

CPENG-1952.pdf

By H. David Jeong

WORD COUNT

17708

TIME SUBMITTED

12-SEP-2016 01:22AM

PAPER ID

24909047

Journal of Computing in Civil Engineering

NLP-based approach to classifying heterogeneous terms for unambiguous exchange of roadway data --Manuscript Draft--

Manuscript Number:	CPENG-1952R1
Full Title:	NLP-based approach to classifying heterogeneous terms for unambiguous exchange of roadway data
Article Type:	Technical Paper 2
Abstract:	<p>The inconsistency of data terminology due to the fragmented nature of the highway industry has imposed big challenges on integrating digital data from distinct sources. The issue of semantic heterogeneity may lead to the lack of common understanding of the same data between the sender and receiver. Explicit semantic relations among terms in digital dictionaries, such as ontologies can enable the meaning of a roadway concept name to be transparent and unambiguously understood by computer systems. However, due to the lack of an effective automated method, current practices of identifying these relations hardly rely on a manual process of knowledge acquisition from domain experts or text documents which is laborious and time-consuming. This paper presents a novel methodology that leverages recent advances in Natural Language Processing (NLP) techniques to extract English-American roadway terms used in different government agencies and their semantic relations from roadway design manuals and specifications. The proposed method includes the following three stages: (1) implementing NLP techniques to detect commonly used technical terms from the highway corpus; (2) utilizing machine learning to learn the semantic similarity among roadway terms using their context data in the corpus; and (3) developing a classification algorithm to identify semantic relation types among technical terms. The key merit in this technique is the automated identification of semantic relations among heterogeneous roadway terms from design guidebooks without reliance on other existing hand-coded semantic resources. The proposed methodology was evaluated by conducting an experiment comparing the automatically-identified synonyms by the proposed system with a human-constructed golden standard dataset obtained from Wikipedia. The result shows that the proposed model achieves a precision of over 80 percent.</p>

1 **NLP-based approach to classifying heterogeneous terms
2 for unambiguous exchange of roadway data**

3 Tuyen Le¹, H. David Jeong²

4 **ABSTRACT**

5 The inconsistency of data terminology due to the fragmented nature of the highway in-
6 dustry has imposed big challenges on integrating digital data from distinct sources. The
7 issue of semantic heterogeneity may lead to the lack of common understanding of the same
8 data between the sender and receiver. Explicit semantic relations among terms in dig-
9 ital dictionaries, such as ontologies can enable the meaning of a roadway concept name
10 to be transparent and unambiguously understood by computer systems. However, due to
11 the lack of an effective automated method, current practices of identifying these relations
12 hardly rely on a manual process of knowledge acquisition from domain experts or text doc-
13 uments which is laborious and time-consuming. This paper presents a novel methodology
14 that leverages recent advances in Natural Language Processing (NLP) techniques to extract
15 English-American roadway terms used in different government agencies and their semantic
16 relations from roadway design manuals and specifications. The proposed method includes
17 the following three stages: (1) implementing NLP techniques to detect commonly used tech-
18 nical terms from the highway corpus; (2) utilizing machine learning to learn the semantic
19 similarity among roadway terms using their context data in the corpus; and (3) developing a
20 classification algorithm to identify semantic relation types among technical terms. The key
21 merit in this technique is the automated identification of semantic relations among heteroge-
22 neous roadway terms from design guidebooks without reliance on other existing hand-coded

2

¹Ph.D. Candidate, Department of Civil, Construction and Environmental Engineering, Iowa State University. Ames, IA 50011, United States. E-mail: ttle@iastate.edu.

²Associate Professor, Department of Civil, Construction and Environmental Engineering, Iowa State University. Ames, IA 50011, United States. E-mail: djeong@iastate.edu.

semantic resources. The proposed methodology was evaluated by conducting an experiment comparing the automatically-identified synonyms by the proposed system with a human-constructed golden standard dataset obtained from Wikipedia. The result shows that the proposed model achieves a precision of over 80 percent.

Keywords: Roadway Data, Data Sharing, Semantic Interoperability, Semantic Relation, Natural Language Processing, Vector Space Model

INTRODUCTION

The implementation of advanced, and computerized technologies such as 3D modeling and Geographic Information System (GIS) throughout the life cycle of a highway project has allowed a large portion of project data to be available in a digital format. The efficiency improvement in sharing these data between project participants and stages, will in turn, translate into increased productivity, efficiency in project delivery and accountability. The interoperability issue has been widely recognized as a key obstacle blocking the flow of digital data throughout the entire project life cycle. The inadequate interoperability cost is estimated to be over \$15.8 billion per year in the U.S. capital facilities industry as reported by the National Institute of Standard and Technology (NIST); and the largest cost item is the laborious work for finding, verifying, and transferring facility and project information into a useful format during the operation and maintenance stage (Gallaher et al. 2004). This finding indicates that the lack of readiness for downstream phases to directly use the transferred digital project data generated from upstream design and construction stages results in high operational costs. Since the roadway sector, which is one of the major domains in the construction industry, has not yet successfully facilitated a high degree of interoperability (Lefler 2014); huge cost savings would be achieved if roadway data is seamlessly shared across project phases and among state and local agencies.

Semantic interoperability, which relates to the issue whereby two computer systems may not share a common understanding of a specific piece of data, is a radical barrier to computer-to-computer data exchange. Due to the fragmented nature of the infrastructure domain,

50 data representation/terminology differs between phases, stakeholders, or geographic regions
51 (counties, states, etc.). Retrieving right pieces of data in such a heterogeneous environment
52 becomes increasingly complex (Karimi et al. 2003). Polysemy and synonymy are two major
53 linguistic obstacles to semantic integration and use of a multitude of data sources (Noy
54 2004). Polysemy refers to cases when a unique term has several distinct meanings. For
55 example, *roadway type* can either mean the classification of roadways by material or function.
56 Walton et al. (2015) suggests the following three reasons for semantic heterogeneity among
57 transportation databases: (1) isolation in definitions among separate sources, (2) temporary
58 of definitions and (3) variety of data collection methods. Synonymy, in contrast, is associated
59 with a set of different terms used to present the same concept. For instance, ‘profile’, ‘crest’,
60 ‘grade-line’ and ‘vertical alignment’ are equivalent terms of the *longitudinal centerline* of a
61 roadway. Under these situations, simply mapping of data names will likely lead to a failure
62 of data extraction, or use of wrong data. Thus, addressing the semantic inconsistency issue
63 30 becomes crucial to ensure a common understanding of the same dataset among software
64 applications and guarantee a proper integration of data from multiple sources.

65 Terminology transparency through digital dictionaries like glossaries, taxonomies, ontolo-
66 gies and data dictionaries is identified as a driver of semantic interoperability (Ouksel and
67 Sheth 1999). Although a plethora of semantic resources have been introduced to the highway
68 sector; as shown in the literature review, their coverages of concepts are still inadequate and
69 the inclusion of multiple names for the same concept is limited. This is because of the reliance
70 on a tedious and time-consuming approach which requires developers to manually gather and
71 translate knowledge from domain experts or text documents into a machine-readable format.
72 Thus, there is a need for computer-aided methods to remove this knowledge acquisition bot-
73 tleneck (Mounce et al. 2010), so that digital dictionaries can be quickly constructed to meet
74 a specific need and to be able to keep up with the growth of terms due to rapid applications
75 of new technologies and knowledge.

76 Recent achievements in accuracy and processing time of advanced Natural Language

77 Processing (NLP) techniques have driven text mining and cognitive recognition research to
78 a new era. There is a rich set of NLP tools that can support various text processing tasks
79 ranging from basic grammar analyses of individual words (Toutanova et al. 2003; Cunningham
80 et al. 2002), and their dependencies (Chen and Manning 2014), to deep learning of
40 meanings (Mikolov et al. 2013; Pennington et al. 2014). These NLP advances offer numerous
81 potentials for the construction industry where most of the domain knowledge resources
82 are in text documents (e.g., design guidelines, specifications). The implementation of NLP
83 will allow for a fast translation of the domain knowledge into a computer-readable format
84 which is required for machine-to-machine based data exchange.

85 This paper presents an NLP-based automated approach to gather commonly used American-
86 English roadway terms in different highway agencies and classify the semantic relations
87 among these heterogeneous terms. This study leverages NLP and machine learning to extract
88 meanings of terms by analyzing their statistical data of context words in various state design
89 manuals. The semantics of roadway terms are represented as vectors in a high dimensional co-
90 ordinate system in which the semantic similarity among terms is quantifiable. The proposed
91 methodology also includes a new classification algorithm that utilizes syntactic rules and
92 cluster analysis to categorize related terms into three different groups that are synonyms,
93 hyponyms, and attributes. A Java package built upon the proposed method and several
94 datasets resulting from the study can be found at [42](https://github.com/tuyenbk/mvdgenerator)
95 <https://github.com/tuyenbk/mvdgenerator>.

96

97 BACKGROUND

1

98 Natural Language Processing

99 NLP is a research area developing techniques that can be used to analyze and derive value
100 information from natural languages like text and speech. Some of the major applications
101 of NLP include language translation, information extraction, opinion mining (Cambria and
102 White 2014). These applications are embodied by a rich set of NLP techniques ranging
103 from grammar processing such as Tokenization (breaking a sentence into individual tokens)

104 (Webster and Kit 1992; Zhao and Kit 2011), Part-of-Speech (POS) tagging (assigning tags,
105 39 adjective, noun, verb, etc. to each token of a sentence) (Toutanova et al. 2003; Cunningham
106 et al. 2002), and Dependency parser (identifying relationships between linguistic units) (Chen
107 and Manning 2014), to the semantic level, for instance word sense disambiguation (Lesk 1986;
108 Yarowsky 1995; Navigli 2009). NLP methods can be classified into two main groups: (1)
109 7 rule-based and (2) machine-learning (ML) based methods. Rule-based systems, which rely
110 solely on hand-coded syntax rules, are not able to fully cover all human rules (Marcus 1995);
111 and their performances, therefore, are relatively low. In contrast, the ML-based approach
112 is independent of languages and linguistic grammars (Costa-Jussa et al. 2012) as linguistics
113 patterns can be quickly learned from even un-annotated training examples. Thanks to
114 its impressive out-performance, NLP research is shifting to statistical ML-based methods
115 (Cambria and White 2014).

116 **Vector Representation of Word Semantics**

117 Measuring semantic similarity, which is one of the important NLP-related research topics,
118 aims to determine how much two linguistic units (e.g., words, phrases, sentences, concepts)
119 are semantically alike. For example, a *railway* might be more similar to a *roadway* than to
120 7 *train*. The state-of-the-art methodology for this task can be divided into two categories that
121 are (1) thesaurus-based methods and (2) vector space models (VSM) (Harispe et al. 2013).
122 The former approach relies on a hand-coded digital dictionary (e.g., WordNet) that formally
123 structures terms through a network of semantic relations. Computational platforms (e.g.,
124 information retrieval) built upon such dictionaries measure the semantic 1 similarity between
125 a given pair of words by computing the length of their connecting path in the hierarchy. This
126 method would be an ideal solution when digital dictionaries are available. However, digital
127 dictionaries are typically hand-crafted; they are therefore not available to many domains
128 (Kolb 2008). The latter method, on the other hand, assesses the meanings of words or
129 phrases by analyzing their occurrence frequencies in natural language text documents. VSM
130 outperforms the dictionary-based method in terms of time saving as a semantic model can

131 be automatically obtained from a text corpus and corpus collecting is much easier than
132 manually constructing a digital dictionary (Turney and Pantel 2010).

133 VSM estimates semantic similarity based on the *distributional model* which represents
29 the meaning of a word through its context (co-occurring words) in the corpus (Erk 2012).

134 13 The distributional model stands on the *distributional hypothesis* that states that two similar
135 terms tend to occur in the same context (Harris 1954). The outcome of this approach is
136 a Vector Space Model (VSM), in which each vector represents a word in the vocabulary.
137 The similarity between semantic units in this model can be represented by the Euclidean
138 distance between the corresponding points (Erk 2012). The conventional method to construct
139 a VSM is to use the ‘word-context’ matrix which shows how frequent a word is the context
140 of one another in a given text corpus. These raw data of frequencies are used to estimate
141 the co-occurrence probabilities. This statistical process results in a new matrix in which
142 5 each row is a vector representation. Pointwise Mutual Information (PMI) (Church and
143 Hanks 1990) or its variant, Positive PMI (PPMI) is a popular method to calculate the
144 co-occurrence probabilities. A more advanced approach uses machine learning to train the
145 1 representation vectors of terms. One example of this line of methodology is the Skip-gram
146 neural network model (Mikolov et al. 2013) which aims to predict the context words of
147 22 a given input word. The training objective is to minimize the overall error between the
148 predicted and the actual context vectors. Glove (Pennington et al. 2014), an alternative
149 machine learning model for building VSM, trains on the global ‘word-context’ matrix with
150 4 the objective that the probability of co-occurrence between two words equals the dot product
151 1 of their vector representations. The major difference between these two models is that Skip-
152 Gram model trains the local context data within a context window, Glove trains on the global
153 co-occurrence statistics. There are contradictive recommendations on the winning model in
154 the literature. The authors of Glove suggested that their model out-performs Skip-Gram and
155 others in the state of the art. However, a number of independent benchmarking experiments
156 have consistently indicated the outperformance of the Skip-gram model to its alternatives.

