

Determining Semantic Similarity among Entity Classes from Different Ontologies

Authors: M. Andrea Rodriguez, Max J. Egenhofer

Presented By: Yaoxuan Wang



- Introduction
- Approach
- Experiments
- Conclusion
- Discussion

- **Motivation**
 - Need new tools that can improve the retrieval and integration of information
- **Previous**
 - Compare concepts from different ontologies are based on an a priori integration of local ontologies
- **This work**
 - Create a computational model to assess semantic similarity among entity classes from unconnected and independent ontologies
 - No integration of ontologies is needed

- Entity Class:
 - refers to concepts that group entities or objects of the real world into classes
 - e.g. building, lake , city, etc.
- Entity Class Representation:
 - **Synonym set**
 - a set of synonym words that denotes an entity class
 - address polysemy and synonymy
 - e.g. bank, depository financial institution
 - **Semantic interrelation**
 - determine general organization of entity classes
 - Hyponymy, i.e. "is-a" relation equal to "inheritance"
 - Meronymy, i.e. "part-whole" relation equal to "composition"
 - **Distinguishing features**
 - distinguish entity classes from the same superclass besides semantic interrelation
 - e.g. hospital vs. apartment building
 - **Classification**
 - functions, what is done to or with instances of a class, e.g. practice of stadium
 - parts, structural elements of a class, e.g. roof, floor of a building
 - attributes, the rest

- Recall:
 - the fraction of similar entity classes that are detected by the model

$$recall = \frac{|A \cap B|}{|A|}, \quad (13a)$$

- Precision:
 - the fraction of entity classes detected by the model that are actually similar

$$precision = \frac{|A \cap B|}{|B|}. \quad (13b)$$

- **A** is the set of similar entity classes
- **B** is the set of similar entity classes calculated by the model
- $| \cdot |$ is the counting measure or cardinality

TABLE 1
Entity_Class Definition in BNF Notation, and an Example of the Definition of *Stadium*

BNF Notation	Example: <i>Stadium</i>
<pre> <entity_class> ::= entity_class { name: {<syn_set>} description: <description> is_a: <is-a> part_of: <part_of> whole_of: <whole_of> parts: <parts> functions: <functions> attributes: <attributes>} <is_a> ::= {} {<pts_entity_classes>} <part_of> ::= {} {<pts_entity_classes>} <whole_of> ::= {} {<pts_entity_classes>} <parts> ::= {} {<syn_sets>} <functions> ::= {} {<syn_sets>} <attributes> ::= {} {<syn_sets>} <syn_sets> ::= {<syn_set>} {<syn_sets>, <syn_set>} <syn_set> ::= <word> <syn_set>, <word> <description> ::= <word> <description> <word> <pts_entity_classes> ::= <pointer> <pt_to_entity_classes>, <pointer> </pre>	<pre> entity_class { name: {stadium,bowl,arena} description: large often unroofed structure in which athletic events are held is_a: {construction*} part_of: {} whole_of: {athletic_field*} parts: { {athletic_field,sports_field,playing_field}, {dressing_room},{foundation}, {midfield},{spectator_stands,stands}, {ticket_office,box_office,ticket_booth} } functions: { {play,compete},{play,practise}, {recreate,play} } attributes: { {architectural_property}, {covered/uncovered},{name}, {lighted/unlighted},{owner_type}, {sports_type},{user_type} } } </pre>

(*x** denotes a pointer to the entity class *x*)

- **Entity Class Representation Comparison**
 - the representation of entity classes across ontologies should share some components
- **Three independent similarity assessments**
 - **Synonym sets**
 - synonyms refer to the same entity class
 - the similarity between an entity class and itself is maximal
 - exploit the general agreement in the use of words
 - detect equivalent words that likely refer to the same entity class
 - BUT only a basic level of similarity assessment
 - e.g. clinic vs. hospital, polysemy
 - **Distinguishing feature**
 - determine how similar entity classes are
 - Only have some common feature, e.g. stadium vs. sports_arena
 - **Semantic relations**
 - whether target entity classes are related to the same set of entity classes
 - e.g. hospital vs. house – building
 - how to compare?
 - a comparison between the **semantic neighborhoods** of entity classes

- **Semantic neighborhood (N)**

$$N(a^o, r) = \{c_i^o\} \text{ such that } \forall i \ d(a^o, c_i^o) \leq r. \quad (1)$$

