1    **Extending Building Information Models Semi-Automatically Using Semantic Natural**

2    **Language Processing Techniques**

3    Jiansong Zhang, S.M.ASCE[1]; and Nora M. El-Gohary, A.M.ASCE[2]

4    **Abstract**

5    Automated compliance checking (ACC) of building designs requires automated extraction of

6    information from building information models (BIMs). However, current Industry Foundation

7    Classes (IFC)-based BIMs provide limited support for ACC, because they lack the necessary

8    information that is needed to perform compliance checking (CC). In this paper, the authors

9    propose a new method for extending the IFC schema to incorporate CC-related information, in

10   an objective and semi-automated manner. The method utilizes semantic natural language

11   processing (NLP) techniques and machine learning techniques to extract concepts from

12   documents that are related to CC (e.g., building codes) and match the extracted concepts to

13   concepts in the IFC class hierarchy. The proposed method includes a set of methods/algorithms

14   that are combined into one computational platform: (1) a method for concept extraction that

15   utilizes pattern-matching-based rules to extract regulatory concepts from CC-related regulatory

16   documents, (2) a method for concept matching and semantic similarity (SS) assessment to select

17   the most related IFC concepts to the extracted regulatory concepts, and (3) a machine learning

18   classification method for predicting the relationship between the extracted regulatory concepts

19   and their most related IFC concepts. The proposed method enables the extension of the IFC

20   schema, in an objective way, using any construction regulatory document. To test and evaluate

21   the proposed method, two chapters were randomly selected from the International Building Code

[1] Graduate Student, Dept. of Civil and Environmental Engineering, Univ. of Illinois at Urbana-Champaign, 205 N. Mathews Ave., Urbana, IL 61801.

[2] Assistant Professor, Dept. of Civil and Environmental Engineering, Univ. of Illinois at Urbana-Champaign, 205 N. Mathews Ave., Urbana, IL 61801 (corresponding author). E-mail:gohary@illinois.edu; Tel: +1-217-333-6620; Fax: +1-217- 265-8039.

22    (IBC) 2006 and 2009.  Chapter 12 of IBC 2006 was used for training/development and Chapter

23    19 of IBC 2009 was used for testing and evaluation. Comparing to manually-developed gold-

24    standards, 91.7% F1-measure, 84.5% adoption rate, and 87.9% precision were achieved for

25    regulatory concept extraction, IFC concept selection, and relationship classification, respectively.

26    **CE Database subject headings:** Project management; Construction management; Building

27    information models; Information management; Computer applications; Artificial intelligence.

28    **Author keywords:** Automated compliance checking; Automated information extraction;

29    Building information modeling (BIM); Natural language processing; Semantic systems;

30    Knowledge modeling; IFC extension; Automated construction management systems.

31    **Introduction**

32    In comparison to the traditional manual process, automated compliance checking (ACC) of

33    construction projects is expected to reduce the time, cost, and errors of compliance checking

34    (Nguyen and Kim 2011; Kasim et al. 2013; Zhang and El-Gohary 2013). ACC has been pursued

35    both in academia and industry. Among the existing construction regulatory ACC efforts, building

36    information models (BIMs) were mostly utilized as the representation of project information

37    (Eastman et al. 2009). Due to the lack of a fully developed all-inclusive BIM data/information

38    schema that can sufficiently represent project information for ACC needs in different areas (e.g.,

39    fire safety, structural safety, and sustainability), existing ACC efforts typically went into one of

40    two directions: either creating their own BIM or extending existing BIMs.

41    One of the important ACC projects, the Construction and Real Estate NETwork (CORENET)

42    project of Singapore (Eastman et al. 2009), developed their own semantic objects in FORNAX

43    library (i.e., a C++ library) to represent building design information. In the U.S., the General

44  Services Administration (GSA) design rule checking efforts defined the BIM modeling

45  requirements in a well-documented building information modeling guide and allowed users to

46  choose their own BIM authoring tool to define building models according to the guide (Eastman

47  et al. 2009). In addition, many of the existing research efforts proposed or implemented the idea

48  of extending BIMs to fulfill their specific information needs. For instance, Nguyen and Kim

49  (2011) and Sinha et al. (2013) extended existing BIMs in Autodesk Revit Architecture by

50  creating new project parameters such as "area of opening in firewall" and "width of opening in

51  firewall"; Kasim et al.  (2013) extended existing BIMs through adding new data items into

52  industry foundation classes (IFC)-represented BIMs directly; Nawari (2011) proposed the

53  development of appropriate Information Delivery Manuals (IDM) and Model View Definitions

54  (MVDs) for the ACC domain to achieve the required level of detail on IFC-represented BIMs;

55  and Tan et al. (2010) extended IFC in eXtensible markup language (ifcXML) to develop an

56  extended building information modeling (EBIM) in XML.

