

# Resolving XML Semantic Ambiguity

## Technical Report – XSDF-TR-2014

Nathalie Charbel

LE2I Lab. UMR-CNRS,  
University of Bourgogne (UB),  
21078 Dijon, France  
nathaliecharbel@gmail.com

Joe Tekli

ECE Dept., Lebanese  
American University (LAU),  
36 Byblos, Lebanon  
joe.tekli@lau.edu.lb

Richard Chbeir

LIUPPA Lab., University of Pau  
& Adour Countries (UPPA),  
64600 Anglet, France  
richard.chbeir@univ-pau.fr

Gilbert Tekli

NOBATEK,  
67 Rue de Mirambeau  
64600 Anglet, France  
gtekli@gmail.com

### ABSTRACT

XML semantic-aware processing has become one of the central issues in Web data management, data processing, and information retrieval. While XML data is semi-structured, yet it remains prone to lexical ambiguity, and thus requires dedicated semantic analysis and sense disambiguation processes to assign well-defined meaning to XML elements and attributes. This becomes crucial in an array of applications ranging over semantic-aware query rewriting, semantic document clustering and classification, schema matching, as well as blog analysis and event detection in social networks and tweets. Most existing approaches in this context: i) ignore the problem of identifying ambiguous XML nodes, ii) only partially consider their structural relations/context, iii) use syntactic information in processing XML data regardless of the semantics involved, and iv) are static in adopting fixed disambiguation constraints thus limiting user involvement. In this paper, we provide a new XML Semantic Disambiguation Framework titled *XSDF* designed to address each of the above motivations, taking as input: an XML document and a general purpose semantic network, and then producing as output a semantically augmented XML tree made of unambiguous semantic concepts. Experiments demonstrate the effectiveness of our approach in comparison with alternative methods.

This technical report contains proofs, computation examples, and experimental results adding to the contributions of the main paper.

### CONTENTS

I. Proofs of Propositions and Lemmas.....	1
II. Computation Examples.....	2
III. Disambiguation Algorithms.....	3
IV. Evaluating XML Node Ambiguity.....	3
V. Time Analysis.....	4

### I. Proofs of Propositions and Lemmas

**Proposition 1:** The ambiguity degree of an XML node  $x$  in tree  $T$  increases when the number of senses of  $x.\ell$  is high in the reference semantic network  $SN$ , or else it decreases such that:

$$Amb_{Polysemy}(x.\ell, SN) = \frac{senses(x.\ell) - 1}{Max(senses(SN)) - 1} \in [0,1] \quad (1)$$

where  $Max(senses(SN))$  is the maximum number of senses of a word/expression in  $SN$  (e.g., in WordNet 2.1 [14],  $Max_{polysemy} = 33$ , for the word “head”) □

---

**Proof of Prop. 1:** Based on formula 1,  $Amb_{Polysemy}$  varies as follows:

- The minimum value  $Amb_{Polysemy} = 0$  is obtained when  $x.\ell = 1$ , i.e.,  $\ell$  has only one sense (e.g., “first name” in WordNet), i.e.,  $x$  is without ambiguity: it always refers to the same meaning.
- The maximum value  $Amb_{Polysemy} = 1$  is obtained when  $x.\ell = Max(polysemy(SN))$ .
- When  $polysemy(x.\ell)$  increases/decreases,  $Amb_{Polysemy}$  follows accordingly such that  $Amb_{Polysemy} \in [0, 1]$  □