158 For example, a comparative study conducted by Levy et al. (2015) on the accuracy in various
159 tasks and golden standards reveals that Skip-gram outperforms Glove in every experiment
160 and is the winner in most of the tasks, especially on the WordSim Similarity dataset. Among
161 these tasks, the best precision of Skip-gram is .793, while PPMI and Glove achieve the highest
162 score of .755 and .725 respectively. The out-performance of Mikolov's model on the similarity
163 task is confirmed in another benchmarking study (Hill et al. 2015) where this model is also
164 found as the winner in most of the tests.

165 The VSM approach has been progressively implemented in recent NLP related studies
166 in the construction industry. Yalcinkaya and Singh (2015) utilized VSM to extract principle
167 research topics related to BIM from a corpus of nearly 1,000 paper abstracts. This approach
168 was also used for information retrieval to search for text documents (Lv and El-Gohary
169 2015) or CAD documents (Hsu 2013). The increasing number of successful use cases in the
170 construction industry has evidently demonstrated that the VSM method can be successfully
171 implemented for identifying the semantic similarity between data labels which is critical to
172 tackle the issue of semantic interoperability in sharing digital data across the life cycle of a
173 highway project. 38

174 Related studies

175 A popular solution to semantic interoperability is to develop taxonomies, ontologies or
176 other forms of digital dictionaries that can provide machine-readable definitions of domain
177 concepts. A plethora of such semantic resources have been developed for the highway in-
178 dustry. However, conventional development methods require significant human efforts on
179 knowledge retrieval, and ontology construction and validation. The pioneer in this line of
180 research is the e-Cognos ontology (Wetherill et al. 2002; Lima et al. 2005) which formulates
181 the execution process of a construction project as an explicitly interactive network of the
182 following principal concepts: Actor, Resources, Products, Processes and Technical Topics.
183 The ontology developers of this project reviewed existing taxonomies (BS61000, UniClass,
184 IFC) and construction specific documents, and interacted with the end users to identify

relevant concepts and their semantic relations. Industry experts were invited to validate e-Cognos through questionnaires on concept names and relations. Since the introduction of the high-level ontology of e-Cognos, a plenty of ontologies have been built for various aspects of a highway project, for instance, highway construction taxonomy (El-Diraby and Kashif 2005; El-Diraby et al. 2005), freight ontology (Seedah et al. 2015a), and the ontology of urban infrastructure products (Osman and Ei-Diraby 2006). Like the e-Cognos project, these studies also relied on domain experts for constructing their semantic products. The limitation regarding time and labor costs of the ad-hoc traditional methodology has created a bottleneck to the progress in enabling semantic interoperability. In addition, the existing ontologies primarily focus on the description of concepts, the heterogeneity of concept names is usually neglected. Therefore, research is needed not only to automate the process of formulating domain concepts but also to incorporate term heterogeneity into ontologies.

Another strategy for semantic interoperability targets at the heterogeneity of concept names rather the concept description as in an ontology model. A few frameworks to assist practitioners in precisely mapping data labels from heterogeneous sources have been introduced for various construction sectors. In the building sector, buildingSMART proposed a novel framework, namely IFD (International Framework for Dictionaries) (ISO 12006-3) for developing a multilingual data schema in which each concept can have multiple names in different languages. With IFD, the identity of a concept is defined by a Global Unique ID (GUID) rather than its name; hence an IFD-based data exchange mechanism is able to eliminate the semantic mismatches due to the name inconsistency (IFD Library Group 2008; Hezik 2008). The buildingSMART data dictionary (bSDD) (buildingSMART 2016) is the first digital library of building concepts that is crafted in the IFD structure. Each concept in bSDD consists a set of synonymy names not only in English but also in computer-coded languages (e.g., IFC) and in other human languages (e.g., French, Norwegian). Therefore, a complete bSDD would enable digital data in regardless of languages to be sharable and unambiguously reusable. Yet, its size remains limited as the identification of these sets

of synonyms is labor and time extensive. In the transportation sector, there has been a shortage of research efforts targeting the heterogeneity of data names at the database level until recently. Seedah et al. (2015b) proposed a role-based classification schema (RBCS) to classify data in freight databases. RBCS defines nine distinct groups of roles that are time (year, month), place (city name, population), commodity (liquid, value), link (roadway name, width), mode (truck, rail), industry (company name, sales), event (accident, number of fatalities), and human (officer, driver age). The authors argue that once the data elements across separate databases are categorized using this standard system, it becomes easier for practitioners to identify the semantic relatedness in their definitions. However, even if RBCS is successfully applied to all freight databases, identifying the exact type of relation (synonym, functional relation) between two data elements in the same category is still a challenging task.

In attempts to reduce laborious work on defining concepts, a few researchers have sought to propose semi-automated and automated methods for identifying semantic relations among technical terms. Abuzir and Abuzir (2002) developed the ThesWB system which utilizes hand-coded syntax patterns to detect lexical relations between civil engineering terms from HTML web pages. The performance of ThesWB was not reported, but it is not likely to be high since rule-based approaches are repeatedly criticized for not being able to capture all the variant ways to present relations among terms in natural language (Marcus 1995; Navigli and Velardi 2010). Rezgui (2007) suggested a more sophisticated approach that is based the statistics of word occurrence rather than predefined rules to extract potential pairs of related terms from domain text documents. This method implements TF-IDF to evaluate the importance degree of a keyword to the examined domain; and analyzes the co-occurrence frequencies using Metric Clusters to assess the potentiality that exists a semantic relation within a given pair of important keywords. These potential relationships are then validated and categorized by domain experts. Since the method detects relations between words occurring in the same sentence, equivalent terms which are used interchangeably

239 could not be captured. In another study to identify semantic relations, Zhang and El-
240 Gohary (2016) proposed a fully automated methodology for both tasks of retrieving related
241 candidate and classifying the relations. This algorithm was reported to achieve an average
242 precision of nearly 90 percent in the relation classification task. However, the algorithm
243 identifies potentially related concepts based on the pre-defined lexical relations provided in
244 WordNet, a generic lexicon that lacks concepts in many construction sectors including the
245 civil infrastructure, it would not be scalable well on matching terms in these domains.

246 As shown in the literature review, there are numerous research efforts in developing
247 ontologies for the highway sector. However, the existing ontologies are mainly hand-coded
248 through the manual processes of knowledge acquisition and translation into a digital format.
249 This ad-hoc approach has created a bottleneck in facilitating the semantic interoperability
250 for the whole industry and as a result, semantic resources for many aspects of a project are
251 still not available. A few efforts have been made to automate the process of constructing or
252 extending existing semantic resources. The most rigorous methodology in the state-of-the-
253 art is the one developed by Zhang and El-Gohary (2016) that is fully-automated with high
254 accuracy. 1 One limitation of this algorithm is the reliance on an existing semantic resource;
255 it, therefore, would not be applicable to such a domain like the infrastructure that is out of
256 the vocabulary scope. Thus, there is a need for an automated approach that can not only
257 allow for a fast development of highway lexicons but also remove the dependence on other
258 existing semantic models.

259 **PROPOSED METHODOLOGY TO AUTOMATED CLASSIFICATION OF**
260 **ROADWAY TERMS**

261 The goal of this research is to propose an NLP-based methodology that can automate the
262 process of extracting roadway technical terms and their semantic relations from American-
263 English roadway documents. 5 As shown in Figure 1, the proposed methodology consists of
264 three major modules that are to: (1) utilize NLP techniques to extract multi-word roadway
265 technical terms from a collected text corpus, (2) train the data obtained form the text corpus

266 1

using the Skip-gram neural network model (Mikolov et al. 2013) to develop a Roadway Vector Space Model (Rd-VSM) that presents the semantics of roadway terms, and (3) develop an algorithm integrating Rd-VSM and various linguistic patterns to classify relations among technical terms (synonyms, hyponyms and attributes). The below sections discuss these steps in detail.

271 **Text corpus collection**

272 In order to capture the heterogeneity of roadway terms, the authors collected a plethora of
273 highway engineering manuals and guidelines from 30 State Departments of Transportation.
274 The content of a written guidance document in the engineering field is commonly presented
275 in various formats such as plain text, tables, and equations. Since the structures of words
276 in tables and equations are not yet supported by the state-of-the-art NLP techniques, they
277 were removed from the text corpus. The removal of these features may slightly reduce the
278 corpus size, and accordingly affects the training dataset; however, it is necessary since words
279 in tables and equations are not organized in the formal structure of a sentence and therefore
280 the NLP algorithm may extract unreal noun phrases. The final outcome of this phase is a
281 plain text corpus consisting of nearly 16 million words. This dataset is utilized to extract
282 multiple-word technical terms which are then trained and transformed into representation
283 vectors.

284 **Multi-word terms extraction**

285 Linguists argue that a technical term is either a noun (e.g., road) or a noun phrase
286 (NP) (e.g., right of way) that frequently occurs in domain text documents (Justeson and
287 Katz 1995). The meaning of a multi-word term may not be directly interpreted from the
288 meanings of its constituents; therefore, it must be treated as an individual word. To meet
289 that requirement, multi-word terms in the corpus need to be detected and replaced with
290 connected blocks of their members. As mentioned, a multiple-word term must be an NP;
291 thus, NPs will be good multi-word term candidates. To detect this type of term, the corpus
292 is first scanned to search for NPs, of which the importance is then evaluated based on their

293 statistics of occurrence. The process of extracting multi-word terms is discussed in detail as
294 follows.

295 *Noun phrase extraction*

296 This research implements the Apache OpenNLP package to find sequences of words that
297 match pre-defined noun phrase patterns. Figure 2 illustrates how noun phrases are extracted
298 from the corpus of highway technical documents. This process includes the following steps.

299 i Word tokenizing: In this step, the text corpus is broken down into individual units
300 (also called tokens) using OpenNLP Tokenizer.

301 ii Part of Speed (POS) tagging: The purpose of this step is to determine the Part of
302 Speech (POS) tag (e.g., NN-noun, JJ-adjective, VB-verb, etc.) for each unit of the
303 tokenized corpus obtained from the previous step. A full set of POS tags can be found
304 in the Penn Treebank (Marcus et al. 1993).

305 iii Noun phrase detection: Table 1 presents the proposed extraction patterns which are
306 modified from the filters suggested by Justeson and Katz (1995) to extract NPs. The
307 tagged corpus is thoroughly scanned to collect sequences matching those patterns. In
308 addition, in order to reduce the discrimination between the syntactic variants of the
309 same term, the collected NPs need to be normalized. The following discuss two types
310 of syntactic variants considered and the proposed normalization methods.

- 311 • Type 1 - Plural forms, for example ‘roadways’ and ‘roadway’. Stemming is a
312 popular process to reduce words to their stems. Despite the fact that, none
313 of the existing algorithms can completely eliminate the errors of over and un-
314 der stemming, they are good enough to not degrade the overall performance of
315 NLP applications (Jivani et al. 2011). This study implements the Pling stem-
316 mer (Suchanek et al. 2006), which stems an English noun to its singular form,
317 to normalize plural nouns in the corpus. One advantage of this algorithm is the

318 utilization of both syntactic rules and the vocabulary in a dictionary; hence the
319 miss- or over-stemming errors that take off a true suffix can be reduced.

- 320 • Type 2 - Preposition noun phrases, for example ‘type of roadway’ and ‘roadway
321 type’. In order to normalize this type of variant, the form with preposition is
322 converted into the non-preposition form by removing the preposition and reversing
323 the order of the remaining portions. For instance, ‘type of roadway’ will become
324 ‘roadway type’.

325 The first column in Table 2 represents several examples of the NP bag retrieved from this
326 phase. Since an NP is not certainly a technical term, those that are clearly unlikely to be a
327 term should be excluded from the candidate list. Occurrence frequency is a key indicator for
328 the importance of a candidate as a technical term tends to repeatedly occur in domain text
329 documents. To eliminate ‘bad’ candidates, a threshold of frequency can be applied. If users
330 choose a high threshold, rare terms would not be captured. This issue can be addressed
331 when the corpus size is extended. In our experiment, with a frequency threshold of 2, the
332 final list of NPs consists of 112,024 items; and it drops to 8,922 when a threshold of 50 is
333 used. Since this research aims at common technical terms, the authors used a threshold of
334 50 to remove possibly meaningless term candidates.