- a^o : the entity class
- c_i^o : the set of entity classes whose distance $d() \leq r$
- r : nonnegative integer
- $d()$: the shortest path in the ontology between two entity classes,
i.e. the smallest number of arcs
- the semantic neighborhood of an entity class also contains itself

- **Similarity of semantic relation**

- NOT based on this path distance
- Path distance \Rightarrow neighborhoods
- the similarity between entity class A and B depends on the similarity of the entity classes in their neighborhoods

$$N(a^o, r) = \{c_i^o\} \text{ such that } \forall i \, d(a^o, c_i^o) \leq r. \quad (1)$$

- a^o = stadium in WordNet ontology
- $r = 1$

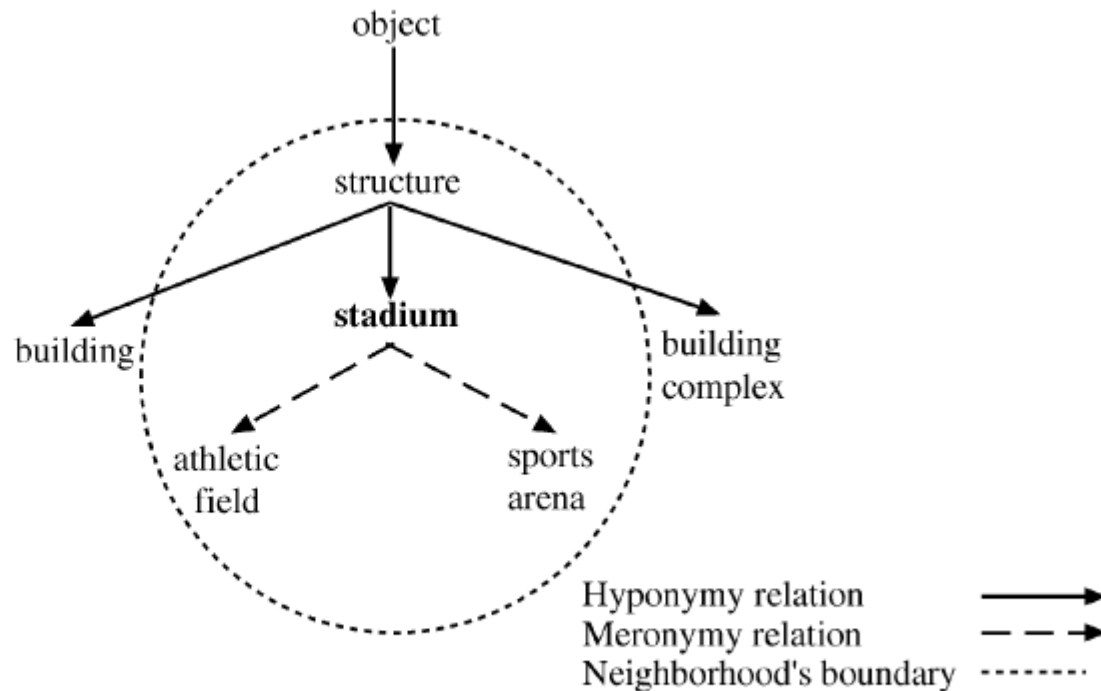


Fig. 1. Example of the immediate semantic neighborhood of *stadium* in a portion of the WordNet ontology.

- **Integrate the similarity assessments**

$$S(a^p, a^q) = \omega_w \cdot S_w(a^p, b^q) + \omega_u \cdot S_u(a^p, b^q) + \omega_n \cdot S_n(a^p, b^q),$$

for $\omega_w, \omega_u, \omega_n \geq 0$.

(2)

- $S()$: the similarity between entity class a in ontology p and b in ontology q .
- S_w : the similarity of synonym sets
- S_u : the similarity of features
- S_n : the similarity of semantic neighborhoods
- ω_w : the weight of the similarity of synonym sets component
- ω_u : the weight of the similarity of features component
- ω_n : the weight of the similarity of semantic neighborhoods component

- A Matching Model to Similarity Assessment

$$S(a, b) = \frac{|A \cap B|}{|A \cap B| + \alpha(a, b)|A/B| + (1 - \alpha(a, b))|B/A|}, \quad (3)$$

for $0 \leq \alpha \leq 1$.