---

**Proposition 2:** The ambiguity degree of an XML node  $x$  in tree  $T$  increases when the distance in number of edges between  $x$  and  $R(T)$  is low, or else it decreases such that:

$$Amb_{Depth}(x, T) = 1 - \frac{x.d}{Max(depth(T))} \in [0,1] \quad (2)$$

where  $Max(depth(T))$  is the maximum depth in  $T$  □

---

**Proof of Prop. 2:** Following formula 2,  $Amb_{Depth}$  varies as follows:

- The maximum value  $Amb_{Depth} = 1$  is obtained when  $x.d = 0$ , i.e., when  $x = R(T)$ .
- The minimum value  $Amb_{Depth} = 0$  is obtained when  $x.d = Max(depth(T))$ , i.e., when  $x$  is one of the farthest nodes from  $R(T)$ : one of the deepest (most specific) leaf nodes in  $T$ 's hierarchy.
- When  $x.d$  increases/decreases,  $Amb_{Depth}$  follows inversely such that  $Amb_{Depth} \in [0, 1]$  □

---

**Proposition 3:** The ambiguity degree of an XML node  $x$  in tree  $T$  increases when the number of children nodes of  $x$  having distinct node labels, designated as  $\overline{x.f}$ , is low, or else it decreases:

$$Amb_{Density}(x, T) = 1 - \frac{\overline{x.f}}{Max(\overline{fan-out}(T))} \in [0,1] \quad (3)$$

where  $Max(\overline{fan-out}(T))$  is the maximum number of children nodes with distinct node labels in  $T$ . We identify this factor as node density factor to distinguish it from traditional node fan-out: number of children nodes (regardless of label, cf. **Error! Reference source not found.**) □

---

**Proof of Prop. 3:** Following formula 3,  $Amb_{Density}$  varies as follows:

- The maximum value  $Amb_{Density} = 1$  is obtained when  $\overline{x.f} = 0$ , i.e., when  $x$  is a leaf node and does not have any children nodes (to provide hints concerning  $x$ 's meaning).

- The maximum value  $Amb_{Density} = 0$  is obtained when  $\overline{x.f} = \text{Max}(\text{fan-out}(T))$ , i.e., when  $x$  has the largest number of children nodes with distinct labels in  $T$ . In other words, it has the highest possible number of hints regarding its meaning in  $T$ .
- When  $\overline{x.f}$  increases/decreases,  $Amb_{Density}$  follows inversely such that  $Amb_{Density} \in [0, 1] \square$

**Lemma 1:** The ambiguity degree measure  $Amb\_Deg$  in **Error! Reference source not found.** varies in accordance with Propositions 1-3, and conforms to Assumptions 1-4 (refer to main paper)  $\square$

**Proof of Lemma 1:** Following formula 4 in the main paper,  $Amb\_Deg$  varies as follows:

- The minimum value  $Amb\_Deg = 0$  is obtained when the label of node  $x$  has only one possible sense, i.e., when  $Amb_{Polysem}(x.l, SN) = 0$ , regardless of its depth and density factors, and regardless of parameter weights (Assumption 4).
- The maximum value  $Amb\_Deg = 1$  is obtained when: i)  $x.l$  has the maximum number of senses in  $SN$  (e.g., 33 senses in WordNet [16]), ii)  $x$  is the root node of  $T$ ,  $x.d = 0$ , and iii)  $x$  has the minimum number of children nodes with distinct labels in  $T$ ,  $\overline{x.f} = 0$ , regardless of  $w_{Depth}$  and  $w_{Density}$  parameter values (Assumptions 1-3, and Lemmas 1-3)
- The value of  $Amb\_Deg$  increases with: i) the increase in  $x.l$ 's polysemy in  $SN$ , ii) the decrease in  $x$ 's depth, and iii) the decrease in  $x$ 's density in the XDT. Inversely, the value of  $Amb\_Deg$  decreases with: i) the decrease in  $x.l$ 's polysemy, ii) the increase in  $x$ 's depth, and iii) the increase in  $x$ 's density (Assumptions 1-3, and Lemmas 1-3).

