

NLP-based approach to classify heterogeneous terms for unambiguous exchange of roadway data

Tuyen Le ¹, H. David Jeong ²

ABSTRACT

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 in digital dictionaries, such as lexicons, and ontologies can enable the meaning of a concept name to be unambiguously understood by computer systems. However, the current manual process of identifying these relations to support developing semantic resources for the highway sector is laborious and time-consuming due to the lack of an effective automated method. This paper presents a novel methodology that leverages recent advances in Natural Language Processing (NLP) techniques to extract English-American roadway terms and their semantic relations from various design guidelines used in different government agencies. Natural Language Processing (NLP) techniques and the C-value method are first implemented to extract commonly used technical terms from a corpus of roadway design guidelines collected from across the State Departments of Transportation. A model for measuring semantic similarity is then trained on the data of context words in the corpus using the Skip-gram neural network model. This semantic model is then utilized by a proposed term classification algorithm that measures the semantic similarity between terms and assigns relation types (synonyms, hyponyms, and attributes) to pairs of related terms. The proposed methodology was evaluated by conducting

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.

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, Semantic Relation, Data Sharing, Semantic Interoperability, NLP, 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 through the entire project life cycle. The inadequate interoperability cost is estimated of 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 through 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, data representation/terminology differs between phases, stakeholders, or geographic regions

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

Terminology transparency through digital dictionaries like glossaries, taxonomies, ontologies and data dictionaries is identified as a driver of semantic interoperability (Ouksel and Sheth 1999). Unfortunately, although, a plethora of semantic resources have been introduced for the highway sector; as shown in the literature review, their coverages of concepts are still inadequate and the inclusion of multiple names to the same concept is still 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. There is a need for computer-aided methods to remove this knowledge acquisition bottleneck (Mounce et al. 2010), such that digital dictionaries can be quickly constructed to meet a specific need and to keep up with the sustainable growth of terms arisen along with new knowledge and technologies.

Recent achievements in accuracy and processing time of advanced Natural Language

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

This paper presents an NLP-based automated approach to gather commonly used American-English roadway terms in different state agencies and classify the semantic relations among these heterogeneous terms by analyzing their occurrence in roadway text documents. In order to achieve that goal, Natural Language Processing (NLP) techniques and the C-value method (Frantzi et al. 2000) are used to detect technical terms from a corpus of roadway design guidelines collected from across the State Departments of Transportation. A model for measuring semantic similarity is then trained on the dataset of words and their contexts in the corpus using the Skip-gram neural network model (Mikolov et al. 2013). This semantic model is then utilized by a proposed term classification algorithm that measures the semantic similarity between terms and assigns a relation type (synonyms, hyponyms, and attributes) to each pair of related terms. A Java package and several datasets result from the study can be found at <https://github.com/tuyenbk/mvdgenerator>.

BACKGROUND

Natural Language Processing

NLP is a research area developing techniques that can be used to analyze and derive value information from natural languages like text and speech. Some of the major applications of NLP include language translation, information extraction, opinion mining (Cambria and White 2014). These applications are embodied by a rich set of NLP techniques ranging

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

Vector Representation of Word Semantics

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

131 outperforms the dictionary-based method in terms of time saving as a semantic model can
132 be automatically obtained from a text corpus and corpus collecting is much easier than
133 manually constructing a digital dictionary (Turney and Pantel 2010).

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

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

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

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 development methods require numerous 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 con-

cepts 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 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.

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-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 level for the whole industry when semantic resources for many aspects of a project are still not available. A few efforts have been made to automate the process of constructing or extending existing semantic resources. The most rigorous methodology in the state-of-the-art is the one developed by Zhang and El-Gohary (2016) that is fully-automated with high accuracy. One limitation of this algorithm is the reliance on an existing semantic resource; it, therefore, would not be applicable to such a domain like the infrastructure that is out of the vocabulary scope. Thus, there is a need for an automated approach that can not only allow for a fast development of highway lexicons but also remove the dependence on other existing semantic models.

PROPOSED METHODOLOGY TO AUTOMATED CLASSIFICATION OF HIGHWAY TERMS

The goal of this research is to propose an NLP-based methodology that can automate the process of extracting roadway technical terms and their semantic relations from American-English roadway documents. As shown in Figure 1, the proposed methodology consists of three major modules that are to: (1) utilize NLP techniques to extract multi-word roadway

technical terms from a collected text corpus, (2) train the data obtained from text corpus 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.

