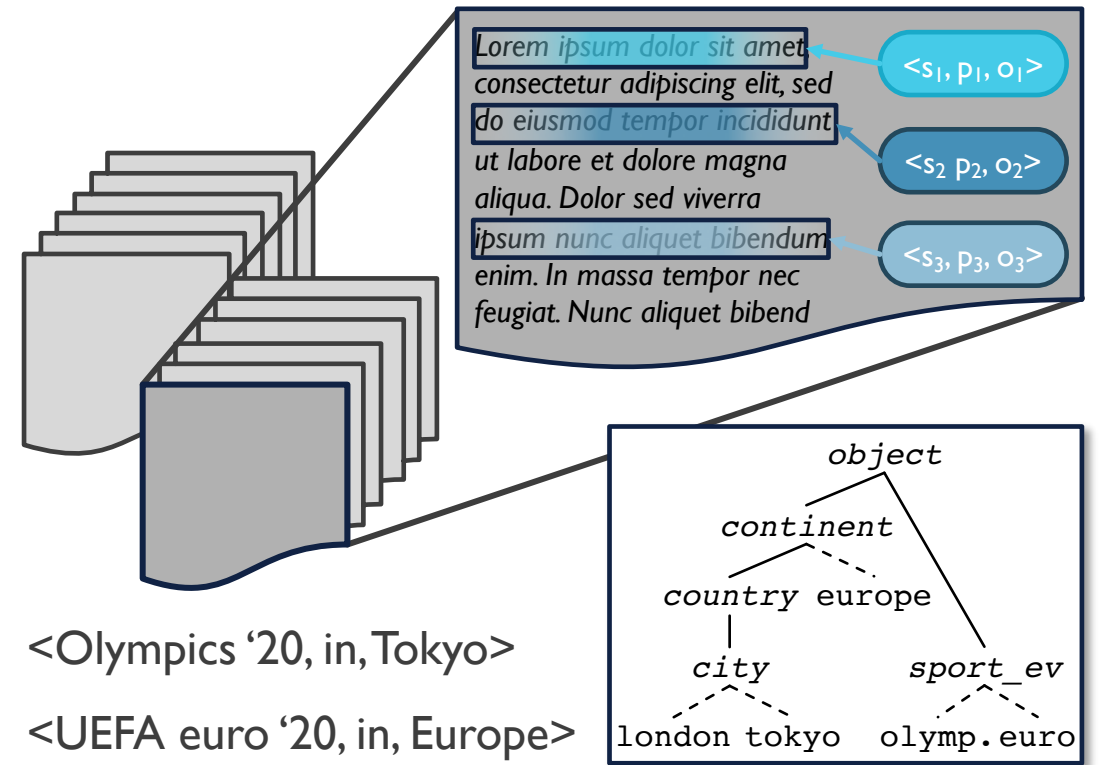
- Complementary documents on
  - Covid-19 spread in cities

- ❖ Words very different compared to corpus
  - ❖ Different (topic / SCD) distributions
  - ❖ Likely to be classified as unrelated
    - ❖ Fig.: Similarity values of complements and unrelated documents for corpus enrichment



# NEW PROBLEM: HOW TO GRASP COMPLEMENTARITY ON A FORMAL LEVEL?

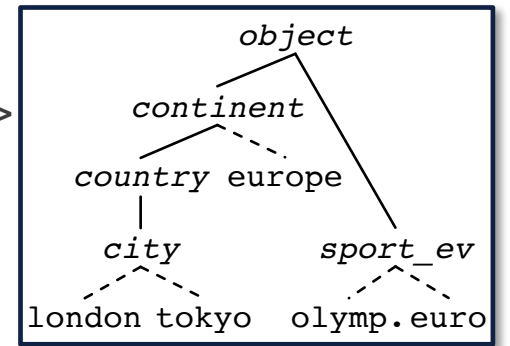
- Use SCDs specifically in the form of **SPO triples**
  - SPO triple: <subject, predicate, object>
  - Extract for any document, e.g., with OpenIE tools together with a taxonomy
    - Hierarchy of concepts
    - Dictionary of synonyms
- Allow for grasping complementarity on a semantic level by
  - Looking at shared concepts in the SPO triples
  - While also accounting for hierarchy and synonyms