**Lemma 2:** The context vector weight measure  $w_{\overline{V_d(x)}}(l_r)$  in **Error! Reference source not found.** varies in accordance with Assumptions 5 and 6 (refer to main paper for definition and assumptions)  $\square$

**Proof of Lemma 2:** Based on formula 5 in the main paper,  $w_{\overline{V_d(x)}}(l_r)$  varies as follows:

- $w_{\overline{V_d(x)}}(l_r)$  increases with  $Freq(l_r, S_d(x))$ , which increases with  $Struct(x_i, S_d(x))$ :
  - The value of  $Struct(x_i, S_d(x))$  is inversely proportional to the distance between the target node  $x$  and context node  $x_i$  (following Assumption 1).
    - The minimum value  $Struct(x_i, S_d(x)) = \frac{1}{d+1}$  is reached when  $x_i \in R_d(x)$  where  $R_d(x)$  is the outer-most ring in sphere  $S_d(x)$ .
    - The maximum value  $Struct(x_i, S_d(x)) = 1$  is reached when processing the target node itself, i.e.,  $x_i = x$ .
  - The value of  $Freq(l_r, S_d(x))$  is proportional to the number of occurrences of nodes having the same label  $x_i.l = l_r$  (following Assumption 2).

- The minimum value  $Freq(l_r, S_d(x)) = \frac{1}{d+1}$  is reached when there is only one context node having label  $l_r$  and occurring on the outer-most ring  $R_d(x)$  of the sphere neighborhood, more formally:
$$\exists x_i \in S_d(x) \ / \ (x_i.l = l_r) \wedge (x_i \in R_d(x)) \wedge (\forall x_j \in S_d(x) / x_j.l \neq l_r)$$

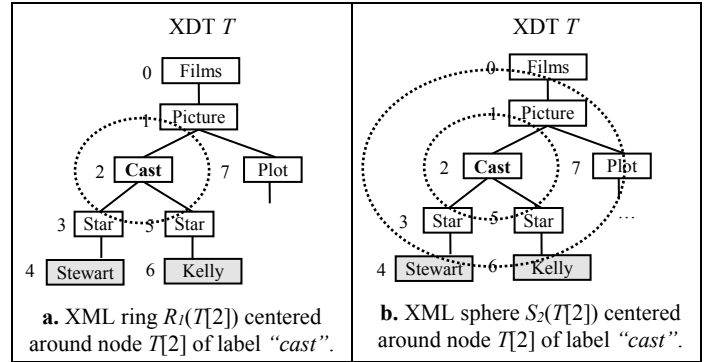
- The maximum value of  $Freq(l_r, S_d(x)) = \frac{|S_d(x)|+1}{2}$  is reached when all context nodes have the same label  $l_r$  and appear on the inner-most ring  $R_1(x)$  of the sphere neighborhood, more formally:
$$\forall x_i \in S_d(x) \ / \ (x_i.l = l_r) \wedge (x_i \in R_1(x)).$$
 Here,
$$Freq(l_r, S_d(x)) = 1 + \frac{1}{2} \times (|S_d(x)| - 1) = \frac{|S_d(x)| + 1}{2}$$

Consequently:

- The minimum value  $w_{\overline{V_d(x)}}(l_r) = 0$  is obtained when no nodes of label  $l_r$  occur in the sphere context of target node  $x$ .
- The maximum value  $w_{\overline{V_d(x)}}(l_r) = 1$  is obtained when maximum frequency is obtained, since the weight score is normalized using maximum frequency, i.e.,  $\frac{|S_d(x)|+1}{2} \square$

## II. Computation Examples

Consider XML document trees in Figure 1 (reported from the main paper), where the XML sphere  $S_2(T[2])$  is centered around node  $T[2]$  of label "cast" with radius 2 consists of: ring  $R_1(T[2])$  of radius 1 comprises nodes  $T[1]$  ("picture"),  $T[3]$  ("star") and  $T[5]$  ("star"), and ring  $R_2(T[2])$  of radius 2 comprises nodes  $T[0]$  ("Films"),  $T[4]$  ("Stewart"),  $T[6]$  ("Kelly"), and  $T[7]$  ("Plot").