Text corpus collection

In order to develop a domain text corpus for the highway sector, the authors collected a plethora of highway engineering manuals and guidelines from the Federal Department of Transportation (DOT) and from 22 State DOTs. The content of a written guidance document in the engineering field is commonly presented in various formats such as plain text, tables, and equations. Since the structure of words in tables and equations are not yet supported by the state-of-the-art NLP techniques, they were removed from the text corpus. The removal of these features slightly reduces the corpus size, and accordingly affects the training dataset; however, it is necessary since words in tables and equations are not organized in the formal structure of a sentence and therefore the NLP algorithm may extract unreal noun phrases. The final outcome of this phase is a plain text corpus consisting of nearly 16 million words. This dataset is utilized to extract multiple-word technical terms which are then trained and transformed into representation vectors.

Multi-word terms extraction

Linguists argue that a technical term is either a noun (e.g., road) or a noun phrase (NP) (e.g., right of way) that frequently occurs in domain text documents (Justeson and Katz 1995). The meaning of a multi-word term may not be directly interpreted from the meanings of its constituents; therefore, it must be treated as an individual word. To ensure that, multi-word terms in the corpus need to be detected and replaced with connected blocks of their members. As mentioned, a multiple-word term must be a noun phrase; thus, noun phrases will be good multi-word term candidates. To detect this type of terms, the corpus is

first scanned to search for NPs; and the importance of each of these NPs is then evaluated based on the statistics of occurrence frequencies. The process of extracting multi-word terms is discussed in detail as follows.

Noun phrase extraction

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

- i Word tokenizing: In this step, the text corpus is broken down into individual units (also called tokens) using OpenNLP Tokenizer.
- ii Part of Speech (POS) tagging: The purpose of this step is to determine the Part of Speech (POS) tag (e.g., NN-noun, JJ-adjective, VB-verb, etc.) for each unit of the tokenized corpus obtained from the previous step. A full set of POS tags can be found in the Penn Treebank (Marcus et al. 1993).
- iii Noun phrase detection: Table 1 presents the proposed extraction patterns which are modified from the filters suggested by Justeson and Katz (1995) to extract NPs. The tagged corpus is thoroughly scanned to collect sequences matching those patterns. In addition, in order to reduce the discrimination between the syntactic variants of the same term, the collected NPs need to be normalized. The following discuss two types of syntactic variants considered and the proposed normalization methods.
 - Type 1 - Plural forms, for example ‘roadways’ and ‘roadway’. Stemming is a popular process to reduce words to their stems. Despite the fact that, none of the existing algorithms can completely eliminate the errors of over and under stemming, they are good enough to not degrade the overall performance of NLP application (Jivani et al. 2011). This study implements the Pling stemmer (Suchanek et al. 2006), which stems an English noun to its singular form, to normalize plural nouns

in the corpus. One advantage of this algorithm is the utilization of both syntactic rules and the vocabulary in a dictionary; hence the miss- or over-stemming errors that take off a true suffix can be reduced.

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

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

Multi-word term candidate ranking and selection

Multi-word term definition varies between authors, and there is a lack of formal and widely accepted rules to define if an NP is a multi-word term (Frantzi et al. 2000). There are a number of methods proposed for estimating termhood (the degree that a linguistic unit is a domain-technical concept), such as TF-IDF (Sparck Jones 1972; Salton and Buckley 1988), C-Value (Frantzi et al. 2000), Termex (Sclano and Velardi 2007). These methods are based on the occurrence frequencies of NPs in the corpus. Among these methods, Termex outperformed other methods on the Wikipedia corpus, and C-Value was the best on the

GENIA medical corpus (Zhang et al. 2008). This result indicates that the C-value method is more suitable for term extraction from a domain corpus rather than a generic corpus. For this reason, the C-value has been widely used to extract domain terms in the biomedical field, for instance studies performed by Ananiadou et al. (2000), Lossio-Ventura et al. (2013), and Nenadić et al. (2002). Since the corpus used in this study was mainly collected from technical domain documents, C-value would be the most suitable for the termhood determination task. The C-value measure, as formulated in Equation 1, suggests that the longer an NP is, the more likely that is a term; and the more frequently it appears in the domain corpus, the more likely it will be a domain term.