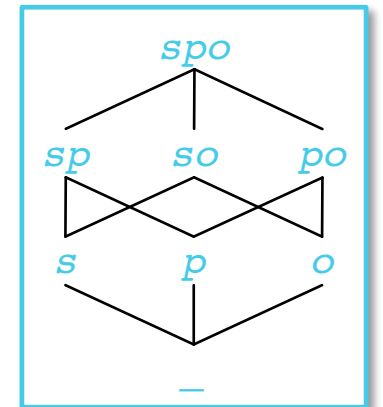
# A FORMAL DEFINITION: COMPLEMENTARY SCDS

- Let  $x^\uparrow$  refer to the concept or meaning of  $x$
- Seven types of complementarity between SCDS  $t_i, t_j$ 
  1.  $s$   $t_i = \langle s^\uparrow, p, o \rangle, t_j = \langle s^\uparrow, p, o \rangle$
  2.  $p$   $t_i = \langle s, p^\uparrow, o \rangle, t_j = \langle s, p^\uparrow, o \rangle$
  3.  $o$   $t_i = \langle s, p, o^\uparrow \rangle, t_j = \langle s, p, o^\uparrow \rangle$
  4.  $sp$   $t_i = \langle s^\uparrow, p^\uparrow, o \rangle, t_j = \langle s^\uparrow, p^\uparrow, o \rangle$
  5.  $so$   $t_i = \langle s^\uparrow, p, o^\uparrow \rangle, t_j = \langle s^\uparrow, p, o^\uparrow \rangle$
  6.  $po$   $t_i = \langle s, p^\uparrow, o^\uparrow \rangle, t_j = \langle s, p^\uparrow, o^\uparrow \rangle$
  7.  $spo$   $t_i = \langle s^\uparrow, p^\uparrow, o^\uparrow \rangle, t_j = \langle s^\uparrow, p^\uparrow, o^\uparrow \rangle$
  - Types get less strict  $\rightarrow$  Order in lattice

- $t_1: \langle \text{Olympics '20, in, Tokyo} \rangle$
- $t_2: \langle \text{UEFA euro '20, in, Europe} \rangle$
- $t_3: \langle \text{Covid-19, in, Tokyo} \rangle$
- $t_4: \langle \text{Covid-19, in, London} \rangle$



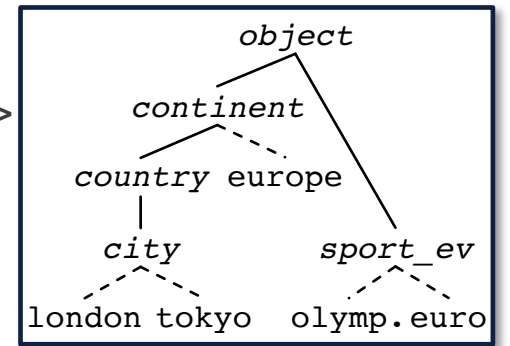
- $t_1, t_3$   $s$ -complementary
  - $s_1, s_3$  share *object*;  $p_1 = p_3; o_1 = o_3$
  - and  $sp, so, spo$ -complementary
- $t_1, t_4$   $so$ -complementary (+ $spo$ )
  - Same holds for  $t_2, t_3$  and  $t_2, t_4$
- If deleting *object* in taxonomy,  
no complementarity!



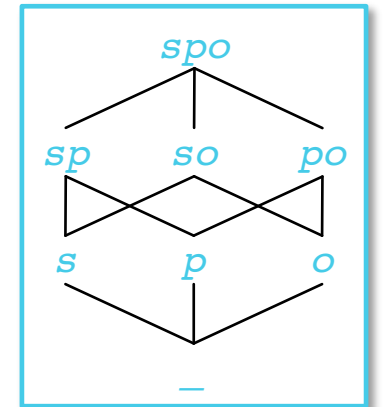
# A FORMAL DEFINITION: COMPLEMENTARY DOCUMENTS – PRELIMINARIES

- Let  $\mathfrak{C}_x(t_i, t_j)$ ,  $x \in \mathcal{X} = \{s, p, o, sp, so, po, spo\}$  be an **indicator function**
  - Returns 1 if  $t_i, t_j$   $x$ -complementary; otherwise 0
  - $\mathfrak{C}_x$  is symmetric, i.e.,  $\mathfrak{C}_x(t_i, t_j) = \mathfrak{C}_x(t_j, t_i)$
- Assign **weights**  $w_x$ ,  $\sum_{x \in \mathcal{X}} w_x = 1$ , to complementarity types  $x$  to encode which complementarity interested in
  - Depends on corpus composition and desired outcome
  - $w_x = 0$  if  $x$ -complementarity uninteresting
    - $w_{sp} = 1$ , rest 0: complementary SCDs specific to object  $o$
  - $w_x \neq 0$  only for types of same level in lattice sensible
    - $w_x = \frac{1}{3}$  for  $x \in \{sp, so, po\}$ , rest 0