**Figure 1.** Sample XML (ring and) sphere neighborhoods.

	Cast	Picture	Star
$\overline{V_1(T[2])}$	0.4	0.2	0.4

	Cast	Picture	Star	Films	Stewart	Kelly	Plot
$\overline{V_2(T[2])}$	0.25	0.1667	0.3334	0.0835	0.0835	0.0835	0.0835

**Figure 2.** Sample sphere context vectors based on the sphere neighborhoods in Figure 1.

Considering  $\overline{V_1(T[2])}$  in Figure 2:

- $w_{\overline{V_1("Cast")}}("Cast") = \frac{2 \times (1)}{(4) + 1} = 0.4$  given that: i) "Cast" is the label of  $S_l(T[2])$ 's target (center) node  $T[2]$ , i.e.,  $Struct(T[2], S_l(T[2])) = 1$ , ii)  $T[2]$  is the only node occurrence of label "cast" in  $S_l(T[2])$ , i.e.,  $Freq("Cast", S_l(T[2])) = 1$ , and iii)  $|S_l(T[2])| = 4$ .
- $w_{\overline{V_1("Picture")}}("Picture") = \frac{2 \times (0.5)}{(4) + 1} = 0.2$  given: i)  $S_l(T[2])$  contains one node  $T[1]$  of label "Picture" positioned at distance = 1 from the target node  $T[2]$ , i.e.,  $Struct(T[3], S_l(T[2])) = 1 - \frac{1}{2} = 0.5$ , ii)  $T[2]$  is the only node occurrence of label "Picture" in  $S_l(T[2])$ , i.e.,  $Freq("Picture", S_l(T[2])) = 0.5$ , and iii)  $|S_l(T[2])| = 4$ .

Likewise, considering vector  $\overline{V_2(T[2])}$ :  $w_{\overline{V_2("Star")}}("Star") = \frac{2 \times (0.6667 + 0.6667)}{(7) + 1} = 0.3335$  given: i)  $S_l(T[2])$  contains two nodes  $T[3]$  and  $T[5]$  of label "Star" at distance = 1 from the target node  $T[2]$ , i.e.,  $Struct(T[3], S_l(T[2])) = 1 - \frac{1}{3} = 0.6667$ , ii)  $T[3]$  and  $T[5]$  are the only node occurrences of label "Star"  $S_l(T[2])$ , i.e.,  $Freq("Star", S_l(T[2])) = 0.6667 + 0.6667 = 1.3334$  and iii)  $|S_l(T[2])| = 7$ .

### III. Disambiguation Algorithms

Algorithm $XSD_{Concept}$	
<b>Input:</b> $x, d$	// Target XML node $x$ and sphere neighborhood radius $d$
$T$	// Source XML tree
$SN$	// Reference weighted semantic network
<b>Output:</b> $c$	// Semantic concept representing the sense (meaning) of $x, \ell$
Begin	1
Generate $S_d(x) \in T$ and $\overline{V_d(x)}$	// Context vector of $x$ in $T$
For each $s_p \in \overline{SN} \rightarrow x, \ell$	// For each sense of the target node label
{	4
For each $x_i \in S_d(x)$	// Processing senses of context nodes
{	6
For each $s_j^i \in \overline{SN} \rightarrow x_i, \ell$	{ $Sim\_Score(s_j^i) = Sim(s_p, s_j^i)$ }
$Max\_Score(x_i, \ell) = \text{Max}_{s_j^i \rightarrow x_i, \ell} (Sim\_Score(s_j^i))$	8
}	9
$Concept\_Score(s_p) = \sum_{x_i \in S_d(x)} \frac{Max\_Score(x_i, \ell) \times w_{\overline{V_d(x)}}(x_i, \ell)}{ S_d(x) }$	10
}	11
$c = s_k / Concept\_Score(s_k) = \text{Max}_{s_r \rightarrow x, \ell} (Concept\_Score(s_p))$	12
Return $c$	// Sense (concept) with maximum score
End	14

Figure 3. Algorithm  $XSD_{Concept}$  for concept-based XML Semantic Disambiguation.