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

Where:

a is a candidate noun phrase

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

f is the frequency of *a* in the corpus

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

P(Ta) is the size of Ta set.

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

To automatically remove candidates that are unlikely to be real terms, a threshold C-value can be used. However, doing this may eliminate the real terms that appear in the bottom due to their low frequencies. Manual evaluation of the entire candidate list would

avoid the removal of real terms with low C-values. To minimize both laborious work and the number of true terms wrongly discarded, the authors suggest the following method to identify the threshold value. The ranked list of candidates is divided into groups of 100 items. A graduate student with a civil engineering background was asked to utilize a bottom-up approach to evaluate group by group and stop at which the percentage of actual terms achieved 80 percent. Users can choose a higher percentage in cases that extracting of rare terms is critical. Table 2 illustrates the evaluation results for several excerpts of the extracted term candidates. The precision values, which represent the percentages of real terms in these groups, are presented in Figure 3. As shown in the figure, precision values are less than 80 percent for groups with c-values less than 70. This value is set as the threshold for the acceptance of term candidates. The final selected list is comprised of nearly 6,000 multi-word roadway technical terms.

Construction of term space model

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

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

‘vertical-alignment’.

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

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

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

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

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

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

The learning model includes three major parameters that are *frequency threshold*, *hidden layer size* and *window size* (see Table 3). To eliminate those data points with low frequencies of occurrence that are unlikely to be technical terms, Word2vec allows for the use of *frequency threshold*. Any word with the rate lower than the limit will be ignored. Radim (2014) suggests a range of (0-100) depending on the data set size. Setting this parameter high will enhance the accuracy, but many true technical terms would be out of vocabulary. A preliminary study based on the preliminary corpus with only several millions of words shows that with the frequency of 20, there are very few non-technical terms involved in the training dataset. Hence, with the larger dataset to be collected, this parameter can be higher and up to around 50. The second important parameter is *layer size* which determines the number of nodes in the hidden layer. This parameter highly affects the training accuracy and processing time. A larger layer size is better in terms of accuracy, but this will be paid off by the running time. A reasonable figuration for this parameter is from tens to hundreds (Radim 2014). The final major parameter, *context window size*, decides how many context words to be considered. Google recommends a size of 10 for the Skip-gram model (Google Inc. 2016). These parameters are subject to be changed so that the best model can be achieved. The effects of these parameters on the model performance are discussed in Section 4.

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

The similarity between terms in the Rd-VSM model can be measured by the angle between two word representation vectors (Equation 2) or the distance between two word points

(Equation 3). Figure 5 illustrates the clustering of terms by their distances. In this figure, an *inlet* can be inferred to be more similar to an *outlet* (blue) than a *sidewalk* (green). Using this technique, the most similar terms for a given term can be obtained. Table 4 shows a partial ranked list of the nearest terms of ‘roadway’ in order of similarity score.

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

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

Where: n is the hidden layer size.

Semantic relation classification

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

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

Attributes and hyponyms

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

Synonyms

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

By the end of the synonym recognition phase, the algorithm will generate a list of classified nearest terms. Table 6 shows one example of the output generated by the algorithm.

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 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

Algorithm 1 Semantic relation classification algorithm

```
1: Inputs: term  $t$ , list of nearest terms  $N$ , full list of 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
```

the system. The F score, which is a combined measure of precision and recall, presents the overall performance of a system.

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

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

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

488 Table 7 shows the performance with various training model settings. The parameters of
489 the baseline model are 50, 100 and 5 respectively for frequency threshold, hidden layer and
490 window size. The authors changed the configuration of these parameters one by one to
491 evaluate their effects to the model performance. While changing a certain parameter, other
492 parameters are kept unchanged compared to their values in the base model. As presented

in the table, the model performance is not significantly sensitive to the changes of training parameters. The increase of window size to 10 or 15 resulted in the best model which has a precision of 81% and an F-measure of 65%. The changes of other parameters did not improve the performance. Especially, the increase of frequency threshold value has negative impact. This result confirms the reasonable selection of the frequency threshold to eliminate unlikely term candidates in the NP extraction phase.