335 *Multi-word term candidate ranking and selection*

336 Multi-word term definition varies between authors, and there is a lack of formal and
337 widely accepted rules to define if an NP is a multi-word term (Frantzi et al. 2000). There are
338 a number of methods proposed for estimating termhood (the degree that a linguistic unit
339 is a domain-technical concept), such as TF-IDF (Sparck Jones 1972; Salton and Buckley
340 1988), C-Value (Frantzi et al. 2000), Termex (Sciano and Velardi 2007). These methods are
341 based on the occurrence frequencies of NPs in the corpus. Among these methods, Termex
342 outperforms other methods on the Wikipedia corpus, and C-Value is the best on the GENIA
343 medical corpus (Zhang et al. 2008). This result indicates that the C-value method is more

36

344 suitable for term extraction from a domain corpus rather than a generic corpus. For this
 345 reason, the C-value has been widely used to extract domain terms in the biomedical field,
 346 for instance studies performed by Ananiadou et al. (2000), Lossio-Ventura et al. (2013), and
 347 Nenadić et al. (2002). Since the corpus used in this study was mainly collected from technical
 348 domain documents, C-value would be the most suitable for the termhood determination task.
 349 The C-value measure, as formulated in Equation 1, suggests that the longer an NP is, the
 350 more likely that is a term; and the more frequently it appears in the domain corpus, the
 351 more likely it will be a domain term.

$$352 \quad C-value(a) = \begin{cases} 353 \quad 8 \\ log_2|a| \cdot f(a), & \text{if } a \text{ is not nested} \\ log_2|a|(f(a) - \frac{1}{P(T_a)} \sum_{b \in T_a} f(b)), & \text{otherwise} \end{cases} \quad (1)$$

353 Where:

354 **a** is a candidate noun phrase

355 **|a|** is the length of noun phrase **a**

356 **f** is the frequency of **a** in the corpus

357 **Ta** is the set of extracted noun phrases that contain **a**

358 **P(Ta)** is the size of **Ta** set.

359 The term extraction process above results in a dataset containing the detected terms
 360 along with their c-value termhood scores. These term candidates are ranked by C-value, and
 361 the ones that have negative C-values are discarded.

362 To automatically remove candidates that are unlikely to be real terms, a threshold C-
 363 value can be used. However, doing this may eliminate the real terms that appear in the
 364 bottom due to their low frequencies. Manual evaluation of the entire candidate list would
 365 avoid the removal of real terms with low C-values. To minimize both laborious work and

366 the number of true terms wrongly discarded, the authors suggest the following method to
367 identify the threshold value. The ranked list of candidates is divided into groups of around
368 200 items. A graduate student with a civil engineering background was asked to utilize
369 a bottom-up approach to evaluate group by group and stop at which the percentage of
370 actual terms achieved 80 percent. Users can choose a higher percentage limit in cases where
371 the accuracy is critical. This will increase manual evaluation effort. Table 2 illustrates the
372 evaluation results for several excerpts of the extracted term candidates. The precision values,
373 which represent the percentages of real terms in these groups, are presented in Figure 3. As
374 shown in the figure, precision values are less than 80 percent for groups with c-values less
375 than 50. This value is set as the threshold for the acceptance of term candidates. The final
376 selected list is comprised of nearly 8,000 multi-word roadway technical terms.

377 Construction of term space model

378 This step aims at converting the vocabulary in the roadway corpus into a vector space
379 model, namely Rd-VSM. Skip-gram (Mikolov et al. 2013), which is an un-supervised machine
380 model, is employed to learn the semantic similarity among words in the text corpus. The
381 Skip-Gram model requires a set of training data in which the input data is a linguistic unit
382 (word or term), and the output data is a set of context words that appear around the input
383 unit in the corpus. In order to collect this training dataset, the tokenized and stemmed
384 highway corpus is scanned to capture instances of terms and their corresponding context
385 words. Each occurrence of a word will correspondingly generate a data point in the training
386 dataset.

387 Before collecting the training dataset, an additional step is needed to handle the issue
388 related to multi-word terms. Since document scanning is on a word-by-word basis, the corpus
389 must be adjusted so that multi-word terms can be treated as single words. To fulfill that
390 requirement, the white spaces within a multi-word term are replaced with minus (-) symbols
391 to connect its individual words into a single unit. For instance, ‘vertical alignment’ becomes
392 ‘vertical-alignment’.

393 The number of context words to be collected is dependent on the window size that limits
394 12 how many words to the left and the right of the target word. In the example sentence below,
395 the context of the term ‘roadway’ with the window size of 5 will be the following word set
396 {bike, lane, width, on, a, with, no, curb, and, gutter}. Any context word that is in the stop
397 list (a list that contains frequent words in English such as ‘a’, ‘an’, and ‘the’ that have little
398 meaning) will be neglected from the context set. In this example, the adjusted context set
399 of ‘roadway’ is {bike, lane, width, curb, gutter}.

400 "The minimum [bike lane width on a roadway with no curb and gutter] is 4 feet ."

401 The semantic similarity is trained using the Word2Vec module in the Apache Spark MLlib
402 package (Apache.org 2016), an emerging open-source engine, which is based on the Skip-
403 gram neural network model (Mikolov et al. 2013). Figure 4 illustrates the learning network
404 when the context set includes only one word, where V and N respectively denote the corpus
405 vocabulary and hidden layer size. In this model, a word in the corpus vocabulary is encoded
406 as a ‘one-hot’ vector which is a vector in which only one element at the index of the word in
407 the vocabulary is set one, and all other items are zero. For example, the one-hot vector of
408 k^{th} word in the vocabulary with the size of V will be $\{x_1 = 0, x_2 = 0, \dots, x_k = 1, \dots, x_V = 0\}$.
409 The outcome of this machine learning process is a set of word representation vectors in an
410 N -dimension coordinate system. The similarity among these vectors represents the similarity
411 in context between the corresponding words. The bullets below explain how the predicted
412 context vector of k^{th} word is computed using the parameter matrices resulted from the
413 learning process. As we can see, the similarity between two predicted context vectors depends
414 only on the similarity between their corresponding input representation vectors; thus, these
415 vectors are used to represent the semantics of words.

- 416 • k^{th} word: $[x_k]_{1..V} = [x_1 = 0, x_2 = 0, \dots, x_k = 1, \dots, x_V = 0]$ which is an one-hot vector.
- 417 • Hidden vector: $[h]_{1..N} = [x_k]_{1..V} \cdot W_{V..N} = [w_{k1}, w_{k2}, \dots, w_{kN}] = v_{wk}$ which is equivalent
418 to the k^{th} row of the W matrix since the input vector is a ‘one-hot’ vector. The v_{wk}

419 vector is called the input *representation vector* of the k^{th} word.

- 420 • Predicted context vector: $[y_k]_{1..V} = v_{wk} \cdot W'_{N..V}$.

421 The learning model includes three major parameters that are *frequency threshold*, *hidden*
422 *layer size* and *window size* (see Table 3). To eliminate those data points with low frequen-
423 cies of occurrence that are unlikely to be technical terms, Word2Vsec allows for the use of
424 *frequency threshold*. Any word with the rate lower than the limit will be ignored. Radim
425 (2014) suggests a range of (0-100) depending on the data set size. Setting this parameter
426 high will enhance the accuracy, but many true technical terms would be out of vocabulary.
427 A preliminary study based on the preliminary corpus with only several millions of words
428 shows that with the frequency of 20, there are very few non-technical terms involved in the
429 training dataset. Hence, with the larger dataset to be collected, this parameter can be higher
430 and up to around 50. The second important parameter is *layer size* which determines the
431 number of nodes in the hidden layer. This parameter highly affects the training accuracy
432 and processing time. A larger layer size is better in terms of accuracy, but this will be paid
433 off by the running time. A reasonable configuration for this parameter is from tens to hundreds
434 (Radim 2014). The final major parameter, *context window size*, decides how many context
435 words to be considered. Google recommends a size of 10 for the Skip-gram model (Google
436 Inc. 2016). These parameters are subject to be changed so that the best model can be
437 achieved. The effects of these parameters on the model performance are discussed in Section
438 4.

439 Figure 5 presents the Rd-VSM vector space model derived from the training process when
440 the parameters, *frequency threshold*, *hidden layer size* and *window size* are set 50, 300 and
441 10 respectively. In this model, each word in the highway corpus is represented as a vector in
442 a high dimensional space. Since the representation vectors are in a multi-dimensional space;
443 to present the space in 2D graph, PCA (Principle Component Analysis) is used to reduce
444 the dimension size to two.

445 The similarity between terms in the Rd-VSM model can be measured by the angle be-

446 tween two word representation vectors (Equation 2) or the distance between two word points
447 (Equation 3). Figure 5 illustrates the clustering of terms by their distances. In this figure,
448 an *inlet* can be inferred to be more similar to an *outlet* (blue) than a *sidewalk* (green). Using
449 this technique, the most similar terms for a given term can be obtained. Table 4 shows a
450 partial ranked list of the nearest terms of ‘roadway’ in order of similarity score.

451

$$\text{cosine_similarity} = \frac{A \cdot B}{\|A\| \|B\|} \quad (2)$$

452

$$\text{dis_similarity} = \sqrt{(xA_1 - xB_1)^2 + (xA_2 - xB_2)^2 + \dots + (xA_n - xB_n)^2} \quad (3)$$

453

454 Where: n is the hidden layer size.

455 **Semantic relation classification**

456 The purpose of this module is to design an algorithm for automated classification of
457 the semantic relations among the roadway technical terms. This study considers three core
458 relation types of a semantic resource that are: synonym (meaning equivalence), hypernym-
459 hyponym (also known as IS-A or parent-child relation), attribute (concept property) (Jiang
460 and Conrath 1997; Lee et al. 2013). The following describes the fundamental logic behind
461 the designed algorithm. Two terms that relate to each other through these semantic relations
462 would have a high similarity score. Therefore, the top nearest terms resulted from Rd-VSM
463 would be a great starting point for detecting relations between technical terms. For example,
464 in the list of the nearest terms of ‘roadway’ (see Table 4), true synonyms are ‘highway’ (rank
465 1), ‘traveled-way’ (2) and ‘road’ (4); attributes include ‘roadway-section’ (3), ‘roadway-
466 shoulder’ (12); and ‘adjacent-roadway’ (7) and ‘undivided’ (37) are hyponyms which show
467 different types of roadway.

468 Algorithm 1 shows the designed pseudo code for classifying the nearest terms of a given
469 target term. The algorithm utilizes linguistic rules and clustering analysis to organize the
470 nearest list into the following three groups: (1) attribute, (2) hyponym, and (3) synonym.
471 The algorithm first detects terms belonging to the first two categories using linguistic pat-

472 terns, and employs cluster analysis for the last group.

473 *Attributes and hyponyms*

474 The filter rules to detect these relations are presented in Table 5. For a multi-word term
475 matching pattern 1, we can infer that *Noun1* is an attribute of concept *Noun2*; and *Noun2*
476 is an attribute of *Noun1* in the pattern 2. Pattern 3 is for detecting hyponyms where the
477 matched NP is a hyponym of its *Noun2* component.

478 *Synonyms*

479 After the words in the first two categorized are classified, the remained nearest words
480 will fall into the third group. However, some of them may have far or even no relation
481 with the target word. In order to address this issue, this framework employs the K-mean
482 clustering algorithm (MacQueen 1967) to split the remained list into multiple layers based
483 on the similarity score. Those terms in the last layers are unlikely to be synonyms; and
484 thus, are removed from the classified list. Only the terms in the top cluster are kept and
485 categorized as synonyms.

486 By the end of the synonym recognition phase, the algorithm will generate a list of classified
487 nearest terms for a given input word. Table 6 shows one example of the output generated
488 by the algorithm.

489 **14 PERFORMANCE EVALUATION**

This section presents a performance evaluation of the proposed system on the ability to identify synonyms. In this experiment, a gold standard is used. The gold standard consists of 70 sets of synonyms (both single and multi-word terms) which were examined and extracted from a Wikipedia transportation glossary (Wikipedia 2016). The developed algorithm was employed to find the synonym for a given input term. The automatically identified synonym is the nearest word in the synonym lexical group. The evaluation outcome returns “true” if the automatically identified synonym belongs to the actual synonym set of the tested term in the golden standard, or “false” if it does not. The answer will be “N/A” if the target

Algorithm 1 Semantic relation classification algorithm

```
1: Inputs: term  $t$ , list of nearest terms  $N$ , list of multi-word terms  $F$ 
2: Output: Classified list of terms  $C$ 
3: procedure TERM CLASSIFICATION PROCEDURE
4:    $Att \leftarrow$  list of attributes
5:    $Hyp \leftarrow$  list of hyponyms
6:    $Syn \leftarrow$  list of synonyms
7:    $w \leftarrow null$ 
8:   for all  $n \in N$  do
9:     if  $n$  contains  $t$  then
10:       $w \leftarrow n$ 
11:    else
12:      for all  $f \in F$  do
13:        if  $f$  contains both  $n$  and  $t$  then
14:           $w \leftarrow f$ 
15:          Break for
16:        if  $w$  matches Attribute pattern then
17:          add  $w$  to  $Att$ 
18:        else if  $w$  matches Hyponym pattern then
19:          add  $w$  to  $Hyp$ 
20:        else
21:          add  $w$  to  $Syn$ 
22:   Cluster  $Syn$  and discard low relevant terms
```

15

term is out of the model vocabulary. The performance was evaluated using the following three measures including precision, recall, and f-measure. Precision refers the accuracy in the conclusions made by the system, and recall reflects the coverage of domain terms of the system. The F score, which is a combined measure of precision and recall, presents the overall performance of a system.

$$Precision = \frac{\text{number of correctly detected synonyms}}{\text{total detected synonyms}} \quad (4)$$

$$Recall = \frac{\text{number of correctly detected synonyms}}{\text{total tested terms}} \quad (5)$$

$$F - measure = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall} \quad (6)$$

490 Table 7 shows the performance with various training model settings. The parameters of
491 the baseline model are 50, 100 and 5 respectively for frequency threshold, hidden layer and

492 window size. The authors changed the configuration of these parameters one by one to
493 evaluate their effects to the model performance. While changing a certain parameter, other
494 parameters are kept unchanged compared to their values in the base model. As presented
34
495 in the table, the model performance is not significantly sensitive to the changes of training
496 parameters. The increase of window size to 10 or 15 resulted in the best model which has
1
497 a precision of 81% and an F-measure of 65%. The changes of other parameters did not
498 improve the performance. Especially, the increase of frequency threshold value from 50 to
499 75 has negative impact to all measures. This result confirms the reasonable selection of the
500 frequency threshold to eliminate unlikely term candidates in the NP extraction phase.