Algorithm $XSD_{Context}$	
<b>Input:</b> $x, d$	// Target XML node $x$ and sphere neighborhood radius $d$
$T$	// Source XML tree
$SN$	// Reference semantic network
<b>Output:</b> $c$	// Semantic concept representing the sense (meaning) of $x, \ell$
Begin	1
Generate $S_d(x) \in T$ and $\overline{V_d(x)}$	// Context vector of $x$ in $T$
For each $s_p \in \overline{SN} \rightarrow x, \ell$	// For each sense of the target node label
{	4
Generate $S_d(s_p) \in SN$ and $\overline{V_d(s_p)}$	// Context vector of $s_p$ in $SN$
$Context\_Score(s_p) = \frac{\overline{V_d(x)} \cdot \overline{V_d(s_p)}}{\  \overline{V_d(x)} \  \times \  \overline{V_d(s_p)} \ }$	// Context similarity
}	7
$c = s_k / Context\_Score(s_k) = \text{Max}_{s_r \rightarrow x, \ell} (Context\_Score(s_p))$	8
Return $c$	// Sense (concept) with maximum score
End	10

Figure 4. Algorithm  $XSD_{Context}$  for context-based XML Semantic Disambiguation.

### IV. Evaluating XML Node Ambiguity

Detailed manual and system ambiguity ratings concerning the three extreme correlations scores (*maximum*, closest to *null*, and *minimum* scores) highlighted in bold in Table 1 (report from the main paper), are shown in Figures 5-8 respectively.

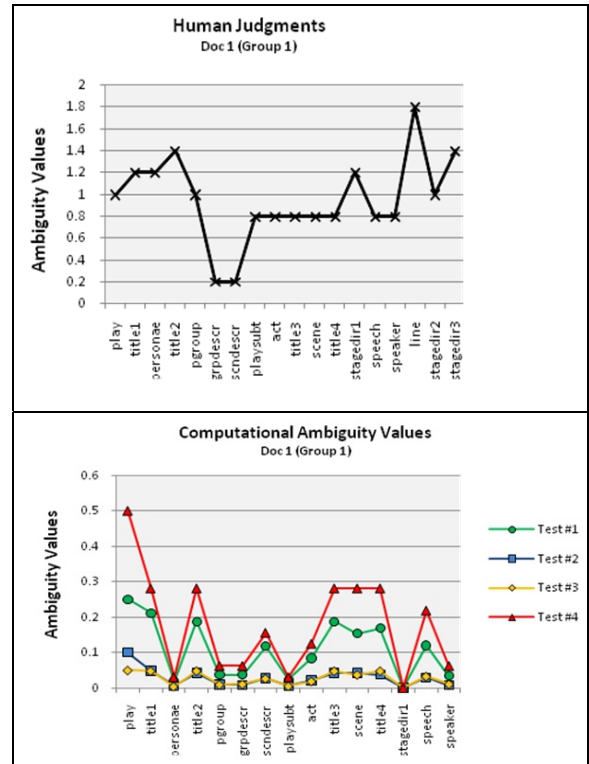
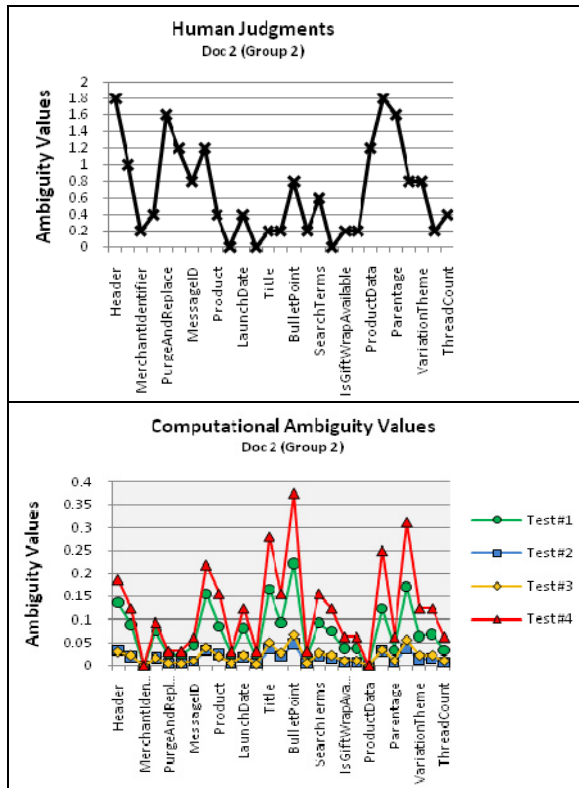