- $t_1$ : <Olympics '20, in, Tokyo>
- $t_2$ : <UEFA euro '20, in, Europe>
- $t_3$ : <Covid-19, in, Tokyo>
- $t_4$ : <Covid-19, in, London>



- Pairs  $(t_1, t_3)$ ,  $(t_1, t_4)$ ,  $(t_2, t_3)$ ,  $(t_2, t_4)$ 
  - $\mathfrak{C}_s(t_1, t_3) = 1$ , for other pairs 0
  - $\mathfrak{C}_{sp}(t_1, t_3) = 1$ , for other pairs 0
  - $\mathfrak{C}_{so} = \mathfrak{C}_{spo} = 1$  for all pairs
  - $\mathfrak{C}_{po} = 0$  for all pairs





# A FORMAL DEFINITION: COMPLEMENTARY DOCUMENTS – PRELIMINARIES

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

- 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:

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

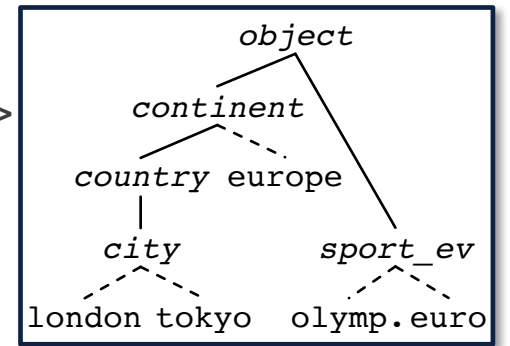
- $c$  is symmetric, i.e.,  $c(d', d) = c(d, d')$
- E.g., with  $w_{sp} = 1$ , rest 0:

$$c(d_1, d_2) = 1 + 0 + 0 + 0 = 1$$

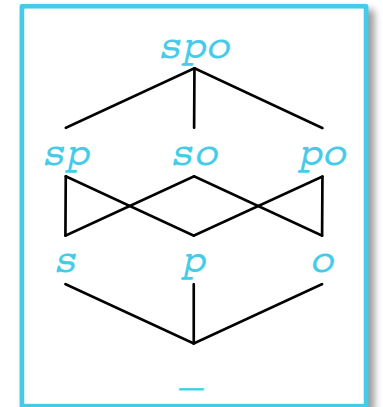
- E.g., with  $w_x = \frac{1}{3}$  for  $x \in \{sp, so, po\}$ , rest 0:

$$c(d_1, d_2) = \frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 0 + 3 \left( \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 0 \right) = 1.667$$

- $d_1$ 
  - $t_1$ : <Olympics '20, in, Tokyo>
  - $t_2$ : <UEFA euro '20, in, Europe>
- $d_2$ 
  - $t_3$ : <Covid-19, in, Tokyo>
  - $t_4$ : <Covid-19, in, London>



- Pairs  $(t_1, t_3)$ ,  $(t_1, t_4)$ ,  $(t_2, t_3)$ ,  $(t_2, t_4)$ 
  - $\mathfrak{C}_s(t_1, t_3) = 1$ , for other pairs 0
  - $\mathfrak{C}_{sp}(t_1, t_3) = 1$ , for other pairs 0
  - $\mathfrak{C}_{so} = \mathfrak{C}_{spo} = 1$  for all pairs
  - $\mathfrak{C}_{po} = 0$  for all pairs