- $S(a, b)$: the similarity of entity classes a and b calculated by matching model
- A, B : description sets of a and b , i.e. synonym sets, sets of distinguishing features, and sets of entity classes in the semantic neighborhood
- $A \cap B$: intersection of A and B
- A / B : difference of A and B
- $| \cdot |$: the cardinality of a set
- $\alpha(a, b)$: defines the relative importance of the noncommon characteristics
- a : referred to as the **target**
- b : referred to as the **base**
- the matching model is not forced to satisfy metric properties
 - e.g. $S(\text{building}, \text{office_building}) \neq S(\text{office_building}, \text{building})$

- A Matching Model (2)

$$\alpha(a^p, b^q) = \begin{cases} \frac{\text{depth}(a^p)}{\text{depth}(a^p) + \text{depth}(b^q)} & \text{depth}(a^p) \leq \text{depth}(b^q) \\ 1 - \frac{\text{depth}(a^p)}{\text{depth}(a^p) + \text{depth}(b^q)} & \text{depth}(a^p) > \text{depth}(b^q). \end{cases} \quad (4)$$

- In a cross-ontology, two independent ontologies are connected by making each of their roots a direct descendant of an imaginary and more general entity class, called *anything*.
- $\text{depth}()$: the shortest path from the entity class to the imaginary root
 - reflects the degree of granularity upon which the ontology was designed
- value of α are greater than 0 and less than or equal to 0.5

- A Matching Approach (3)

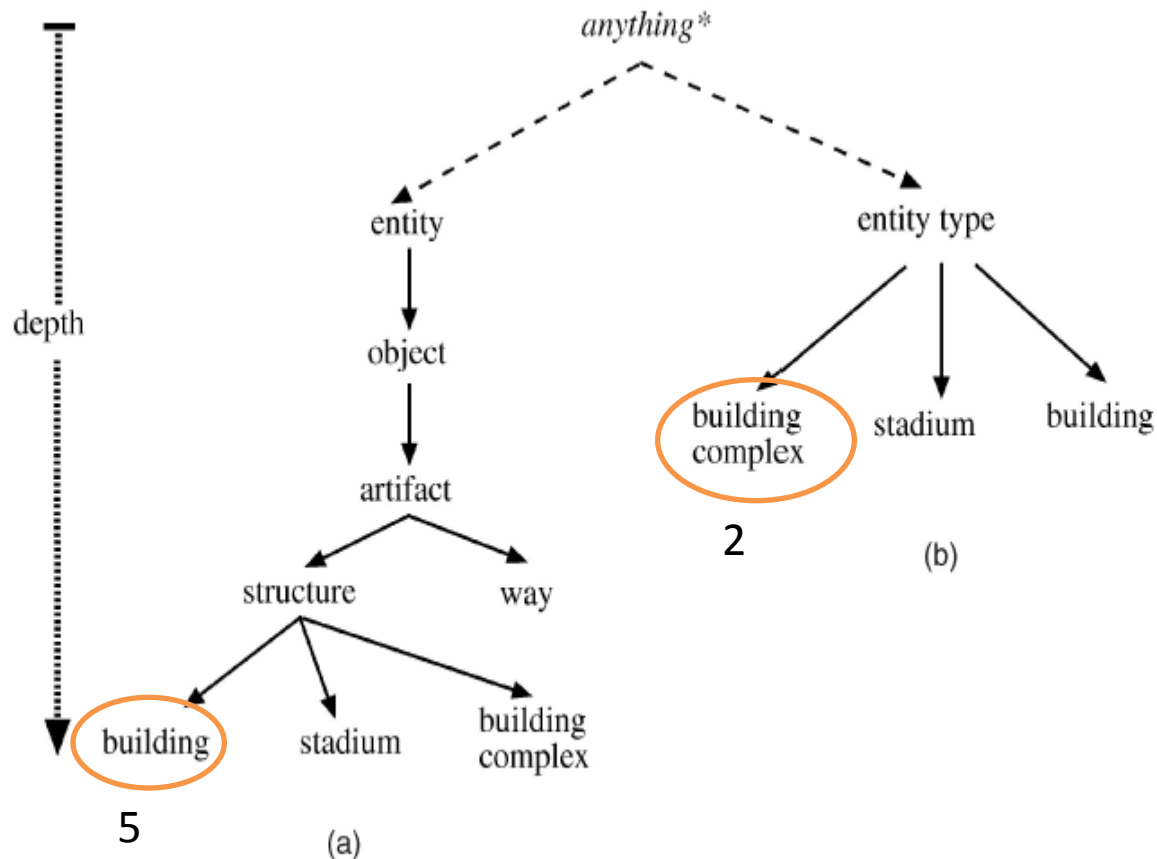


Fig. 2. Connecting independent ontologies: (a) partial WordNet ontology and (b) partial SDTS ontology. (*Anything** denotes an imaginary root.)