Figure 5. Manual and system generated average ambiguity degrees highlighting maximum correlation with documents of data-set 1 of Group 1. The x axis represents node labels (tag names and/or values) being disambiguated.

**Table 1.** Correlation between human ratings and system generated ambiguity degrees (reported from the main paper).

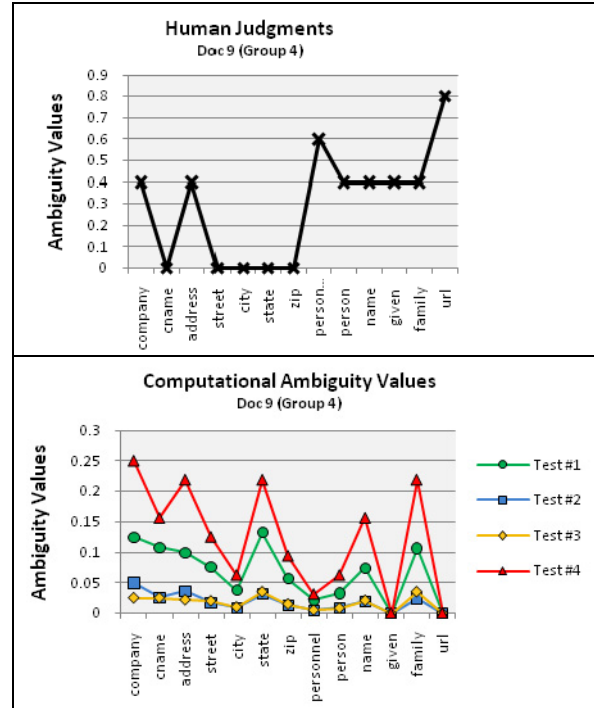
		Test #1 <i>All factors</i>	Test #2 <i>Polysemy</i>	Test #3 <i>Depth</i>	Test #4 <i>Density</i>
<b>Group 1</b>	Data-set 1	0.394	0.411	0.335	<b>0.439</b>
<b>Group 2</b>	Data-set 2	<b>0.017</b>	0.181	0.243	0.139
<b>Group 3</b>	Dataset 3	-0.087	-0.139	-0.071	-0.138
	Data-set 4	0.408	0.438	0.390	0.398
	Data-set 5	-0.184	-0.185	-0.131	-0.235
<b>Group 4</b>	Data-set 6	-0.284	-0.291	-0.243	-0.316
	Data-set 7	-0.177	-0.190	-0.254	-0.143
	Data-set 8	-0.119	-0.025	0.033	-0.156
	Data-set 9	-0.452	-0.301	-0.251	<b>-0.456</b>
	Data-set 10	-0.258	0.180	0.412	0.276



**Figure 6.** Manual and system generated average ambiguity degrees highlighting closest to zero (null) correlation with documents of data-set 2 of Group 2.