501 The proposed model was also compared with the generic WordNet database. Table
4
502 8 presents the comparison of performance between the proposed framework (with the 50-
503 100-10 setting) and WordNet. As shown, the present system outperforms WordNet in all
504 measures, and the combined F-measure is significantly improved (65% compared to 52%).
505 The biggest contribution to the improvement of the overall F-measure is the recall value
506 which represents a better coverage of roadway vocabulary.

507 DISCUSSIONS

508 This paper proposes an NLP based methodology to assist professionals in extracting
509 roadway terms and their semantic relations from text documents. A key contribution to the
510 body of knowledge is the novel framework with a new algorithm that allows for automated
511 detection of technical terms and their relations without reliance on existing hand-coded dic-
512 tionaries as used by previous researchers such as Zhang and El-Gohary (2016). The present
513 framework is not to completely eliminate human involvement, but is expected to significantly
514 reduce manual efforts and become an enabling tool that can help researchers in the high-
515 way domain quickly develop supporting ontologies and other forms of semantic resources
516 for their specific use cases. With respect to the facilitating semantic interoperability for the
517 infrastructure sector, the findings of this study would accelerate the process of removing
518 the current bottleneck in extensive machine readable dictionaries which are required for an

519 unambiguous data sharing, integration or exchange.

520 The semantic similarity model and the relation classification algorithm developed in this
521 study are also expected to become fundamental resources for a variety of NLP related stud-
522 ies in the highway domain. NLP based platforms can utilize these resources for term sense
523 analysis which is crucial for text mining to extract meaningful information from text docu-
524 ments, information retrieval, or natural language based human-machine interaction. Some
525 specific examples of these potential applications are as follows. Information retrieval systems
526 can use the semantic relations provided by the algorithm to classify project documents by
527 relevant topics by analyzing the relatedness between the index keywords in those documents.
528 In addition, questionnaire designers can utilize the system to search for synonyms so that
529 appropriate terms can be selected for specific groups of potential respondents who might be
530 from multiple disciplines or regions. Another application is that query systems for extracting
531 data from 3D engineered models would be able to find alternative ways to query data when
532 users' keywords do not match any entity in the database. Since users have different ways and
533 keywords to query data, the ability to recognize synonyms and related concepts of a query
534 system would provide flexibility to the end user. Also, the synonym detection function would
535 enable the matching data items such as (e.g., cost, productivity, etc.) when integrating data
536 from distinct departments or states to develop a national database. Moreover, this study
537 can make fundamentally transform to the way human interacts with a machine as technical
538 terms which are a basic unit of human language can be precisely understood by computer
539 systems. Instead of using computer languages, the end user can use natural language to
540 communicate with computer systems.

541 The current study has a number of limitations. First, the highway corpus is still relatively
542 small with only 16 million words, compared to the corpus sizes in other domains with billions
543 of words. Since the recall value largely depends on the corpus size, the expansion of the
544 highway corpus size by adding more documents from other state agencies and disciplines
545 (e.g., survey, construction, operation and maintenance) would enable more technical terms

546 to be covered in the vector space model; and consequently, the recall would be improved.
547 Secondly, the number of semantic relation categories is limited to only three types of relations
548 that are attributes, hyponyms and synonyms. There are other important semantic relations
549 that are not considered such as hypernyms, siblings, functional associations, etc. Including
550 these relations would reduce incorrect synonym matching, which will enhance the precision
551 value, for those cases that a word does not have any equivalent term. Third, this study only
552 targets at the synonymy issue, the issue of polysemy is not yet addressed. Further research
553 is needed to detect different senses of a roadway term. One potential solution is to apply
554 cluster analysis on the instances of context to determine the possibility that a term would
555 have multiple meanings.

556 **CONCLUSIONS**

557 Data manipulation from multiple sources is a challenging task in highway asset man-
558 agement due to the inconsistency of data format and terminology. The contribution of this
559 study is an NLP-based approach to automated classification of semantic relations among
560 roadway technical terms based on their word occurrences in domain text documents. This
561 research employs advanced NLP techniques to extract technical terms from a highway text
562 corpus which is composed of 16 million words built on a collection of design manuals from 30
563 State DOTs across the U.S. Machine learning is used to train the semantic similarity between
564 technical terms. An algorithm is designed to classify the nearest terms resulted from the
565 semantic similarity model into distinct groups according to their lexical relationships.

566 The developed system has been evaluated by comparing the results obtained from the
567 computational model and a man-crafted gold standard. The result shows an accuracy of
568 over 80 percent. The best model is associated with the training parameters of 50, 100 and
569 10 respectively for frequency threshold, hidden layer size, and window size. Although a
570 significant improvement has been made in comparison with an existing thesaurus database,
571 the overall performance is not relatively high. This might be due to the limited ³² size of the
572 training data. Future research will be conducted to expand the highway corpus to further

573 disciplines such as asset management, and transportation operation.

574 The proposed automated methodology for detecting semantic relations from manuals is
575 expected to significantly reduce human efforts in developing a semantic resource for a specific
576 use case within the highway domain and become an enabler for semantic interoperability in
577 this domain. The research also opens a new gate for computational tools regarding natural
578 language processing in the highway sector. The developed system would enable computer
579 systems to understand terms and consequently transform the way human interacts with a
580 computer by allowing users to use natural language.

581 REFERENCES

- 582 Abuzir, Y. and Abuzir, M. O. (2002). "Constructing the civil engineering thesaurus (cet)
583 using theswb." *Computing in Civil Engineering*.
- 584 Ananiadou, S., Albert, S., and Schuhmann, D. (2000). "Evaluation of automatic term recog-
585 nition of nuclear receptors from medline." *Genome Informatics*, 11, 450–451.
- 586 Apache.org (2016). "Machine learning library (mlib), <[https://spark.apache.org/docs/1.1.0/mllib-](https://spark.apache.org/docs/1.1.0/mllib-guide.html)
587 [guide.html](https://spark.apache.org/docs/1.1.0/mllib-guide.html)>.
- 588 buildingSMART (2016). "buildingsmart data dictionary, <<http://bsdd.buildingsmart.org/>>.
589 (Accessed: March 15, 2016).
- 590 Cambria, E. and White, B. (2014). "Jumping nlp curves: a review of natural language
591 processing research [review article]." *Computational Intelligence Magazine, IEEE*, 9(2),
592 48–57.
- 593 Chen, D. and Manning, C. D. (2014). "A fast and accurate dependency parser using neural
594 networks." *EMNLP*, 740–750.
- 595 Church, K. W. and Hanks, P. (1990). "Word association norms, mutual information, and
596 lexicography." *Computational linguistics*, 16(1), 22–29.
- 597 Costa-Jussa, M. R., Farrús, M., Mariño, J. B., and Fonollosa, J. A. (2012). "Study and
598 comparison of rule-based and statistical catalan-spanish machine translation systems."
599 *Computing and Informatics*, 31(2), 245–270.

- 600 Cunningham, H., Maynard, D., Bontcheva, K., and Tablan, V. (2002). "Gate: an architecture
601 for development of robust hlt applications." *Proceedings of the 40th annual meeting on*
602 *association for computational linguistics*, Association for Computational Linguistics, 168–
603 175.
- 604 El-Diraby, T. and Kashif, K. (2005). "Distributed ontology architecture for knowledge man-
605 agement in highway construction." *Journal of Construction Engineering and Management*,
606 131(5), 591–603.
- 607 El-Diraby, T., Lima, C., and Feis, B. (2005). "Domain taxonomy for construction concepts:
608 Toward a formal ontology for construction knowledge." *Journal of Computing in Civil*
609 *Engineering*, 19(4), 394–406.
- 610 Erk, K. (2012). "Vector space models of word meaning and phrase meaning: A survey."
611 *Language and Linguistics Compass*, 6(10), 635–653.
- 612 Frantzi, K., Ananiadou, S., and Mima, H. (2000). "Automatic recognition of multi-word
613 terms: the c-value/nc-value method." *International Journal on Digital Libraries*, 3(2),
614 115–130.
- 615 Gallaher, M. P., O'Connor, A. C., Dettbarn, J. L., and Gilday, L. T. (2004). *Cost analysis*
616 *of inadequate interoperability in the U.S. capital facilities industry*. U.S. Department of
617 Commerce Technology Administration, National Institute of Standards and Technology,
618 Gaithersburg, MD.
- 619 Google Inc. (2016). "word2vec, <<https://code.google.com/archive/p/word2vec/>>." (accessed
620 May 12, 2016).
- 621 Harispe, S., Ranwez, S., Janaqi, S., and Montmain, J. (2013). "Semantic measures for the
622 comparison of units of language, concepts or instances from text and knowledge base
623 analysis." *arXiv preprint arXiv:1310.1285*.
- 624 Harris, Z. S. (1954). "Distributional structure." *Word*.
- 625 Hezik, M. (2008). "Ifd library background and history." *The IFD Library/IDM/IFC/MVD*
626 *Workshop*.

- 627 Hill, F., Reichart, R., and Korhonen, A. (2015). "Simlex-999: Evaluating semantic models
628 with (genuine) similarity estimation." *Computational Linguistics*, 41(4), 665–695.
- 629 Hsu, J.-y. (2013). "Content-based text mining technique for retrieval of cad documents."
630 *Automation in Construction*, 31, 65–74.
- 631 IFD Library Group (2008). "Ifd library white paper. Accessed: 2015-07-06.
- 632 Jiang, J. J. and Conrath, D. W. (1997). "Semantic similarity based on corpus statistics and
633 lexical taxonomy." *arXiv preprint cmp-lg/9709008*.
- 634 Jivani, A. G. et al. (2011). "A comparative study of stemming algorithms." *Int. J. Comp.
635 Tech. Appl.*, 2(6), 1930–1938.
- 636 Justeson, J. S. and Katz, S. M. (1995). "Technical terminology: some linguistic properties
637 and an algorithm for identification in text." *Natural Language Engineering*, 1(01), 9–27.
- 638 Karimi, H. A., Akinci, B., Boukamp, F., and Peachavanish, R. (2003). "Semantic interoper-
639 ability in infrastructure systems." *Information Technology*, 42–42.
- 640 Kolb, P. (2008). "Disco: A multilingual database of distributionally similar words." *Proceed-
641 ings of KONVENS-2008, Berlin*.
- 642 Lee, T., Wang, Z., Wang, H., and Hwang, S.-w. (2013). "Attribute extraction and scor-
643 ing: A probabilistic approach." *Data Engineering (ICDE), 2013 IEEE 29th International
644 Conference on*, IEEE, 194–205.
- 645 Lefler, N. X. (2014). "Nchrp synthesis 458: Roadway safety data interoperability between
646 local and state agencies." *Report no.*, Transportation Research Board.
- 647 Lesk, M. (1986). "Automatic sense disambiguation using machine readable dictionaries: how
648 to tell a pine cone from an ice cream cone." *Proceedings of the 5th annual international
649 conference on Systems documentation*, ACM, 24–26.
- 650 Levy, O., Goldberg, Y., and Dagan, I. (2015). "Improving distributional similarity with
651 lessons learned from word embeddings." *Transactions of the Association for Computational
652 Linguistics*, 3, 211–225.
- 653 Lima, C., El-Diraby, T., and Stephens, J. (2005). "Ontology-based optimization of knowledge

- 654 management in e-construction." *Journal of IT in Construction*, 10, 305–327.
- 655 Lossio-Ventura, J. A., Jonquet, C., Roche, M., and Teisseire, M. (2013). "Combining c-
- 656 value and keyword extraction methods for biomedical terms extraction." *LBM'2013: 5th*
- 657 *International Symposium on Languages in Biology and Medicine*.
- 658 Lv, X. and El-Gohary, N. M. (2015). "Semantic annotation for context-aware information
- 659 retrieval for supporting the environmental review of transportation projects." *Computing*
- 660 *in Civil Engineering 2015*, ASCE, 165–172.
- 661 MacQueen, J. (1967). "Some methods for classification and analysis of multivariate observa-
- 662 tions." *Proceedings of the fifth Berkeley symposium on mathematical statistics and proba-*
- 663 *bility*, Vol. 1, Oakland, CA, USA., 281–297.
- 664 Marcus, M. (1995). "New trends in natural language processing: statistical natural language
- 665 processing." *Proceedings of the National Academy of Sciences*, 92(22), 10052–10059.
- 666 Marcus, M. P., Marcinkiewicz, M. A., and Santorini, B. (1993). "Building a large annotated
- 667 corpus of english: The penn treebank." *Computational linguistics*, 19(2), 313–330.
- 668 Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). "Efficient estimation of word
- 669 representations in vector space." *arXiv preprint arXiv:1301.3781*.
- 670 Mounce, S., Brewster, C., Ashley, R., and Hurley, L. (2010). "Knowledge management for
- 671 more sustainable water systems." *Journal of information technology in construction*, 15,
- 672 140–148.
- 673 Navigli, R. (2009). "Word sense disambiguation: A survey." *ACM Computing Surveys*
- 674 (*CSUR*), 41(2), 10.
- 675 Navigli, R. and Velardi, P. (2010). "Learning word-class lattices for definition and hypernym
- 676 extraction." *Proceedings of the 48th Annual Meeting of the Association for Computational*
- 677 *Linguistics*, Association for Computational Linguistics, 1318–1327.
- 678 Nenadić, G., Spasić, I., and Ananiadou, S. (2002). "Automatic acronym acquisition and term
- 679 variation management within domain-specific texts." *Third International Conference on*
- 680 *Language Resources and Evaluation (LREC2002)*, 2155–2162.