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

DISCUSSIONS

This paper proposes an NLP based methodology to assist professionals in extracting roadway terms and their semantic relations from text documents. A key contribution to the body of knowledge is the novel framework with a new algorithm that allows for automated detection of technical terms and their relations without reliance on existing hand-coded dictionaries as used by previous researchers such as Zhang and El-Gohary (2016). The present framework is not to completely eliminate the human interfere, but is expected to become an enabling tool that can help researchers in the domain quickly develop supporting ontologies and other forms of semantic resources for their specific use cases. With respects to the facilitation of semantic interoperability for the infrastructure sector, the implications of this study would accelerate the process of removing the current bottleneck in extensive machine readable dictionaries which are required for an unambiguous data sharing, integration or exchange.

The semantic similarity model and the relation classification algorithm developed in this study are expected to become fundamental resources for a variety of NLP related studies in

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

The current study has a number of limitations. First, the highway corpus is still relatively small with only 16 million words, compared to the corpus sizes in other domains with billions of words. Since the recall value largely depends on the corpus size, the expansion of the highway corpus size by adding more documents from other state agencies and disciplines (e.g., survey, construction, operation and maintenance) would enable more technical terms to be covered in the vector space model; and consequently, the recall would be improved. Secondly, the number of semantic relation categories is limited to only three types of relations that are attributes, hyponyms and synonyms. There are other important semantic relations

that are not considered such as hypernyms, siblings, functional associations, etc. Including these relations would enhance the precision value when incorrect answers such as saying a functionally associated term is a synonym can be reduced. Third, this study only targets at the synonymy issue, the issue of polysemy is not yet addressed. Further research is needed to detect different senses of a roadway term. One potential solution is to apply cluster analysis on the instances of context to determine the possibility that a term would have multiple meanings.

CONCLUSIONS

Data manipulation from multiple sources is a challenging task in highway asset management due to the inconsistency of data format and terminology. The contribution of this study is an NLP-based approach to automated classification of semantic relations among roadway technical terms based on their word occurrences in the domain text documents. This research employs advanced NLP techniques to extract technical terms from a highway text corpus which is composed of 16 million words built on a collection of design manuals from 22 State DOTs across the U.S. Machine learning is used to train the semantic similarity between technical terms. An algorithm is designed to classify the nearest terms resulted from the semantic similarity model into distinct groups according to their lexical relationships.

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

The proposed methodology for detecting semantic relations, with the enhanced automation and scalability level, is expected to significantly reduce human efforts in developing a

semantic resource for a specific use case within the highway domain and become an enabler for semantic interoperability in this domain. The research also opens a new gate for computational tools regarding natural language processing in the highway sector. The developed system would enable computer systems to understand terms and consequently transform the way human interacts with a computer by allowing users to use natural language.

REFERENCES

- Abuzir, Y. and Abuzir, M. O. (2002). “Constructing the civil engineering thesaurus (cet) using the thesweb.” *Computing in Civil Engineering*.
- Ananiadou, S., Albert, S., and Schuhmann, D. (2000). “Evaluation of automatic term recognition of nuclear receptors from medline.” *Genome Informatics*, 11, 450–451.
- Apache.org (2016). “Machine learning library (mllib), <<https://spark.apache.org/docs/1.1.0/mllib-guide.html>> (March).
- buildingSMART (2016). “buildingsmart data dictionary, <<http://bsdd.buildingsmart.org/>>. Accessed: March 15, 2016.
- Cambria, E. and White, B. (2014). “Jumping nlp curves: a review of natural language processing research [review article].” *Computational Intelligence Magazine, IEEE*, 9(2), 48–57.
- Chen, D. and Manning, C. D. (2014). “A fast and accurate dependency parser using neural networks.” *EMNLP*, 740–750.
- Church, K. W. and Hanks, P. (1990). “Word association norms, mutual information, and lexicography.” *Computational linguistics*, 16(1), 22–29.
- Costa-Jussa, M. R., Farrús, M., Mariño, J. B., and Fonollosa, J. A. (2012). “Study and comparison of rule-based and statistical catalan-spanish machine translation systems.” *Computing and Informatics*, 31(2), 245–270.
- Cunningham, H., Maynard, D., Bontcheva, K., and Tablan, V. (2002). “Gate: an architecture for development of robust hlt applications.” *Proceedings of the 40th annual meeting on*

association for computational linguistics, Association for Computational Linguistics, 168–175.

El-Diraby, T. and Kashif, K. (2005). “Distributed ontology architecture for knowledge management in highway construction.” *Journal of Construction Engineering and Management*, 131(5), 591–603.