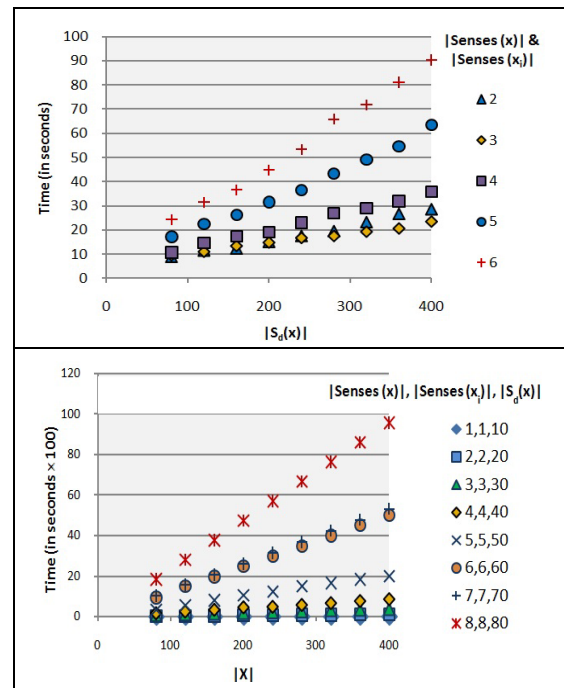
<pre>&lt;?xml version="1.0"?&gt; &lt;company&gt;   &lt;cname id="id3_4"&gt;Informix&lt;/cname&gt;   &lt;address&gt;     &lt;street&gt;123 6th Ave W&lt;/street&gt;     &lt;city&gt;Portland&lt;/city&gt;     &lt;state&gt;OR&lt;/state&gt;     &lt;zip&gt;54678&lt;/zip&gt;   &lt;/address&gt;   &lt;personnel&gt;     &lt;person&gt;       &lt;name&gt;         &lt;given&gt;Fran&lt;/given&gt;         &lt;family&gt;Car&lt;/family&gt;       &lt;/name&gt;       &lt;url&gt;http://null&lt;/url&gt;     &lt;/person&gt;   &lt;/personnel&gt; &lt;/company&gt;</pre>	<pre>&lt;?xml encoding="ISO-8859-1"?&gt; &lt;!ELEMENT company (address, cname,   personnel)&gt; &lt;ATTLIST company id ID #REQUIRED&gt; &lt;!ELEMENT address (street, city, state, zip)&gt; &lt;!ELEMENT personnel (person)&gt; &lt;!ELEMENT person (name, email?, url?)&gt; &lt;!ELEMENT family (#PCDATA)&gt; &lt;!ELEMENT given (#PCDATA)&gt; &lt;!ELEMENT name (family?given?)&gt; &lt;!ELEMENT cname (#PCDATA)&gt; &lt;!ELEMENT email (#PCDATA)&gt; &lt;!ELEMENT street (#PCDATA)&gt; &lt;!ELEMENT city (#PCDATA)&gt; &lt;!ELEMENT state (#PCDATA)&gt; &lt;!ELEMENT zip (#PCDATA)&gt; &lt;!ELEMENT url (#PCDATA)&gt;</pre>
--	--

**Figure 7.** Sample XML document from data-set 9 of Group 4, with corresponding grammar.



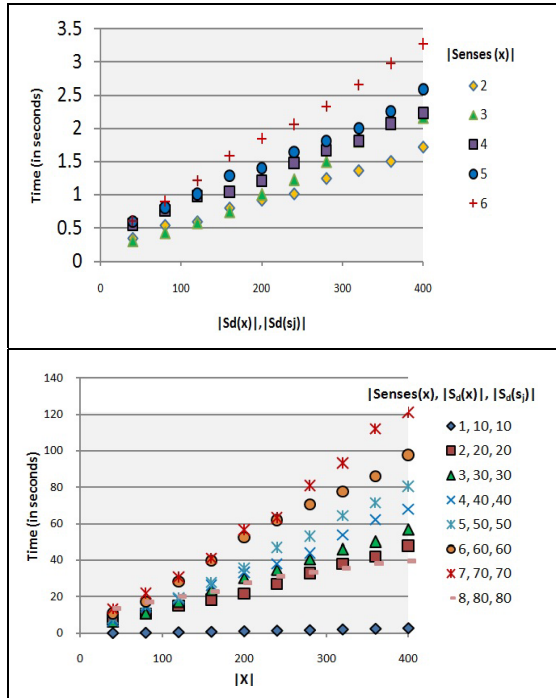
**Figure 8.** Manual and system generated average ambiguity degrees highlighting minimum correlation with documents of data-set 9 of Group 4 (Figure 7). The x axis represents node labels (tag names and/or data values).