- 681 Noy, N. F. (2004). "Semantic integration: a survey of ontology-based approaches." *ACM*
682 *Sigmod Record*, 33(4), 65–70.
- 683 Osman, H. and Ei-Diraby, T. (2006). "Ontological modeling of infrastructure products and
684 related concepts." *Transportation Research Record: Journal of the Transportation Research*
685 *Board*, 1984(-1), 159–167.
- 686 Ouksel, A. M. and Sheth, A. (1999). "Semantic interoperability in global information sys-
687 tems." *ACM Sigmod Record*, 28(1), 5–12.
- 688 Pennington, J., Socher, R., and Manning, C. D. (2014). "Glove: Global vectors for word
689 representation." *Empirical Methods in Natural Language Processing (EMNLP)*, 1532–1543,
690 <<http://www.aclweb.org/anthology/D14-1162>>.
- 691 Radim, R. (2014). "Word2vec tutorial, <<http://rare-technologies.com/word2vec-tutorial/>>."
- 692 Rezgui, Y. (2007). "Text-based domain ontology building using tf-idf and metric clusters
693 techniques." *The Knowledge Engineering Review*, 22(04), 379–403.
- 694 Salton, G. and Buckley, C. (1988). "Term-weighting approaches in automatic text retrieval."
695 *Information processing & management*, 24(5), 513–523.
- 696 Sclano, F. and Velardi, P. (2007). *Termextractor: a web application to learn the shared*
697 *terminology of emergent web communities*. Springer, 287–290.
- 698 Seedah, D. P., Choubassi, C., and Leite, F. (2015a). "Ontology for querying heteroge-
699 neous data sources in freight transportation." *Journal of Computing in Civil Engineering*,
700 04015069.
- 701 Seedah, D. P., Sankaran, B., and O'Brien, W. J. (2015b). "Approach to classifying freight
702 data elements across multiple data sources." *Transportation Research Record: Journal of*
703 *the Transportation Research Board*, (2529), 56–65.
- 704 Sparck Jones, K. (1972). "A statistical interpretation of term specificity and its application
705 in retrieval." *Journal of documentation*, 28(1), 11–21.
- 706 Suchanek, F. M., Ifrim, G., and Weikum, G. (2006). "Leila: Learning to extract informa-
707 tion by linguistic analysis." *Proceedings of the 2nd Workshop on Ontology Learning and*

- 708 *Population: Bridging the Gap between Text and Knowledge*, 18–25.
- 709 Toutanova, K., Klein, D., Manning, C. D., and Singer, Y. (2003). “Feature-rich part-of-
- 710 speech tagging with a cyclic dependency network.” *Proceedings of the 2003 Conference of*
- 711 *the North American Chapter of the Association for Computational Linguistics on Human*
- 712 *Language Technology-Volume 1*, Association for Computational Linguistics, 173–180.
- 713 Turney, P. D. and Pantel, P. (2010). “From frequency to meaning: Vector space models of
- 714 semantics.” *Journal of artificial intelligence research*, 37(1), 141–188.
- 715 Walton, C. M., Seedah, D. P., Choubassi, C., Wu, H., Ehlert, A., Harrison, R., Loftus-Otway,
- 716 L., Harvey, J., Meyer, J., Calhoun, J., et al. (2015). *Implementing the freight transportation*
- 717 *data architecture: Data element dictionary*. Number Project NCFRP-47.
- 718 Webster, J. J. and Kit, C. (1992). “Tokenization as the initial phase in nlp.” *Proceedings of*
- 719 *the 14th conference on Computational linguistics- Volume 4*, Association for Computational
- 720 Linguistics, 1106–1110.
- 721 Wetherill, M., Rezgui, Y., Lima, C., and Zarli, A. (2002). “Knowledge management for
- 722 the construction industry: the e-cognos project.” *Journal of Information Technology in*
- 723 *Construction (ITCon)*, 7, 183–196.
- 724 Wikipedia (2016). “Glossary of road transportation terms. Accessed: April 11, 2016.
- 725 Yalcinkaya, M. and Singh, V. (2015). “Patterns and trends in building information modeling
- 726 (bim) research: A latent semantic analysis.” *Automation in Construction*, 59, 68–80.
- 727 Yarowsky, D. (1995). “Unsupervised word sense disambiguation rivaling supervised methods.”
- 728 *Proceedings of the 33rd annual meeting on Association for Computational Linguistics*,
- 729 Association for Computational Linguistics, 189–196.
- 730 Zhang, J. and El-Gohary, N. (2016). “Extending building information models semiautomat-
- 731 ically using semantic natural language processing techniques.” *Journal of Computing in*
- 732 *Civil Engineering*, C4016004.
- 733 Zhang, Z., Iria, J., Brewster, C., and Ciravegna, F. (2008). “A comparative evaluation of
- 734 term recognition algorithms.” *LREC*.

- 735 Zhao, H. and Kit, C. (2011). "Integrating unsupervised and supervised word segmentation:
736 The role of goodness measures." *Information Sciences*, 181(1), 163–183.

737 **List of Figures**

738 1	Overview of the proposed methodology	39
739 2	Linguistic processing procedure to detect NPs	40
740 3	Multi-word term extraction evaluation	41
741 4	Skip-gram model	42
742 5	Roadway term space model (Rd-VSM)	43

TABLE 1: Term candidate filters

Pattern	Examples
(Adj N)*N	road, roadway shoulder, vertical alignment
(Adj N)*N Prep (of/in) (Adj N)*N	right of way, type of roadway

Note: |, * respectively denote ‘and/or’, and ‘zero or more’.

TABLE 2: Excerpts of the extracted candidate terms

Term candidate	Termhood	real term?
sight distance	9435.314	yes
design speed	9052.556	yes
additional information	1829.0	no
typical section	1801.0	yes
basis of payment	1762.478	no

TABLE 3: Skip-gram model parameters

Parameter	Value
Frequency threshold	50-100
Hidden layer size	100-500
Context window size	5,10,15

TABLE 4: Examples of top nearest words

Term	Nearests	Cosine	Rank
roadway	highway	0.588	1
	traveled-way	0.583	2
	roadway-section	0.577	3
	road	0.533	4
	traffic-lane	0.524	5
	separating	0.522	6
	adjacent-roadway	0.519	7
	travel-way	0.517	8
	entire-roadway	0.513	9
...
	roadway-shoulder	0.505	12
	roadway-cross-section	0.491	18
	undivided	0.452	37
	mainline-roadway	0.450	42

TABLE 5: Patterns to extract attributes and hyponyms

Relation	Pattern	Example
Attribute	Noun1 of Noun2	the width of the road
	Noun1 Noun2	road width, project cost
Hypernym-hyponym	Noun1 Noun2	vertical alignment isA alignment

TABLE 6: An example classified list of nearest terms

Term	Relation Group	Nearests	Cosine	Rank
roadway	Synonym	highway	0.588	1
		traveled-way	0.583	2
		road	0.533	4
		traffic-lane	0.524	5
		travel-way	0.517	8
	Attribute	separating	0.522	6
		roadway-section	0.577	3
		roadway-shoulder	0.505	12
		roadway-cross-section	0.491	18
	Hyponym	adjacent-roadway	0.519	7
		entire-roadway	0.513	9
		undivided	0.452	37
		mainline-roadway	0.450	42

TABLE 7: Performance of the synonym matching task with various training settings

Parameter changed	Model	Precision (%)	Recall(%)	F (%)
Baseline	50-100-5	79	53	63
Window size	50-100-10	81	54	65
	50-100-15	81	54	65
Frequency threshold	<u>75</u> -100-5	74	50	60
	<u>100</u> -100-5	77	51	62
Hidden layer size	50- <u>200</u> -5	79	53	63

TABLE 8: Comparison of synonym matching performance between WordNet and proposed system

Lexicon	Precision (%)	Recall(%)	F (%)
Wordnet	76	40	52
Proposed system	81	54	65

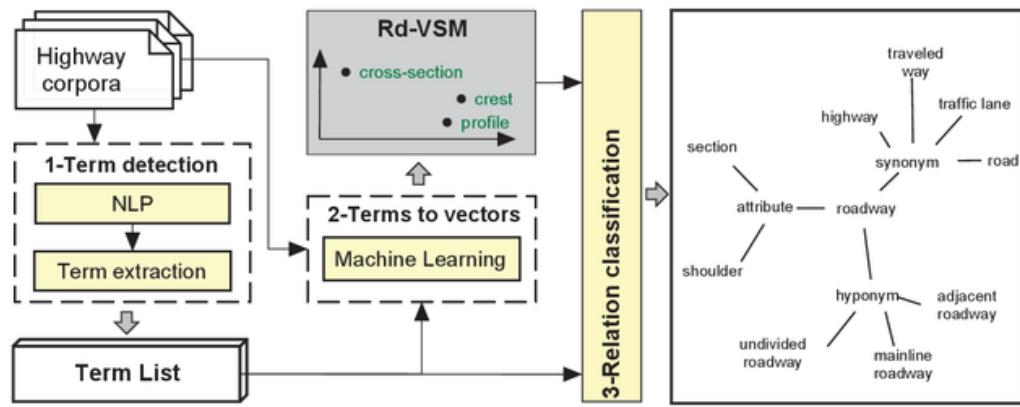


FIG. 1: Overview of the proposed methodology

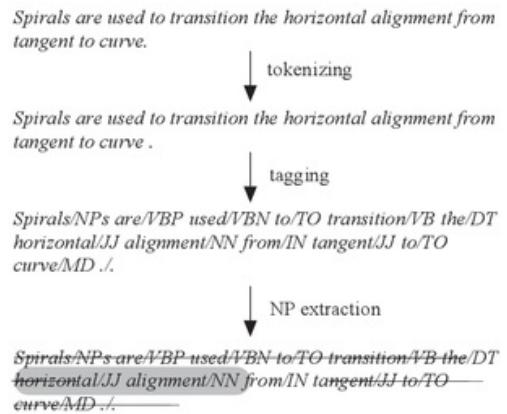


FIG. 2: Linguistic processing procedure to detect NPs

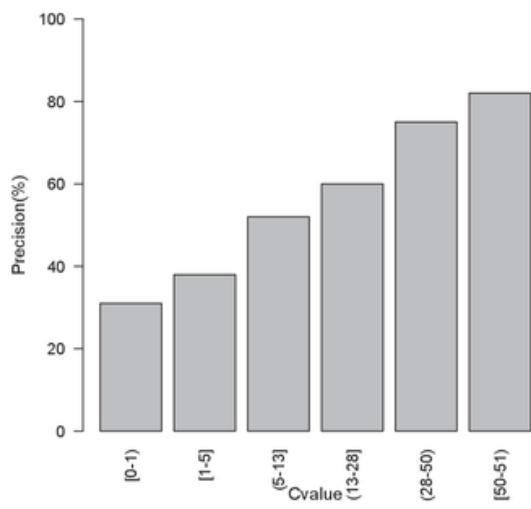


FIG. 3: Multi-word term extraction evaluation

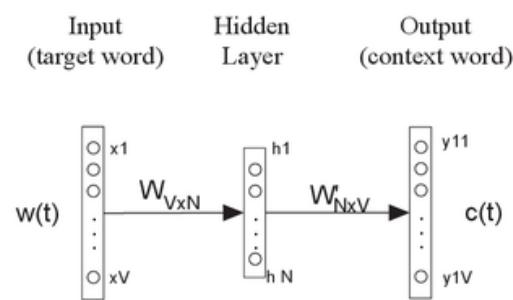


FIG. 4: Skip-gram model

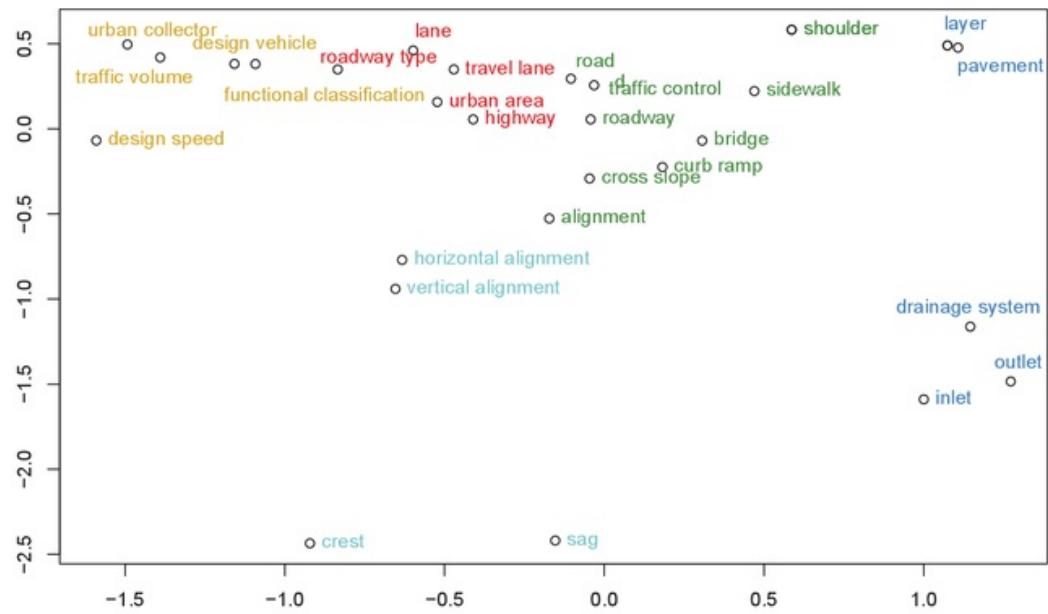


FIG. 5: Roadway term space model (Rd-VSM)

ASCE Journals Sizing Worksheet

Please complete this form for all new manuscripts

September 12, 2016

This worksheet will automatically calculate the total number of printed pages your article will occupy in the journal.