- Word Matching

$$S_t(\textit{building}^w, \textit{building_complex}^s) = \frac{|\{\textit{building}\}|}{|\{\textit{building}\}| + 0.28|\{\}\| + 0.72|\{\textit{complex}\}|} = \frac{1}{1.72} = 0.58.$$

2/(5+2)

(5)

- **Feature Matching**

$$S_u(a^p, b^q) = \omega_p \cdot S_p(a^p, b^q) + \omega_f \cdot S_f(a^p, b^q) + \omega_a \cdot S_a(a^p, b^q)$$

for ω_p, ω_f , and $\omega_a \geq 0$ and $\omega_p + \omega_f + \omega_a = 1.0$.

(6)

- $S_u()$: the similarity between entity class a in ontology p and b in ontology q .
- S_p, S_f, S_a : the similarity of parts, functions, attributes
- $\omega_p, \omega_f, \omega_a$: the weight of the similarity of each specification component
- represented by a synonym set
- **strict string matching**, only if represented by the same word or by synonym sets that intersect
- two distinguishing features are equivalent if the intersection of their synonym sets is not empty

TABLE 2
Entity_Class Definition of *Stadium* in WS and WordNet

Stadium (WS)	Stadium (WordNet)
entity_class { name: {stadium,bowl,arena} description: large often unroofed structure in which athletic events are held is_a: { <i>construction</i> *} part_of: {} whole_of: { <i>athletic_field</i> *} parts: { {athletic_field,sports_field,playing_field}, {dressing_room},{foundation}, {midfield},{spectator_stands,stands}, {ticket_office,box_office,ticket_booth}} functions: { {play,compete},{play,practise}, {recreate,play}} attributes: { {architectural_property}, {covered/uncovered},{name}, {lighted/unlighted},{owner_type}, {sports_type},{user_type}} }	entity_class { name: {stadium,bowl,arena} description: large often unroofed structure in which athletic events are held is_a: { <i>construction</i> *} part_of: {} whole_of: { <i>athletic_field</i> *, <i>sports_arena</i> *} parts: { {athletic_field,sports_field,playing_field}, {foundation},{midfield},{plate}, {sports_arena,field_house},{stands}, {structural_elements}, {standing_room},{tiered_seats}} functions: {} attributes: {} }

(*x** denotes a pointer to the entity class *x*)

$$\begin{aligned}
 X &= stadium^{ws}.parts \cap stadium^w.parts \\
 &= \{\{athletic_field, playing_field, field\}, \{foundation\}, \\
 &\quad \{midfield\}, \{stands\}\}.
 \end{aligned}
 \tag{7}$$

$$\begin{aligned}
 Y &= stadium^w.parts - stadium^{ws}.parts \\
 &= \{\{plate\}, \{sports_area, field_house\}, \\
 &\quad \{standing_room\}, \{structural_elements\}, \\
 &\quad \{tiered_seats\}\}
 \end{aligned}
 \tag{8a}$$

$$\begin{aligned}
 Z &= stadium^{ws}.parts - stadium^w.parts \\
 &= \{\{dressing_room\}, \\
 &\quad \{ticket_office, box_office, ticket_booth\}\}.
 \end{aligned}
 \tag{8b}$$

$$\begin{aligned}
 S_u(stadium^w, stadium^{ws}) &= S_p(stadium^w, stadium^{ws}) \\
 &= \frac{|X|}{|X| + 0.45|Y| + 0.55|Z|} \\
 &= \frac{4}{4 + 0.45 * 5 + 0.55 * 2} = 0.54.
 \end{aligned}
 \tag{9}$$