## V. Time Analysis



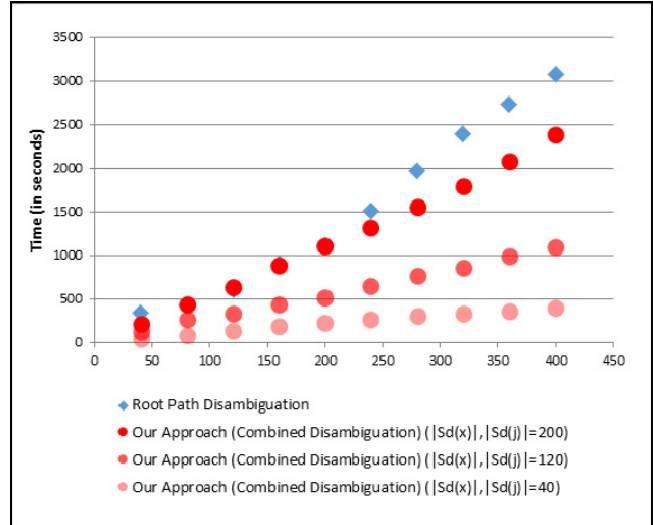
**Figure 9.** Timing results regarding our *concept-based* approach (when executing algorithm *XSDConcepts*, cf. Figure 3).

We evaluated the efficiency of our approach in terms of execution time. Results highlight the linear complexities of both our concept-based approach i.e.,  $O(|senses(x,\ell)| \times |S_d(x)| \times |senses(x_i,\ell)|)$  (Figure 9) and our context-based approach, i.e.,  $O(|senses(x,\ell)| \times (|S_d(x)| + |S_d(s_p)|))$  (Figure 10). Time is also linear w.r.t. the number of target nodes being disambiguate, designated as  $|X|$ . The complete experimental study (along with test documents and the detailed properties of XML nodes being disambiguated, e.g., label polysemy, node depth, density, etc.) is available online<sup>1</sup>.



**Figure 10.** Timing results regarding our *context-based* approach (when executing algorithm  $XSD_{Context}$ , cf. Figure 4).

We have also compared the time complexity of our approach using different configurations, with one of its most recent predecessors. Results in Figure 11 show closely correlated and even reduced time results (depending on the configuration used, e.g., smaller context radiuses in XML document and/or in the reference semantic network). This means that our approach was able to produce improved disambiguation quality while preserving (and sometimes reducing) execution time levels in comparison with its alternatives<sup>2</sup>.



**Figure 11.** Comparing time results with existing *Root Path Disambiguation* [1] approach.

## REFERENCES

- [1] Tagarelli A. *et al.*, *Word Sense Disambiguation for XML Structure Feature Generation*. In Proc. of ESWC, 2009, pp. 143–157.
- [2] Mandreoli F. *et al.*, *Versatile Structural Disambiguation for Semantic-Aware Applications*. In Proc. of Inter. CIKM Conf., 2005. pp. 209-216.

<sup>1</sup> <http://sigappfr.acm.org/Projects/XSDF/>

<sup>2</sup> Note that we did not compare execution time with the *Versatile Structure Disambiguation* approach [2] since we were unable to acquire the system implementation from the authors. We used the authors’ online version of the prototype, which is relatively slow due to network access, and thus could not use it to evaluate processing time.