Please fill in all fields in green below. If you do not know your Manuscript Number, you may leave that field blank.

Length Limits:

Technical Paper/Case Study = 8 pgs. Forum = 4 pgs.

Technical Note = 3 pgs. Discussion/Closure = 2 pgs.

Manuscript number:	CPENG-1952
Journal name:	Journal of Computing in Civil Engineering
Corresponding author name:	H. David Jeong
Email address:	djeong@iastate.edu

*Information on the maximum allowed length for each article type can be found online at:
<http://www.asce.org/Content.aspx?id=29559>*

Number of pages in your manuscript: 29

- Please include figure captions when indicating the size of your manuscript.

Number of figure pages: 5

- Manuscripts should use 12 pt. font, double spaced, with 1 inch margins.

Number of table pages: 8

Estimated article pages: 10

Note: The total displayed above is only an estimate. Final page count will depend on a number of factors, including the size of your figures and tables, and the number of display equations in your manuscript.

*Additional author resources can be found online using the ASCE Author Guide located at:
<http://www.asce.org/Content.aspx?id=18107>*

1
2
3
4
5
6
7

Response to Reviewers' Comments

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Manuscript title: NLP-based approach to classifying heterogeneous terms for unambiguous exchange of roadway data

Manuscript #: CPENG-1952

The authors would like to thank the editor and reviewers for their time and effort in reviewing our manuscript. We hope that the revised manuscript is suitable for publication and we look forward to your response.

Editor's comments

1. "The manuscript has received two reviews. Both reviewers see the potential for publication, but have significant concerns (especially about the novelty and contribution). The authors are, thus, requested to revise the manuscript based on the reviews and resubmit it for another round of review. Please respond point-by-point to each review comment. In responding to a comment, please note any changes made in the manuscript and refer to the respective line numbers. When revising the manuscript, please better highlight your contribution, better justify the advancement in the body of knowledge (especially in comparison to recently published papers), and better describe the implications of the work and how successful the project is."

Response: The authors have revised the manuscript in accordance with the recommendations from the reviewers. Please see the responses below to each reviewer's comments. With respects to highlighting the research contribution, we have revised the manuscript with the aim of providing readers with a clearer discussions and justification for the research value. The following are examples of revised paragraphs discussing the contribution of the research throughout the revised manuscript.

Introduction section, page 3, line 65-75, "Terminology transparency through digital dictionaries like glossaries, taxonomies, ontologies and data dictionaries is identified as a driver of semantic interoperability (Ouksel and Sheth 1999). Although a plethora of semantic resources have been introduced to the highway sector, their coverages of terms are still far inadequate and the inclusion of multiple names for the same concept is limited. This is because of the reliance on a tedious and time-consuming approach which requires developers to manually gather and translate knowledge from domain experts or text documents into a machine-readable format. Thus, there is a need for computer-aided methods to remove this knowledge acquisition bottleneck (Mounce et al. 2010), so that digital dictionaries can be quickly constructed to meet a specific need and to keep up with the growth of new terms due to rapid applications of new technologies and knowledge."

Related studies section, page 10, line 246-258, "As shown in the literature review, there are numerous research efforts in developing ontologies for the highway sector. However, the existing ontologies are mainly hand-coded through manual processes of knowledge

1
2
3
4 acquisition and formally describing them in a digital format. This ad-hoc approach has
5 created a bottleneck in facilitating the semantic interoperability for the whole industry
6 and as a result, semantic resources for many aspects of a project are still not available. A
7 few efforts have been made to automate the process of constructing or extending existing
8 semantic resources. The most rigorous methodology in the state-of-the-art is the one
9 developed by Zhang and El-Gohary (2016) that is fully-automated with high accuracy. One
10 limitation of this algorithm is the reliance on an existing semantic resource; it, therefore,
11 would not be applicable to such a domain like infrastructure that is out of the vocabulary
12 scope. Thus, there is a need for an automated approach that can not only allow for fast
13 development of highway lexicons but also remove the dependence on other existing
14 semantic models.”
15
16

17
18 Discussions section, Page 21, line 508-519, “This paper proposes an NLP based
19 methodology to assist professionals in extracting roadway terms and their semantic
20 relations from text documents. A key contribution to the body of knowledge is the novel
21 approach with a new algorithm that allows for automated detection of technical terms and
22 their relations without reliance on existing hand-coded dictionaries as used by previous
23 researchers such as Zhang and El-Gohary (2016). The present framework is not to
24 completely eliminate the human involvement, but is expected to become an enabling tool
25 that can help researchers in the domain quickly develop supporting ontologies and other
26 forms of semantic resources for their specific use cases. With respect to facilitating
27 semantic interoperability for the infrastructure sector, the findings of this study would
28 accelerate the process of removing the current bottleneck in extensive machine readable
29 dictionaries which are required for an unambiguous data sharing, integration or
30 exchange.”
31
32

33 **Reviewer #1:**
34
35

- 36 1. “There are a number of small grammatical errors in the paper and it would benefit from a final
37 English language check (e.g., use of ‘to the same’ rather than ‘of the same’).”
38

39 *Response: Thank you for your comment. The authors have gone through the manuscript
40 and corrected the grammar and spelling errors.*
41

- 42 2. “On page 3 you state that research to address the issue of terminology inconsistency has been
43 very limited, but you also point to major decade long efforts such as bSDD. I think you are more
44 concerned with research into automated approaches, and this should be made explicit here.”
45

46 *Response: Yes, this manuscript focuses on the automated approach. The authors have
47 revised the manuscript to highlight the research purpose with higher clarity, as follows.*
48

49 *Page 3, line 65-75, “Terminology transparency through digital dictionaries like
50 glossaries, taxonomies, ontologies and data dictionaries is identified as a driver of
51 semantic interoperability (Ouksel and Sheth 1999). Although a plethora of semantic
52 models exist, they are not yet widely adopted in the industry due to the lack of standardization
53 and interoperability between different systems. Therefore, there is a need for an automated
54 approach that can facilitate the development of semantic resources for the infrastructure
55 sector. This paper proposes an NLP based methodology to assist professionals in extracting
56 roadway terms and their semantic relations from text documents. The methodology is
57 novel and automated, allowing for the detection of technical terms and their relations without
58 reliance on existing hand-coded dictionaries. The findings of this study would accelerate
59 the process of removing the current bottleneck in extensive machine readable dictionaries
60 which are required for an unambiguous data sharing, integration or exchange.”*
61
62

1
2
3
4 resources have been introduced to the highway sector, their coverages of terms are still
5 far inadequate for a large number of disciplines, and processes across the highway project
6 life cycle. This is because of the reliance on a tedious and time-consuming approach which
7 requires developers to manually gather and translate knowledge from domain experts or
8 text documents into a machine-readable format. Thus, there is a need for computer-aided
9 methods to remove this knowledge acquisition bottleneck (Mounce et al. 2010), so that
10 digital dictionaries can be quickly constructed to meet a specific need and to keep up with
11 the growth of terms due to rapid applications of new technologies and knowledge.”
12
13
14

- 15
16 3. “On page 4 you present no evidence that the growth of terms is exponential - and it is hard to
17 imagine that this could be the case.”
18

19 *Response: We have toned down the sentence. The revised sentences now reads*
20
21

22 “....digital dictionaries can be quickly constructed to meet a specific need and to keep up
23 with the growth of terms due to rapid applications of new technologies and knowledge.”
24
25
26
27

- 28 4. “On page 6 you could present more explanation as to why reliance on digital dictionaries is
29 becoming a bottleneck.”
30

31 *Response: The authors have included the following information and reference to explain
32 why reliance on digital dictionaries is becoming a bottleneck.*
33
34

35 *Page 5, line 126-128, “However, digital dictionaries are typically hand-crafted; they are
36 therefore not available to many domains (Kolb 2008).”*
37
38

- 39 5. “On page 7 it would be useful to present some information on the performance of the techniques
40 and system that you mention here.”
41

42 *Response: The authors have added several sentences to further discuss the performance of
43 those methods mentioned in the manuscript, as follows.*
44
45

46 *Page 7, line 158-164, “For example, the results from a comparative study conducted by
47 Levy et al. (2015) on the accuracy in various tasks and golden standards reveal that Skip-
48 gram outperforms Glove in every experiment and is the winner in most of the tasks,
49 especially on the WordSim Similarity dataset. Among these tasks, the best precision of Skip-
50 gram is .793, while PPMI and Glove achieve the highest score of .755 and .725
51 respectively.”*
52
53
54

- 55 6. “On page 8 isn't it more important whether efforts like bSDD are complete rather than how long
56 they take to create?”
57
58
59
60
61
62
63
64
65

1
2
3
4 *Response: The focus of this paper is to reduce the reliance on human efforts in developing*
5 *dictionaries such as bSDD. The authors have revised the paragraph with new sentences to*
6 *better describe the idea, as follows.*
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

Page 7 and 8, line 175-196, “A popular solution to semantic interoperability is to develop taxonomies, ontologies or other forms of digital dictionaries that can provide machine-readable definitions of domain concepts. A plethora of such semantic resources have been developed for the highway industry. However, conventional methods require significant human efforts on knowledge retrieval, ontology construction and validation. The pioneer in this line of research is the e-COGNOS ontology (Wetherill et al. 2002; Lima et al. 2005) which formulates the execution process of a construction project as an explicitly interactive network of the principal concepts: Actor, Resources, Products, Processes and Technical Topics. The ontology developers of this project reviewed existing taxonomies (BS61000, UniClass, IFC) and construction specific documents, and interacted with the end users to identify relevant concepts and their semantic relations. Industry experts were invited to validate the developed ontology through questionnaires on concept names and relations. Since the introduction of the high-level ontology of e-Cognos, a plenty of ontologies have been built for various aspects of the life cycle of a highway project, for instance, highway construction taxonomy (El-Diraby and Kashif 2005; El-Diraby et al. 2005), freight ontology (Seedah et al. 2015a), and the ontology of urban infrastructure products (Osman and El-Diraby 2006). Like the e-Cognos project, these studies also relied on domain experts for the construction of their semantic products. The limitation regarding time and labor costs of the ad-hoc traditional methodology has created a bottleneck to the progress in enabling semantic interoperability. In addition, the existing ontologies primarily focus on the description of concepts, and the heterogeneity of concept names is usually neglected. Therefore, research is needed not only to automate the process of formulating domain concepts but also to incorporate term heterogeneity into ontologies.”

40
41 7. “On page 9 how do we know that the coverage of this corpora is sufficient and complete? How
42 about terms from corpora in other English language speaking countries? Or corpora from countries
43 which use languages other than English? Need to scope your work better here.”
44
45
46
47
48
49

Response: Thank you for your comment. The present framework is limited to American-English text documents. The authors have revised the associated sentence regarding this issue to clearly define the scope of the work, as follows.

50
51 *Page 10, line 261-263, “The goal of this research is to propose an NLP-based methodology*
52 *that can automate the process of extracting roadway technical terms and their semantic*
53 *relations from American-English roadway documents.”*
54
55
56
57

58 8. “On page 10 you should discuss the implication of removing tables and equations from your
59 corpus.”
60
61
62
63
64
65

1
2
3
4 *Response: The authors have included the following sentences to discuss the implication of*
5 *removing tables and equations.*

6
7 *Page 11, line 277-280, "The removal of these features may slightly reduce the corpus size,*
8 *and accordingly affect the training dataset; however, it is necessary since words in tables*
9 *and equations are not organized in the formal structure of a sentence and therefore the*
10 *NLP algorithm may extract unreal noun phrases."*

11
12
13 9. "On page 14 you should explain what a 'one-hot vector' is."