- **Semantic-Neighborhood Matching**
 - compare entity classes in semantic neighborhoods based on **word** or **feature matching**

$$S_n(a^p, b^q, r) = \frac{|a^p \cap_n b^q|}{|a^p \cap_n b^q| + \alpha(a^p, b^q) \cdot \delta(a^p, a^p \cap_n b^q, r) + (1 - \alpha(a^p, b^q)) \cdot \delta(b^q, a^p \cap_n b^q, r)}$$

with $\delta(a^p, a^p \cap_n b^q, r) =$

$$\begin{cases} |N(a^p, r)| - |a^p \cap_n b^q| & \text{if } |N(a^p, r)| > |a^p \cap_n b^q| \\ 0 & \text{otherwise.} \end{cases} \quad (10)$$

• Numerator

- $|a^p \cap_n b^q|$ is the approximate cardinality of the set intersection between these semantic neighborhoods
- a_i^p and b_j^q are entity classes in the semantic neighborhood of a^p and b^q
- n and m are the numbers of entity classes in the corresponding semantic neighborhoods

$$|a^p \cap_n b^q| = \left[\sum_{i \leq n} \max_{j \leq m} S(a_i^p, b_j^q) \right] - \varphi S(a^p, b^q), \text{ where}$$

$$\varphi = \begin{cases} 1 & \text{if } S(a^p, b^q) = \max_{j \leq m} S(a^p, b_j^q) \\ & \text{and} \\ 0 & \text{otherwise} \end{cases}$$

$$s(a_i^p, b_j^q) = \omega'_l S_l(a_i^p, b_j^q) + \omega'_u S_u(a_i^p, b_j^q) \text{ with } 0 < \omega'_l + \omega'_u \leq 1.$$

(11)

- **Cross-Ontology Evaluation**

- Address the quality of results of similarity assessments
- Design new experiments
 - have this matching model for similarity evaluations
 - use available ontologies, WordNet and SDTS
 - use human-subject testing
- Create a new ontology WS from the combination of WordNet and SDTS to exploit a more complete definition of entity classes
- **Two type of experiments**
 - searching for equivalent or most similar entity classes across ontologies
 - useful for ontology integration
 - ranking the similarity between entity classes in two ontologies
 - useful for information retrieval
 - e.g. stadium and athletic field

TABLE 3
Characteristics of the Specification Componenets of SDTS, WordNet, and WS

Characteristics	SDTS	WordNet	WS
Words			
Synonymy		√	√
Polysemy	√	√	√
Relations			
Is-a	√	√	√
Part-of		√	√
Whole-of		√	√
Features			
Parts		√	√
Functions			√
Attributes	√		√

- Combinations:
 - identical ontologies (1-2)
 - ontology and sub ontology (3)
 - overlapping ontologies (4)
 - different ontologies (5)
- Assessment measures:
 - use **recall** and **precision**
 - what entity classes are in fact similar?
 - equivalent or most similar entity classes, e.g. building vs. (building, building complex)

TABLE 4
Cases of Cross-Ontology Evaluations

Case	Ontology-Ontology	Description
1	WordNet-WordNet	Same ontology with is-a and part-whole relations
2	SDTS-SDTS	Same ontology with is-a relations and attributes
3	WordNet-WordNet*	Subset with same specification components
4	WordNet*-WS	Overlapping semantic relations and attributes
5	WordNet*-SDTS*	Different ontologies and specification components

(Symbol * denotes small subsets of the initial ontology)

TABLE 5
Recall and Precision of Evaluations with Threshold Equal to 75 Percent

Case	Weights (%)			Recall (%)	Precision (%)
	ω_w	ω_s	ω_r		
WordNet-WordNet	50	0	50	100	97
WordNet-WordNet	0	100	0	48	10
SDTS-SDTS	50	0	50	100	100
SDTS-SDTS	0	0	100	100	1
WordNet-WordNet*	50	0	50	99	98
WordNet-WordNet*	0	50	50	28	14
WordNet*-WS	100	0	0	100	78
WordNet*-WS	50	0	50	55	98
WordNet*-WS	0	50	50	0	0
WordNet*-SDTS*	100	0	0	100	42
WordNet*-SDTS*	50	0	50	50	92
WordNet*-SDTS*	0	100	0	0	0

(Symbol * denotes small subsets of the entire ontology)

- Conclusions:
 - ontologies share more components, the model produces more accurate results
 - the synonym sets and semantic neighborhood of entity classes is more similar across ontologies than the distinguishing features