El-Diraby, T., Lima, C., and Feis, B. (2005). “Domain taxonomy for construction concepts: Toward a formal ontology for construction knowledge.” *Journal of Computing in Civil Engineering*, 19(4), 394–406.

Erk, K. (2012). “Vector space models of word meaning and phrase meaning: A survey.” *Language and Linguistics Compass*, 6(10), 635–653.

Frantzi, K., Ananiadou, S., and Mima, H. (2000). “Automatic recognition of multi-word terms: the c-value/nc-value method.” *International Journal on Digital Libraries*, 3(2), 115–130.

Gallaher, M. P., O’Connor, A. C., Dettbarn, J. L., and Gilday, L. T. (2004). *Cost analysis of inadequate interoperability in the U.S. capital facilities industry*. U.S. Department of Commerce Technology Administration, National Institute of Standards and Technology, Gaithersburg, MD.

Google Inc. (2016). “word2vec, <<https://code.google.com/archive/p/word2vec/>>. (accessed May 12, 2016).

Harispe, S., Ranwez, S., Janaqi, S., and Montmain, J. (2013). “Semantic measures for the comparison of units of language, concepts or instances from text and knowledge base analysis.” *arXiv preprint arXiv:1310.1285*.

Harris, Z. S. (1954). “Distributional structure.” *Word*.

Hezik, M. (2008). “Ifd library background and history.” *The IFD Library/IDM/IFC/MVD Workshop*.

Hill, F., Reichart, R., and Korhonen, A. (2015). “Simlex-999: Evaluating semantic models with (genuine) similarity estimation.” *Computational Linguistics*, 41(4), 665–695.

- Hsu, J.-y. (2013). “Content-based text mining technique for retrieval of cad documents.” *Automation in Construction*, 31, 65–74.
- IFD Library Group (2008). “Ifd library white paper. Accessed: 2015-07-06.
- Jiang, J. J. and Conrath, D. W. (1997). “Semantic similarity based on corpus statistics and lexical taxonomy.” *arXiv preprint cmp-lg/9709008*.
- Jivani, A. G. et al. (2011). “A comparative study of stemming algorithms.” *Int. J. Comp. Tech. Appl*, 2(6), 1930–1938.
- Justeson, J. S. and Katz, S. M. (1995). “Technical terminology: some linguistic properties and an algorithm for identification in text.” *Natural Language Engineering*, 1(01), 9–27.
- Karimi, H. A., Akinci, B., Boukamp, F., and Peachavanish, R. (2003). “Semantic interoperability in infrastructure systems.” *Information Technology*, 42–42.
- Kolb, P. (2008). “Disco: A multilingual database of distributionally similar words.” *Proceedings of KONVENS-2008, Berlin*.
- Lee, T., Wang, Z., Wang, H., and Hwang, S.-w. (2013). “Attribute extraction and scoring: A probabilistic approach.” *Data Engineering (ICDE), 2013 IEEE 29th International Conference on*, IEEE, 194–205.
- Lefler, N. X. (2014). “Roadway safety data interoperability between local and state agencies.” *Report no.*
- Lesk, M. (1986). “Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone.” *Proceedings of the 5th annual international conference on Systems documentation*, ACM, 24–26.
- Levy, O., Goldberg, Y., and Dagan, I. (2015). “Improving distributional similarity with lessons learned from word embeddings.” *Transactions of the Association for Computational Linguistics*, 3, 211–225.
- Lima, C., El-Diraby, T., and Stephens, J. (2005). “Ontology-based optimization of knowledge management in e-construction.” *Journal of IT in Construction*, 10, 305–327.
- Lossio-Ventura, J. A., Jonquet, C., Roche, M., and Teisseire, M. (2013). “Combining c-