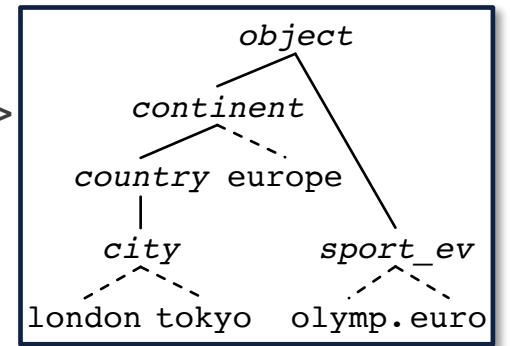


# A FORMAL DEFINITION: COMPLEMENTARY DOCUMENTS – DEFINITION

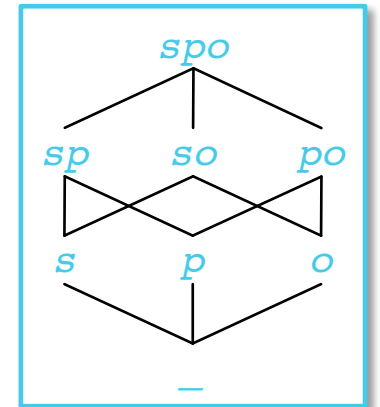
- Given a taxonomy  $\xi$
- Given weights  $w_x, x \in \mathcal{X}$
- Let  $\mathcal{C}_x(t_i, t_j), x \in \mathcal{X}$ , be an indicator function
- Let  $c(d', d)$  denote the complementarity value between documents  $d', d$  with SCDs  $g(d'), g(d)$
- Then, given a threshold  $\theta_d, d'$  is a **complement** to  $d$  if

$$c(d', d) > \theta_d$$

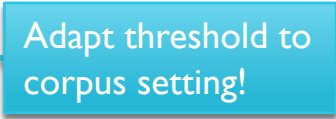
- $d_1$   $t_1$ : <Olympics '20, in, Tokyo>
- $t_2$ : <UEFA euro '20, in, Europe>
- $d_2$   $t_3$ : <Covid-19, in, Tokyo>
- $t_4$ : <Covid-19, in, London>



- $d_2$  complement to  $d_1$ ? ( $\theta_d = 1.5$ )
  - $w_{sp} = 1$ :  
 $c(d_2, d_1) = 1$  ✗
  - $w_x = \frac{1}{3}, x \in \{sp, so, po\}$ :  
 $c(d_2, d_1) = 1.667$  ✓



# USING THE NEW DEFINITION: CORPUS ENRICHMENT – THE COMPLEMENTARITY VERSION

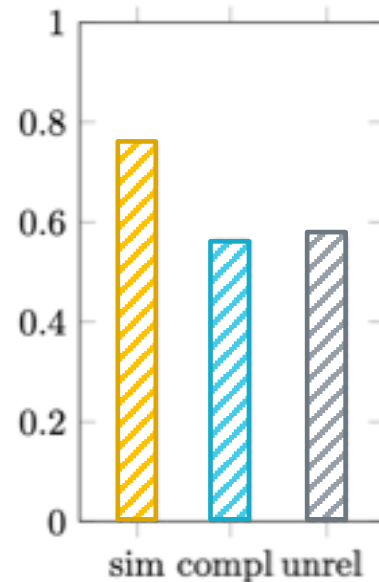
- Inputs:
  - Corpus  $\mathcal{D}$ 
    - Includes SCDs including SPO triples, a taxonomy  $\xi$
  - New document  $d'$
  - Threshold  $\theta_{\mathcal{D}}$  
  - Weights  $\{w_x\}_{x \in \mathcal{X}}$
- Returns  
*Complement = true or Complement = false*
- Goes through all documents  $d \in \mathcal{D}$  and adds up the complementary value  $c(d', d)$  to a corpus value  $c$

```
function extendComplement( $\mathcal{D}, d', \theta_{\mathcal{D}}, \{w_x\}_{x \in \mathcal{X}}$ )  
  if  $g(d') = \emptyset$  then  
    Add SCDs to  $d'$  using OpenIE  
   $c \leftarrow 0$   
  for each  $t_i \in g(d')$  do  
    for each  $d \in \mathcal{D}$  do  
      for each  $t_j \in g(d)$  do  
        for each  $x \in \mathcal{X}$  do  
           $c \leftarrow c + w_x \mathfrak{C}_x(t_i, t_j)$   
  if  $c > \theta_{\mathcal{D}}$  then  
    return true  
  return false
```