- Comparison:
 - an entity class in an ontology with a reduced set of entity classes defined in another ontology.
 - e.g. stadium vs. stadium, athletic field, ballpark, tennis court, commons, building, theater, museum, library, transportation system, house, and sport arena.
- Human-subject test:
 - ask subjects to rank the similarity among the set of entity classes based on the definitions in the WS ontology
 - assume a number of ranks equal to the number of entity classes compared
 - Subjects found that the most similar entity classes to a stadium in decreasing order were: sports arena, ball park, athletic field, tennis court, theater museum, building, commons, library, house and transportation
- Type of Weight setting:
 - $\omega_w, \omega_u, \omega_n$: (0.33, 0.33, 0.33), (0.5, 0, 0.5), (0, 100, 0)

- The ordering of entity classes along the horizontal axis corresponds to the subjects' responses in decreasing order
- SDTS-WS:

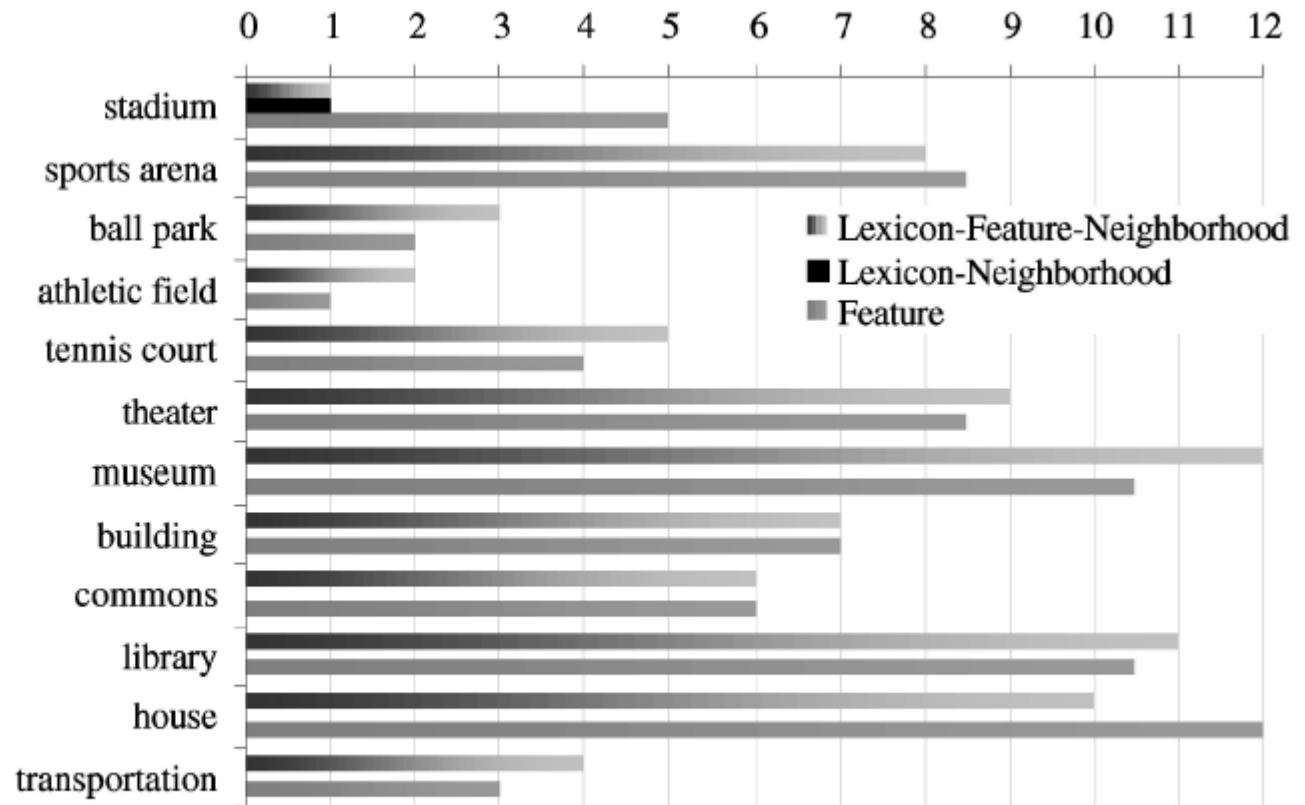


Fig. 3. Evaluations between the definition of *stadium* in SDTS and a set of entity classes defined in WS.