14
15 *Response: The authors have included the following sentences to explain the 'one-hot'*
16 *vector.*

17
18 *Page 16, line 405-408, "In this model, a word in the corpus vocabulary is encoded as a*
19 *"one-hot" vector which is a vector in which only one element at the index of the word in*
20 *the vocabulary is set one, and all other items are zero. For example, the one-hot vector of*
21 *kth word in the vocabulary with the size of V will be {x₁=0, x₂=0, ..., x_k=1, ..., x_V=0}."*

22
23 10. "In the discussion and conclusions there is no consideration of what is good enough in a
24 lexicon. Is your 81% precision sufficient? Is a F-measure of 65% sufficient?"

25
26 *Response: The authors understand that the current performance is not high enough for a*
27 *fully automated process, when 35 percent of the tested terms are out of the vocabulary and*
28 *20% of the answers are incorrect. Therefore, the authors have suggested that the most*
29 *immediate value of this system is significant reduction of human efforts rather than*
30 *completely removing the human involvement. Users can use this proposed framework as*
31 *an assistant tool and the automated result may need a manual review. The following*
32 *statements have been included in the discussions and conclusions sections of the*
33 *manuscript to discuss the above issues.*

34
35 *Discussions section, page 21, line 511-516, "The present framework is not to completely*
36 *eliminate human involvement, but is expected to become an enabling tool that can help*
37 *researchers in the domain quickly develop supporting ontologies and other forms of*
38 *semantic resources for their specific use cases."*

39
40 *Conclusions section, page 23, line 569-573, "Although a significant improvement has*
41 *been made in comparison with an existing thesaurus database, the overall performance is*
42 *not relatively high. This might be due to the size of the training data. Future research will*
43 *be conducted to expand the highway corpus to further disciplines such as asset*
44 *management, and transportation operation."*

45
46 11. "It would be useful to have the final lexicon available online somewhere so that readers of the

47 paper can access and assess the results of this work."

1
2
3
4 *Response: The authors have included a link where the model and datasets from this study*
5 *can be found at the end of the introduction section, as follows.*
6
7
8
9
10

11 *Page 4, line 94-95. "A Java package and several datasets resulting from the study can be*
12 *found at <https://github.com/tuyenbk/mvdgenerator>."*
13

14 **Reviewer #2:**
15
16

- 17 1. "The authors might want to revise the manuscript title to better reflect the scope of research.
18 Particularly, the text corpus came from only roadway design guidelines and does not constitute the
19 entirety of civil infrastructure. While the manuscript places an equal emphasis on the approach as
20 well as on the research product (i.e. InfraLex), the manuscript title should be more carefully crafted
21 into something more accurate"

22 *Response: Thank you for your comment. The authors have revised the title to better reflect*
23 *the research focus which is an automated approach, and the scope within the roadway*
24 *sector as follows "NLP-based approach to classifying heterogeneous terms for*
25 *unambiguous exchange of roadway data."*
26

- 27 2. "Several statements within INTRODUCTION need to be supported.
28 2.1 Page 2, Lines 38-40, "The major cost was time spent...a useful format."
29 2.2 Page 2, Lines 48-49, "Polysemy and synonymy are two ...data sources."
30 2.3 Page 3, Lines 61-63, "However, research to address...very limited."'"
31

32 *Response: The authors have included references to support the above arguments. The*
33 *following are the revised sentences.*
34
35

36 *Page 2, Line 36-40, "The inadequate interoperability cost is estimated to be over \$15.8*
37 *billion per year in the U.S. capital facilities industry as reported by the National Institute*
38 *of Standard and Technology (NIST); and the largest cost item is the laborious work for*
39 *finding, verifying, and transferring facility and project information into a useful format*
40 *during the operation and maintenance stage (Gallaher et al. 2004)."*
41
42

43 *Page 3, Line 52-54, "Polysemy and synonymy are two major linguistic obstacles to*
44 *semantic integration and use of a multitude of data sources (Noy 2004)."*
45
46

47 *Regarding the last statement (page 3, lines 61-63), this is based on the literature review*
48 *conducted by the authors. The following is the revised statement.*
49
50

51 *Page 3, Line 67-69, "Although a plethora of semantic resources have been introduced to*
52 *the highway sector, as shown in the literature review, their coverages of concepts are still*
53 *inadequate and the inclusion of multiple names to the same concept is still limited."*
54
55

- 56 3. "The statement citing the cost of interoperability issue on Page 2 (between lines 35 - 38) does
57 not necessarily apply to the target of the paper, civil infrastructure (or more precisely roadway
58
59

1
2
3
4 design). The citation from NIST was about the capital facilities, which typically refer to buildings
5 and industrial facilities and not horizontal construction.”
6
7

8 *Response: Since the building and roadway sectors share the same interoperability issues,*
9 *the intention of the authors here is to provide a tangible evidence about the cost of*
10 *inadequate interoperability in order to highlight the importance of addressing the same*
11 *issue in the highway sector. The authors have revised the manuscript as follows to explain*
12 *that idea.*

13
14 *Page 2, Line 34-46, “The interoperability issue has been widely recognized as a key*
15 *obstacle blocking the flow of digital data throughout the entire project life cycle. The*
16 *inadequate interoperability cost is estimated to be over \$15.8 billion per year in the U.S.*
17 *capital facilities industry as reported by the National Institute of Standard and Technology*
18 *(NIST); and the largest cost item is the laborious work for finding, verifying, and*
19 *transferring facility and project information into a useful format during the operation and*
20 *maintenance stage (Gallaher et al. 2004). This finding indicates that the lack of readiness*
21 *for downstream phases to directly use the transferred digital project data generated from*
22 *upstream design and construction stages results in high operational costs. Since the*
23 *roadway sector, which is one of the major domains in the construction industry, has not*
24 *yet successfully facilitated a high degree of interoperability (Lefler 2014); huge cost*
25 *savings would be achieved if roadway data is seamlessly shared across project phases and*
26 *among state and local agencies.*

27
28 4. “The authors seemed to suggest that very few research looked into the terminology
29 inconsistency issues in construction whereas a recent paper published in JCCE attempted to
30 address this particular issue in the transportation sector. The paper is entitled “Ontology for
31 Querying Heterogeneous Data Sources in Freight Transportation” and the authors are advised to
32 review it in order to highlight their research uniqueness and contribution.”

33
34 *Response: Thank you for your recommendation. The authors have reviewed this*
35 *recommended article and the following discussion on the gap addressed in that publication*
36 *and others has been added to the revised manuscript.*

37
38 *Page 8, Line 186-196, “Since the introduction of the high-level ontology of e-Cognos, a*
39 *plenty of ontologies have been built for various aspects of the life cycle of a highway*
40 *project, for instance, highway construction taxonomy (El-Diraby and Kashif 2005; El-*
41 *Diraby et al. 2005), freight ontology (Seedah et al. 2015a), and the ontology of urban*
42 *infrastructure products (Osman and El-Diraby 2006). Like the e-Cognos project, these*
43 *studies also relied on domain experts for the construction of their semantic products. The*
44 *limitation regarding time and labor costs of the ad-hoc traditional methodology has*
45 *created a bottleneck to the progress in enabling semantic interoperability. In addition, the*
46 *existing ontologies primarily focus on the description of concepts, the heterogeneity of*
47 *concept names is usually neglected. Therefore, research is needed not only to automate the*

1
2
3
4 *process of formulating domain concepts but also to incorporate term heterogeneity into*
5 *ontologies.*"
6
7
8
9
10
11
12
13
14

5. "The authors used a single source of citation from arXiv to support their claim that Skip-Gram
model outperforms other methods like LSA and so was adopted in the reported research. This is a
critical research design decision and should be more carefully evaluated (and supported). arXiv
does not require peer reviews and its publications are only moderated to check against obvious
policy violations."

15
16 *Response: Thank you for your comment. The authors have added other evidences from*
17 *several peer-review publications to support the selection of this research tool. Below is the*
18 *revised discussion explaining the tool selection.*

19
20 *Page 6-7, Line 154-164, "There are contradictive recommendations on the wining model*
21 *in the literature. The authors of Glove suggested that their model out-performs over Skip-*
22 *Gram and others in the state of the art. However, a number of independent benchmarking*
23 *experiments have consistently indicated the outperformance of the Skip-gram model to its*
24 *alternatives. For example, a comparative study conducted by Levy et al. (2015) on the*
25 *accuracy in various tasks and golden standards reveals that Skip-gram outperforms Glove*
26 *in every experiment and is the winner in most of the tasks, especially on the WordSim*
27 *Similarity dataset. Among these tasks, the best precision of Skip-gram is .793, while PPMI*
28 *and Glove achieve the highest score of .755 and .725 respectively. The out-performance*
29 *of Mikolov's model on the similarity task is confirmed in another benchmarking study (Hill*
30 *et al. 2015) where this model is also found as the winner in most of the tests."*

31
32
33
34
35 6. "How do the results from the step 'Noun phrase detection' look like? Porter stemming (or any
36 stemming algorithm) is known to over or under stem terms at times. How was over or under
37 stemming managed when the authors were generating their bag of noun phrases? Why wouldn't
38 the suffix in the word "undivided" in Table 6 be removed after the authors employed the Porter
39 stemming algorithm?"

40
41 *Response: The authors indeed implemented the Pling Stemmer rather than Porter, but the*
42 *initial manuscript did not update accordingly." Thank you for capturing this inconsistency.*
43 *Since the Pling stemmer specifically stems English plural nouns to its singular form, other*
44 *POS such as adjectives will not be affected. This is the reason why the word "undivided"*
45 *is not affected. In regards to the errors of over and under-stemming, there are evidences*
46 *from the literature review that suggest that even though there are still errors in the existing*
47 *stemmers, they are good enough to not have negative effects on the overall performance of*
48 *NLP applications. Below is the revised discussion on the stemming process in the proposed*
49 *method.*

50
51
52
53
54
55 *Page 12, Line 311-319, "Stemming is a popular process to reduce words to their stems.*
56 *Despite the fact that none of the existing algorithms can completely eliminate the errors of*
57 *over and under stemming, they are good enough to not degrade the overall performance*

1
2
3
4 of NLP application (Jivani et al. 2011). This study implements the Pling stemmer
5 (Suchanek et al. 2006), which stems an English noun to its singular form, to normalize
6 plural nouns in the corpus. One advantage of this algorithm is the utilization of both
7 syntactic rules and the vocabulary in a dictionary; hence the mis- or over-stemming errors
8 that take off a true suffix can be reduced.”
9
10

11 7. “A manual evaluation process was still employed to remove inadequate or meaningless terms
12 from the term candidate list. Could this process become the road block when the authors need to
13 scale up their research?”
14
15

16 *Response: The authors understand that the proposed framework is to reduce human efforts
17 in developing semantic resources rather than completely replace domain experts.
18 Manually reviewing of terms would still be a much easier and less time-consuming task
19 than reviewing a plethora of domain text documents. In addition, the authors suggest a
20 method to reduce both the laborious work and the number of real terms automatically
21 removed (the paragraph discussing this method is included below). Using this method, a
22 domain expert does not need to review the entire list. In our experiment, it took a graduate
23 student only 8 hours to finish his task. However, further research is needed to enhance the
24 accuracy of term extraction, so that human involvement would continue to be reduced.*
25
26
27
28
29

30 **Page 14-15, Line 362-376,** “To automatically remove candidates that are unlikely to be
31 real terms, a threshold C-value can be used. However, doing this may eliminate the real
32 terms that appear in the bottom due to their low frequencies. Manual evaluation of the
33 entire candidate list would avoid the removal of real terms with low C-values. To minimize
34 both laborious work and the number of true terms wrongly discarded, the authors suggest
35 the following method to identify the threshold value. The ranked list of candidates is divided
36 into groups of around 200 items. A graduate student with a civil engineering background
37 was asked to utilize a bottom-up approach to evaluate group by group and stop at which
38 the percentage of actual terms achieved 80 percent. Users can choose a higher percentage
39 limit in cases where the accuracy is critical. Table 2 illustrates the evaluation results for
40 several excerpts of the extracted term candidates. The precision values, which represent
41 the percentages of real terms in these groups, are presented in Figure 3. As shown in the
42 figure, precision values are less than 80 percent for groups with c-values less than 50. This
43 value is set as the threshold for the acceptance of term candidates. The final selected list
44 is comprised of nearly 6,000 multi-word roadway technical terms.”
45
46
47
48
49
50

51 8. “How were the precision rates in Figure 3 determined? These rates seemed to provide critical
52 information for the manual evaluation upon the term candidates and it is unclear how these
53 precision rates were obtained prior to the manual evaluation.”
54

55 *Response: The precision rate is the percentage of real terms in a group of automatically
56 extracted terms. The conclusion of “real” or “unreal” is based on the manual evaluation
57 process. Below is the revised discussion on the evaluation method and how the precision
58 rates are determined.*
59
60
61
62
63
64
65

1
 2
 3
 4 *Page 15, Line 367-373, "The ranked list of candidates is divided into groups of around*
 5 *200 items. A graduate student with a civil engineering background was asked to utilize a*
 6 *bottom-up approach to evaluate group by group and stop at which the percentage of actual*
 7 *terms achieved 80 percent. Users can choose a higher percentage limit in cases where the*
 8 *accuracy is critical. Table 2 illustrates the evaluation results for several excerpts of the*
 9 *extracted term candidates. The precision values, which represent the percentages of real*
 10 *terms in these groups, are presented in Figure 3."*
 11
 12
 13
 14
 15