- value and keyword extraction methods for biomedical terms extraction.” *LBM’2013: 5th International Symposium on Languages in Biology and Medicine*.
- Lv, X. and El-Gohary, N. M. (2015). “Semantic annotation for context-aware information retrieval for supporting the environmental review of transportation projects.” *Computing in Civil Engineering 2015*, ASCE, 165–172.
- MacQueen, J. (1967). “Some methods for classification and analysis of multivariate observations.” *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, Vol. 1, Oakland, CA, USA., 281–297.
- Marcus, M. (1995). “New trends in natural language processing: statistical natural language processing.” *Proceedings of the National Academy of Sciences*, 92(22), 10052–10059.
- Marcus, M. P., Marcinkiewicz, M. A., and Santorini, B. (1993). “Building a large annotated corpus of english: The penn treebank.” *Computational linguistics*, 19(2), 313–330.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). “Efficient estimation of word representations in vector space.” *arXiv preprint arXiv:1301.3781*.
- Mounce, S., Brewster, C., Ashley, R., and Hurley, L. (2010). “Knowledge management for more sustainable water systems.” *Journal of information technology in construction*, 15, 140–148.
- Navigli, R. (2009). “Word sense disambiguation: A survey.” *ACM Computing Surveys (CSUR)*, 41(2), 10.
- Navigli, R. and Velardi, P. (2010). “Learning word-class lattices for definition and hypernym extraction.” *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, 1318–1327.
- Nenadić, G., Spasić, I., and Ananiadou, S. (2002). “Automatic acronym acquisition and term variation management within domain-specific texts.” *Third International Conference on Language Resources and Evaluation (LREC2002)*, 2155–2162.
- Noy, N. F. (2004). “Semantic integration: a survey of ontology-based approaches.” *ACM Sigmod Record*, 33(4), 65–70.

- Osman, H. and Ei-Diraby, T. (2006). “Ontological modeling of infrastructure products and related concepts.” *Transportation Research Record: Journal of the Transportation Research Board*, 1984(-1), 159–167.
- Ouksel, A. M. and Sheth, A. (1999). “Semantic interoperability in global information systems.” *ACM Sigmod Record*, 28(1), 5–12.
- Pennington, J., Socher, R., and Manning, C. D. (2014). “Glove: Global vectors for word representation.” *Empirical Methods in Natural Language Processing (EMNLP)*, 1532–1543, <<http://www.aclweb.org/anthology/D14-1162>>.
- Radim, R. (2014). “Word2vec tutorial, <<http://rare-technologies.com/word2vec-tutorial/>>.
- Rezgui, Y. (2007). “Text-based domain ontology building using tf-idf and metric clusters techniques.” *The Knowledge Engineering Review*, 22(04), 379–403.
- Salton, G. and Buckley, C. (1988). “Term-weighting approaches in automatic text retrieval.” *Information processing & management*, 24(5), 513–523.
- Sclano, F. and Velardi, P. (2007). *Termextractor: a web application to learn the shared terminology of emergent web communities*. Springer, 287–290.
- Seedah, D. P., Choubassi, C., and Leite, F. (2015a). “Ontology for querying heterogeneous data sources in freight transportation.” *Journal of Computing in Civil Engineering*, 04015069.
- Seedah, D. P., Sankaran, B., and O’Brien, W. J. (2015b). “Approach to classifying freight data elements across multiple data sources.” *Transportation Research Record: Journal of the Transportation Research Board*, (2529), 56–65.
- Sparck Jones, K. (1972). “A statistical interpretation of term specificity and its application in retrieval.” *Journal of documentation*, 28(1), 11–21.
- Suchanek, F. M., Ifrim, G., and Weikum, G. (2006). “Leila: Learning to extract information by linguistic analysis.” *Proceedings of the 2nd Workshop on Ontology Learning and Population: Bridging the Gap between Text and Knowledge*, 18–25.
- Toutanova, K., Klein, D., Manning, C. D., and Singer, Y. (2003). “Feature-rich part-of-

- speech tagging with a cyclic dependency network.” *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*, Association for Computational Linguistics, 173–180.
- Turney, P. D. and Pantel, P. (2010). “From frequency to meaning: Vector space models of semantics.” *Journal of artificial intelligence research*, 37(1), 141–188.
- Walton, C. M., Seedah, D. P., Choubassi, C., Wu, H., Ehlert, A., Harrison, R., Loftus-Otway, L., Harvey, J., Meyer, J., Calhoun, J., et al. (2015). *Implementing the freight transportation data architecture: Data element dictionary*. Number Project NCFRP-47.
- Webster, J. J. and Kit, C. (1992). “Tokenization as the initial phase in nlp.” *Proceedings of the 14th conference on Computational linguistics-Volume 4*, Association for Computational Linguistics, 1106–1110.
- Wetherill, M., Rezgui, Y., Lima, C., and Zarli, A. (2002). “Knowledge management for the construction industry: the e-cognos project.
- Wikipedia (2016). “Glossary of road transportation terms. Accessed: April 11, 2016.
- Yalcinkaya, M. and Singh, V. (2015). “Patterns and trends in building information modeling (bim) research: A latent semantic analysis.” *Automation in Construction*, 59, 68–80.
- Yarowsky, D. (1995). “Unsupervised word sense disambiguation rivaling supervised methods.” *Proceedings of the 33rd annual meeting on Association for Computational Linguistics*, Association for Computational Linguistics, 189–196.
- Zhang, J. and El-Gohary, N. (2016). “Extending building information models semiautomatically using semantic natural language processing techniques.” *Journal of Computing in Civil Engineering*, C4016004.
- Zhang, Z., Iria, J., Brewster, C., and Ciravegna, F. (2008). “A comparative evaluation of term recognition algorithms.” *LREC*.
- Zhao, H. and Kit, C. (2011). “Integrating unsupervised and supervised word segmentation: The role of goodness measures.” *Information Sciences*, 181(1), 163–183.