- WordNet-WS:

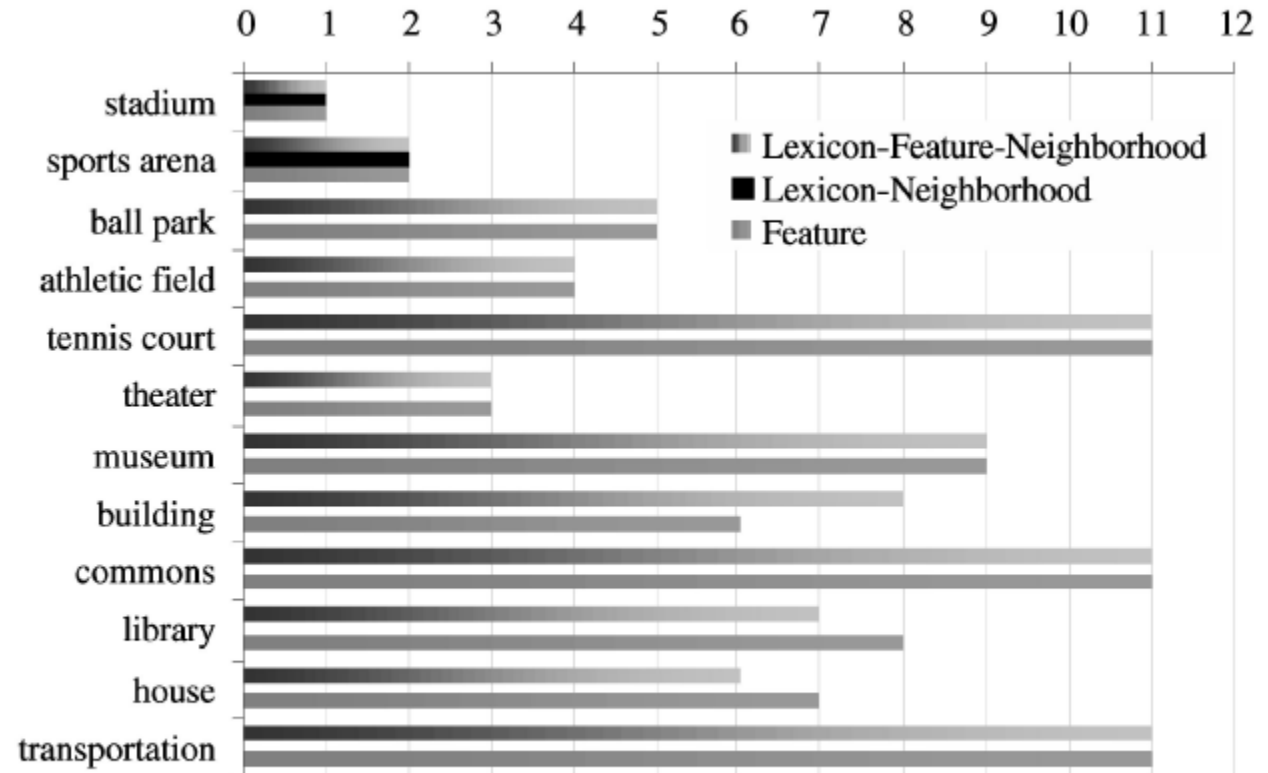


Fig. 4. Evaluations between the definition of *stadium* in WordNet and a set of entity classes defined in WS.

- WS-WS:

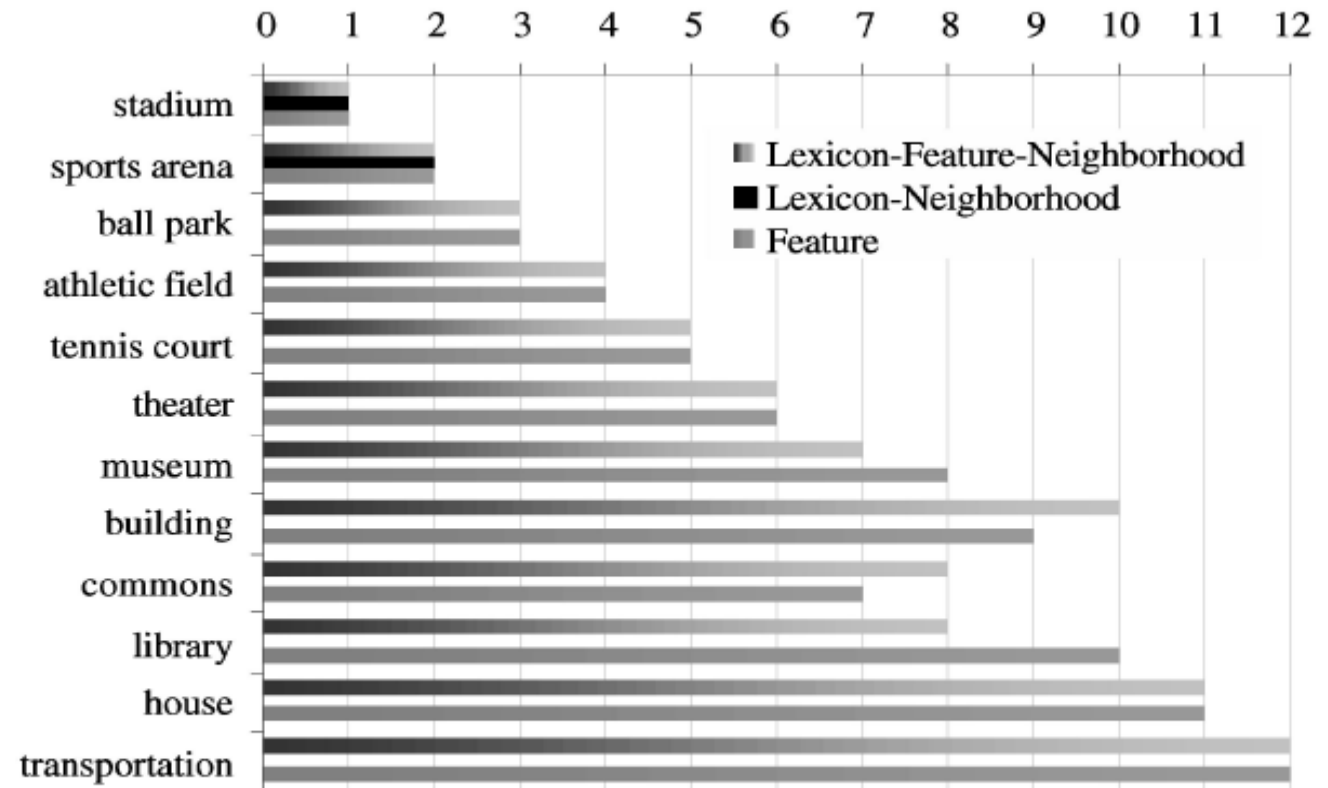


Fig. 5. Evaluations between the definition of *stadium* in WS and a set of entity classes defined in WS.

TABLE 6
Correlation Coefficient for Similarity Ranks in Cross-Ontology Evaluations

Ontologies	Word-Feature-Neighborhood $\omega_w: 33.3, \omega_u: 33.3, \omega_n: 33.3$	Word-Neighborhood $\omega_w: 50, \omega_n: 50$	Feature $\omega_u: 100$
SDTS-WS	0.48	-0.34	0.37
WordNet-WS	0.68	-0.34	0.71
WS-WS	0.96	-0.34	0.97

correlation coefficient

- Conclusions:
 - ontologies share more components, the model produces more accurate results
 - the best combination of weights detected in the former experiment give the worst value of correlation, BECAUSE detect the most similar entity class and nothing else
 - feature matching is important for detecting similar entity classes within an ontology or the similarity of semantically related entity classes across ontologies
- Assignment of weight:
 - ontology characteristics
 - the goal of the similarity assessment
 - i.e. ontology integration vs. information retrieval

- Similarity model:
 - based on the matching process of each component
 - useful as a first step in an ontology integration
- New insight:
 - **Synonym set** and **Semantic neighborhoods**
 - for detecting equivalent or most similar entity
 - **Distinguishing features**
 - for detecting entity classes that are somewhat similar
 - i.e. not synonyms and are located far apart in the hierarchical structure
 - e.g. stadium and athletic field

- Further work needed:
 - Now compare distinguishing feature in terms of a strict string matching between synonym sets refer to those feature, SO **semantic similarity** among features is needed.
 - Parts are also entity classes that could be semantically compared in a recursive process.
 - Verbs could be related by the semantic relation **entailment** (e.g. buy and pay)

Thank You