TABLE 2: Excerpts of the extracted candidate terms

Term candidate	Termhood	real term?
sight distance	9435.314	yes
design speed	9052.556	yes
additional information	1829.0	no
typical section	1801.0	yes
basis of payment	1762.478	no

29
 30 9. "Results in Table 7 are not sensitive to the different parameter settings. Table 7 can be removed
 31 from the manuscript."
 32
 33
 34
 35
 36 *Response: Thank you. The authors have also recognized that there is no consistent relation*
 37 *pattern between the performance and the configuration of these parameters. However, the*
 38 *authors would like to respectfully keep this table in the manuscript to visualize that finding.*
 39
 40

TABLE 7: Performance of the synonym matching task with various training settings

Parameter changed	Model	Precision (%)	Recall(%)	F (%)
Baseline	50-100-5	79	53	63
Window size	50-100-10	81	54	65
	<u>50-100-15</u>	81	54	65
Frequency threshold	<u>75-100-5</u>	74	50	60
	<u>100-100-5</u>	77	51	62
Hidden layer size	50-200-5	79	53	63

52
 53 10. "The body of knowledge does not lie within the method and was more on the research product.
 54 This is a relatively weak contribution"
 55
 56
 57 *Response: The authors have revised the literature review section to focus more on the state-*
 58 *of-the-art approaches and discussions on the contribution of this manuscript to the body*
 59 *of knowledge. Below is the revised body of knowledge section.*
 60
 61
 62
 63
 64
 65

1
2
3
4
5
6
7

Page 7-10, Line 175-258,

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

"Related studies

A popular solution to semantic interoperability is to develop taxonomies, ontologies or other forms of digital dictionaries that can provide machine-readable definitions of domain concepts. A plethora of such semantic resources have been developed for the highway industry. However, conventional methods require significant human efforts on knowledge retrieval, ontology construction and validation. The pioneer in this line of research is the e-COGNOS ontology (Wetherill et al. 2002; Lima et al. 2005) which formulates the execution process of a construction project as an explicitly interactive network of the principal concepts: Actor, Resources, Products, Processes and Technical Topics. The ontology developers of this project reviewed existing taxonomies (BS61000, UniClass, IFC) and construction specific documents, and interacted with the end users to identify relevant concepts and their semantic relations. Industry experts were invited to validate the developed ontology through questionnaires on concept names and relations. Since the introduction of the high-level ontology of e-Cognos, a plenty of ontologies have been built for various aspects of the life cycle of a highway project, for instance, highway construction taxonomy (El-Diraby and Kashif 2005; El-Diraby et al. 2005), freight ontology (Seedah et al. 2015a), 190 and the ontology of urban infrastructure products (Osman and Ei-Diraby 2006). Like the e-Cognos project, these studies also relied on domain experts for the construction of their semantic products. The limitation regarding time and labor costs of the ad-hoc traditional methodology has created a bottleneck to the progress in enabling semantic interoperability. In addition, the existing ontologies primarily focus on the description of concepts, the heterogeneity of concept names is usually neglected. Therefore, research is needed not only to automate the process of formulating domain concepts but also to incorporate term heterogeneity into ontologies.

40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Another strategy for semantic interoperability targets at the heterogeneity of concept names rather the concept description as in an ontology model. A few frameworks to assist practitioners in precisely mapping data labels from heterogeneous sources have been introduced for various construction sectors. In the building sector, buildingSMART proposed a novel framework, namely IFD (International Framework for Dictionaries) (ISO 12006-3) for developing a multilingual data schema in which each concept can have multiple names in different languages. With IFD, the identity of a concept is defined by a Global Unique ID (GUID) rather than its name; hence an IFD-based data exchange mechanism is able to eliminate the semantic mismatches due to the name inconsistency (IFD Library Group ; Hezik 2008).. The buildingSMART data dictionary (bSDD) (buildingSMART 2016) is the first digital library of building concepts that is crafted in the IFD structure. Each concept in bSDD consists a set of synonymy names not only in English but also in computer-coded languages (e.g., IFC-Industry Foundation Classes) and in other human languages (e.g., French, Norwegian). Therefore, a complete bSDD would enable digital data in regardless of languages to be sharable and unambiguously reusable. Yet, its size remains limited as the identification of these sets of synonyms is labor and time

extensive. In the transportation sector, there has been a shortage of research efforts targeting the heterogeneity of data names at the database level until recently. Seedah et al. (2015b) proposed a role-based classification schema (RBCS) to classify data in freight databases. RBCS defines nine distinct groups of roles that are time (year, month), place (city name, population), commodity (liquid, value), link (roadway name, width), mode (truck, rail), industry (company name, sales), event (accident, number of fatalities), and human (officer, driver age). The authors argue that once the data elements across separate databases are categorized using this standard system, it becomes easier for practitioners to identify the semantic relatedness in their definitions. However, even if RBCS is successfully applied to all freight databases, identifying the exact type of relation (synonym, functional relation) between two data elements in the same category is still a challenging task.

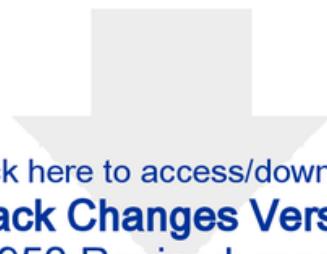
In attempts to reduce laborious work on defining concepts, a few researchers have sought to propose semi-automated and automated methods for identifying semantic relations among technical terms. Abuzir and Abuzir (2002) developed the ThesWB system which utilizes hand-coded syntax patterns to detect lexical relations between civil engineering terms from HTML web pages. The performance of ThesWB was not reported, but it is not likely to be high since rule-based approaches are repeatedly criticized for not being able to capture all the variant ways to present relations among terms in natural language (Marcus 1995; Navigli and Velardi 2010). Rezgui (2007) suggested a more sophisticated approach that is based the statistics of word occurrence rather than predefined rules to extract potential pairs of related terms from domain text documents. This method implements TF-IDF to evaluate the importance degree of a keyword to the examined domain; and analyzes the co-occurrence frequencies using Metric Clusters to assess the potentiality that exists a semantic relation within a given pair of important keywords. These potential relationships are then validated and categorized by domain experts. Since only pairs of terms that occur in the same sentence are considered, equivalent terms which are used interchangeably could not be captured. In another study to identify semantic relations, Zhang and El-Gohary (2016) proposed a fully automated methodology for both tasks of retrieving related candidate and classifying the relations. This algorithm was reported to achieve an average precision of nearly 90 percent in the relation classification task. However, the algorithm identifies potentially related concepts based on the pre-defined lexical relations provided in WordNet, a generic lexicon that lacks concepts in many construction sectors including the civil infrastructure, it would not be scalable well on matching terms in these domains.

As shown in the literature review, there are numerous research efforts in developing ontologies for the highway sector. However, the existing ontologies are mainly hand-coded through manual processes of knowledge acquisition and formally describing them in a digital format. This ad-hoc approach has created a bottleneck in facilitating the semantic interoperability for the whole industry and as a result, semantic resources for many aspects of a project are still not available. A few efforts have been made to automate the process

1
2
3
4 *of constructing or extending existing semantic resources. The most rigorous methodology*
5 *in the state-of-the-art is the one developed by Zhang and El-Gohary (2016) that is fully-*
6 *automated with high accuracy. One limitation of this algorithm is the reliance on an*
7 *existing semantic resource; it, therefore, would not be applicable to such a domain like*
8 *infrastructure that is out of the vocabulary scope. Thus, there is a need for an automated*
9 *approach that can not only allow for fast development of highway lexicons but also remove*
10 *the dependence on other existing semantic models."*

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Track Changes Version



Click here to access/download

Track Changes Version
CPENG-1952 Revised_manuscript.pdf

CPENG-1952.pdf

ORIGINALITY REPORT

9%

SIMILARITY INDEX

PRIMARY SOURCES

- | | | |
|---|--|-----------------|
| 1 | aclweb.org
Internet | 157 words — 2% |
| 2 | "A Framework for the Inter-Connection of Life-Cycle Highway Data Islands", Construction Research Congress 2016, 2016.
Crossref | 88 words — 1% |
| 3 | Evangelos Miliros. "Narrative text classification for automatic key phrase extraction in web document corpora", Proceedings of the seventh ACM international workshop on Web information and data management - WIDM 05 WIDM 05, 2005
Crossref | 33 words — < 1% |
| 4 | Lecture Notes in Computer Science, 2016.
Crossref | 26 words — < 1% |
| 5 | aaai.org
Internet | 26 words — < 1% |
| 6 | Lecture Notes in Computer Science, 2010.
Crossref | 26 words — < 1% |
| 7 | Lecture Notes in Computer Science, 2013.
Crossref | 25 words — < 1% |
| 8 | www.diva-portal.org
Internet | 25 words — < 1% |
| 9 | Zhou, Zhipeng, Yang Miang Goh, and Lijun Shen. "Overview and Analysis of Ontology Studies | 19 words — < 1% |

Supporting Development of the Construction Industry", Journal of Computing in Civil Engineering, 2016.

Crossref

-
- 10 www.esrl.noaa.gov Internet 18 words — < 1 %
- 11 Periñán-Pascual, Carlos, and Eva M. Mestre-Mestre. "DEXTER: Automatic Extraction of Domain-Specific Glossaries for Language Teaching", Procedia - Social and Behavioral Sciences, 2015. Crossref 17 words — < 1 %
- 12 www.deepsky.com Internet 14 words — < 1 %
- 13 etheses.whiterose.ac.uk Internet 14 words — < 1 %
- 14 www.ait.edu.gr Internet 12 words — < 1 %
- 15 Zhang, Jiansong, and Nora M. El-Gohary. "Semantic NLP-Based Information Extraction from Construction Regulatory Documents for Automated Compliance Checking", Journal of Computing in Civil Engineering, 2013. Crossref 11 words — < 1 %
- 16 www.dsi.uniroma1.it Internet 11 words — < 1 %
- 17 Kahya-Ozyirmidokuz, E.. "Analyzing unstructured Facebook social network data through web text mining: a study of online shopping firms in Turkey", Information Development, 2014. Crossref 11 words — < 1 %
- 18 cora.ucc.ie Internet 11 words — < 1 %
- 19 www.umiacs.umd.edu Internet 10 words — < 1 %

- 20 Wächter, Thomas (Prof. Dr. Lawrence Hunter, Prof. Dr. Michael Schroeder and Technische Universität Dresden, Fakultät Informatik). "Semi-automated Ontology Generation for Biocuration and Semantic Search", Saechsische Landesbibliothek- Staats- und Universitaetsbibliothek Dresden, 2011.
Publications
-
- 21 Harris, M.R.. "A term extraction tool for expanding content in the domain of functioning, disability, and health: proof of concept", Journal of Biomedical Informatics, 200308/10
Crossref
-
- 22 Lecture Notes in Computer Science, 2008.
Crossref
- 10 words — < 1%
-
- 23 Dra. Marlen Pérez Díaz. "Computer Simulations for Calculation of Workload Weighed Transmission Curves of Brazilian Shielding Materials", IFMBE Proceedings, 2008
Crossref
- 9 words — < 1%
-
- 24 Lecture Notes in Computer Science, 2012.
Crossref
- 9 words — < 1%
-
- 25 conservancy.umn.edu
Internet
- 9 words — < 1%
-
- 26 eprints.mdx.ac.uk
Internet
- 9 words — < 1%
-
- 27 www.tdx.cat
Internet
- 9 words — < 1%
-
- 28 phd.lib.uni-corvinus.hu
Internet
- 9 words — < 1%
-
- 29 Losh, Molly Gordon, Peter C.. "Quantifying narrative ability in autism spectrum disorder: a computational linguistic analysis of na", Journal of Autism and Developmental
- 9 words — < 1%

30	Lecture Notes in Computer Science, 2011. Crossref	9 words — < 1%
31	Lecture Notes in Computer Science, 2014. Crossref	9 words — < 1%
32	archive.org Internet	8 words — < 1%
33	saffron.insight-centre.org Internet	8 words — < 1%
34	www2.denizyuret.com Internet	8 words — < 1%
35	www.aclweb.org Internet	8 words — < 1%
36	Spasić, Irena, Mark Greenwood, Alun Preece, Nick Francis, and Glyn Elwyn. "FlexiTerm: a flexible term recognition method", Journal of Biomedical Semantics, 2013. Crossref	8 words — < 1%
37	scholar.lib.vt.edu Internet	8 words — < 1%
38	www.bre.co.uk Internet	8 words — < 1%
39	etheses.bham.ac.uk Internet	8 words — < 1%
40	iwcs2015.github.io Internet	8 words — < 1%
41	Osman, H. M., and T. E. El-Diraby. "Knowledge-Enabled Decision Support System for Routing Urban Utilities", Journal of Construction Engineering and Management, 2011.	8 words — < 1%

-
- 42 "Text, Speech, and Dialogue", Springer Nature,
2016 7 words — < 1%
Crossref
-
- 43 Katerina Frantzi. "Automatic recognition of multi-
word terms: the C-value/NC-value method",
International Journal on Digital Libraries, 08/01/2000 6 words — < 1%
Crossref

EXCLUDE QUOTES ON
EXCLUDE BIBLIOGRAPHY ON

EXCLUDE MATCHES OFF