734	List of Tables	
735	1	Term candidate filters 31
736	2	Excerpts of the extracted candidate terms 32
737	3	Skip-gram model parameters 33
738	4	Examples of top nearest words 34
739	5	Patterns to extract attributes and hyponyms 35
740	6	An example in RoadLex 36
741	7	Performance of the synonym matching task with various training settings . . 37
742	8	Comparison of synonym matching performance between WordNet and RoadLex 38

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
<i>Note:</i> , * 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 in RoadLex

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-<u>10</u>	81	54	65
	50-100- <u>15</u>	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 RoadLex

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

743	List of Figures	
744	1 Overview of the proposed methodology	40
745	2 Linguistic processing procedure to detect technical terms	41
746	3 Multi-word term extraction evaluation	42
747	4 Skip-gram model	43
748	5 Highway term space model (Rd-VSM)	44

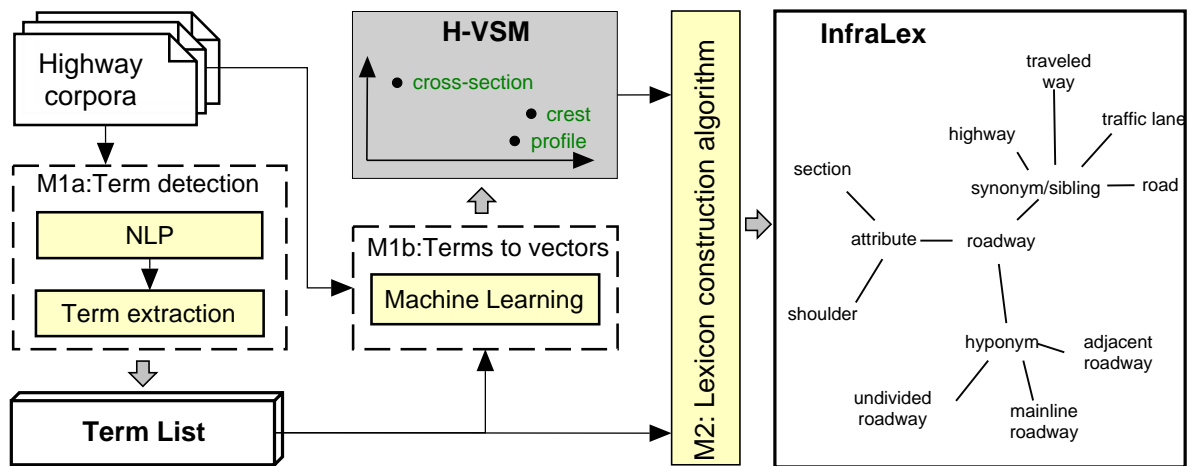


FIG. 1: Overview of the proposed methodology

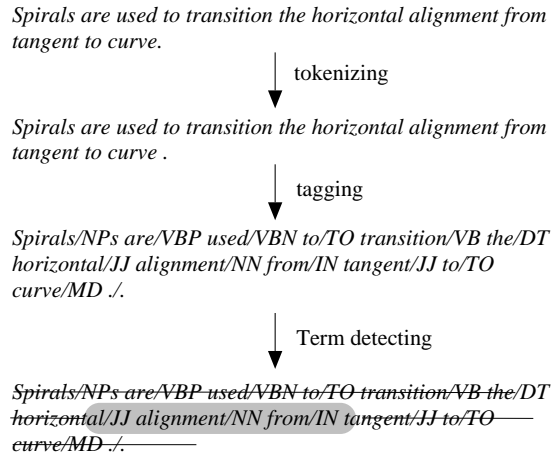


FIG. 2: Linguistic processing procedure to detect technical terms

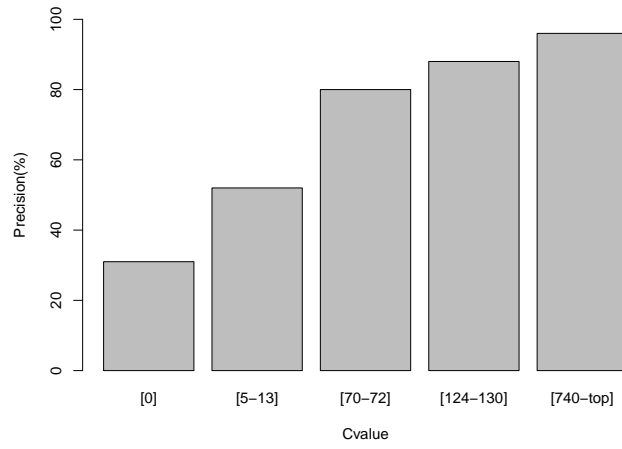


FIG. 3: Multi-word term extraction evaluation

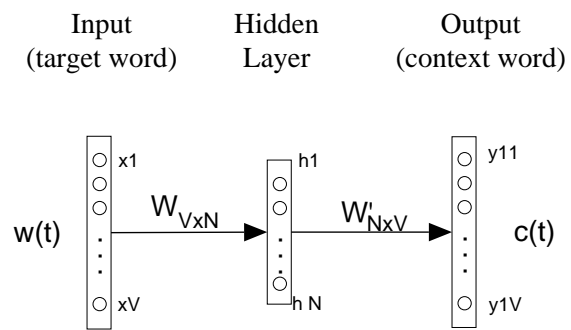


FIG. 4: Skip-gram model

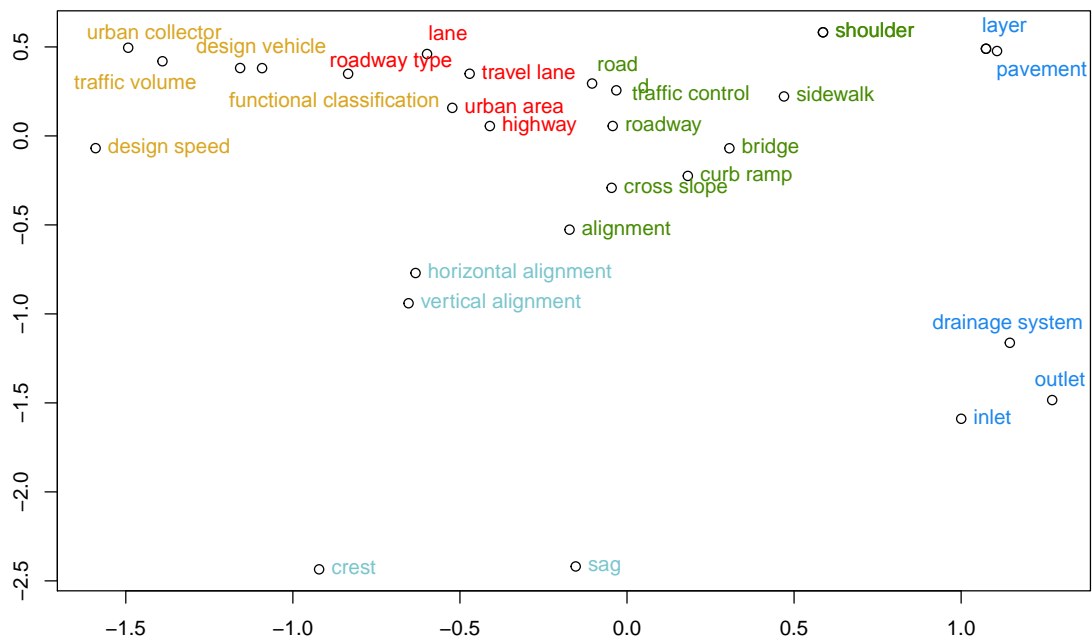


FIG. 5: Highway term space model (Rd-VSM)