# USING THE NEW DEFINITION: CORPUS ENRICHMENT – THE COMPLEMENTARITY VERSION

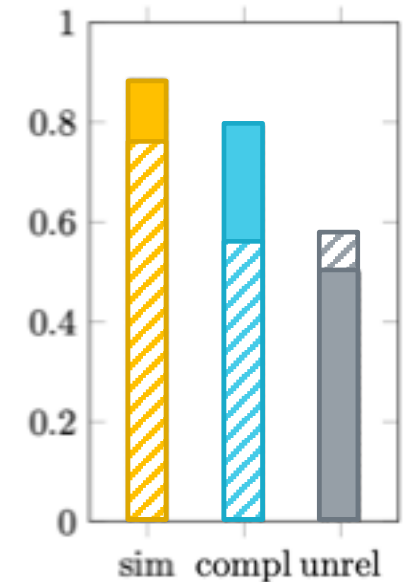
## Similarity-based method

- Classifies both **complements** and **unrelated** documents as **unrelated**
- Identifies **quasi-copies** (sim)



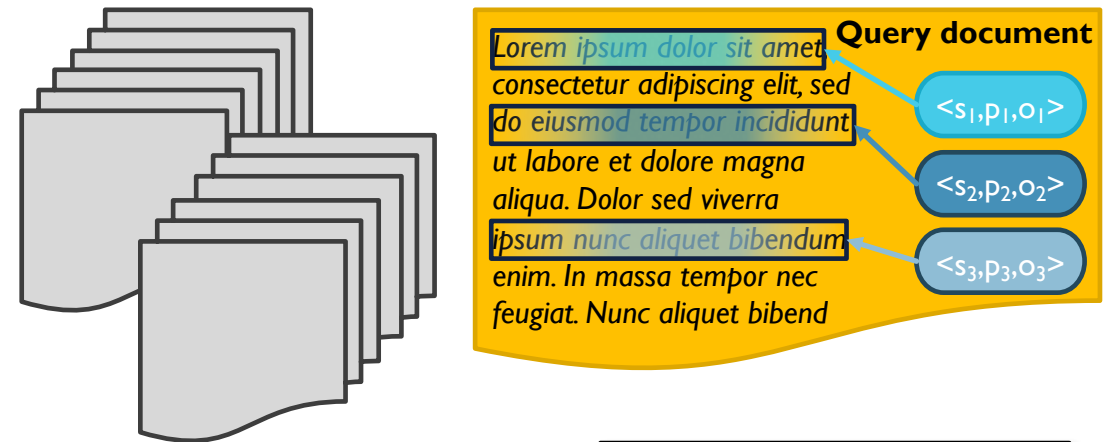
## Complement-based method

- Classifies **complements** as **complements** and **unrelated** documents as **non-complements**
- **Quasi-copies** (sim) have even higher complementarity values than **complements**
  - *Rationale:* Equal SPO triples match all complementarity types



# USING THE NEW DEFINITION: DOCUMENT RETRIEVAL – THE COMPLEMENTARITY VERSION

- Input:
  - Corpus  $\mathcal{D}$ 
    - Includes SCDs including SPO triples, a taxonomy  $\xi$
  - Query document  $d'$
  - Weights  $\{w_x\}_{x \in \mathcal{X}}$
- Extract SPO triples for  $d'$  if none exist
- Compute  $c(d', d)$  for each document in  $d \in \mathcal{D}$
- Return
  - Given threshold  $\theta_d$ : all documents  $d$  with  $c(d', d) > \theta_d$
  - Given a number  $k$ :  $k$  documents with highest  $c(d', d)$



Given threshold  $\theta_d$ :  
 $d_1: c(d', d_1) > \theta_d \rightarrow$  return  
 $d_2: c(d', d_2) \not> \theta_d$   
 $d_3: c(d', d_3) \not> \theta_d$   
 $d_4: c(d', d_4) > \theta_d \rightarrow$  return  
⋮

Ranking,  $k = 3$ :  
 $d_1: c(d', d_1) \rightarrow$  return  
 $d_{12}: c(d', d_{12}) \rightarrow$  return  
 $d_4: c(d', d_4) \rightarrow$  return  
 $d_3: c(d', d_3)$   
⋮

# CONCLUSION: A FORMAL DEFINITION OF COMPLEMENTARITY

- Formal definition of **complementary documents** based on SCDs in the form of SPO triples
  - Grasping complementarity on a semantic level
- Solve tasks using complementarity definition
  - Corpus enrichment, document retrieval
- Future work
  - Include “degree of ancestry” in indicator function
    - Also to distinguish better between copies and complements
  - Deal with uncertainty / unknown concepts