57  These existing efforts to extend BIMs for ACC deepened the understanding of BIM modeling

58  requirements for ACC. However, the model extension methods were mostly ad-hoc and

59  subjective (i.e., relying on subjective developments or extensions by individual software

60  developers and/or researchers); and the resulting models were usually still missing essential

61  compliance checking (CC)-related information that are needed to achieve complete automation

62  in CC (Martins and Monteiro 2013; Niemeijer et al. 2009). In addition, such ad-hoc and

63  subjective developments/extensions lack generality and objectivity, which are essential to full

64  automation of CC at a broader scale. As a result, a more generalized and objective method is

65  needed to extend BIMs for facilitating ACC.

66   To address this gap, in this paper, a new method for extending the IFC schema with regulatory

67   requirement information, in an objective and semi-automated manner, is proposed. The method

68   utilizes semantic natural language processing (NLP) techniques and machine learning (ML)

69   techniques to extract concepts from documents that are related to CC (e.g., building codes) and

70   match the extracted concepts to concepts in the IFC class hierarchy to extend the IFC schema.

71   The method includes developing a set of algorithms/methods and combing them into one

72   computational platform: (1) a pattern matching-based concept extraction method utilizing a set of

73   Part-Of-Speech (POS) patterns to extract regulatory concepts from CC-related documents, (2) a

74   semantic similarity (SS)-based ranking method utilizing a newly-proposed equation to measure

75   SS between concepts to identify the most related IFC concepts to the extracted regulatory

76   concepts, and (3) a ML-based classifier to predict the relationship between the extracted

77   regulatory concepts and their most related IFC concepts. This paper focuses on presenting each

78   method/algorithm and their evaluation results.

79   **Background**

80   ***Building Information Modeling (BIM) and IFC***

81   According to the National Building Information Model Standard Project Committee (National

82   Institute of Building Sciences 2014), a building information model (BIM) is "a digital

83   representation of physical and functional characteristics of a facility. BIM is believed to improve

84   interoperability through structured information and coordinated information flow during a

85   building life cycle and between different disciplines (Hamil 2012). However, although BIM is

86   intended to be fully interoperable, in reality different BIM softwares and platforms are not yet

87   realizing full compatibility and seamless information exchange hitherto, which prevents BIM

88   from realizing its full potential (Young et al. 2009).

89     Standardization is a primary way to improve interoperability. The current main standardization

90     efforts in BIM include Industry Foundation Classes (IFC) and the CIMSteel Integration

91     Standards (CIS/2) (Isikdag et al. 2007). The IFC represent the main data model to describe

92     building and construction industry data. The IFC schema specifications are written using the

93     EXPRESS data definition language (ISO 10303-11 by the ISO TC 184/SC4 committee)

94     (BuildingSmart 2014). The IFC schema is the data exchange standard to facilitate

95     interoperability in the construction industry (BuildingSmart 2014). The CIS/2 is a product model

96     and data exchange file format for structural steel project information (AISC 2014). Both IFC and

97     CIS/2 models are defined using standard STEP description methods, which is the official

98     "Standard for Exchange of Product model data" – ISO 10303. In contrast to CIS/2 which is

99     focusing on modeling information of structural steel framework, IFC schema is designed to

100     cover all subdomains and phases of building and construction industry. Thus, IFC attracted more

101     attention in BIM research. IFC schema is neutral and platform independent. IFC schema is

102     registered as ISO/PAS 16739 and is becoming an official international standard ISO/IS 16739.

103     *Semantic Models and WordNet*

104     In general, "semantics" studies the meanings of the words (Fritz 2006). A semantic model

105     defines data/information entities and relationships between the entities (Hanis and Noller 2011).

106     Specialization and decomposition are two main hierarchical relations in a semantic model.

107     Specialization refers to the relationship between a superclass and its subclass (Dietrich and

108     Urban 2011). Decomposition refers to the relationship between a whole object (class) and the

109     parts that belong to the object (class) (Klas and Schrefl 1995; Böhms et al. 2009).

110     Ontology is the most widely explored and adopted type of semantic model. It is defined as "an

111     explicit specification of a conceptualization", where "a conceptualization is an abstract

112   simplified view of the world that we wish to represent for some purpose" (Gruber 1995; El-

113   Gohary and El-Diraby 2010). An ontology models domain knowledge in the form of concept

114   hierarchies, relationships (between concepts), and axioms to help define the semantic meaning of

115   the conceptualization (El-Gohary and El-Diraby 2010). Ontological models were shown to

116   improve the performance of various information processing tasks such as text classification (e.g.,

117   Zhou and El-Gohary 2014) and information extraction (e.g., Zhang and El-Gohary 2013; Soysal

118   et al. 2010).

119   WordNet was frequently utilized in semantic research efforts. It is a large lexical database of

120   English where the four types of POS words (nouns, verbs, adjectives, and adverbs) are grouped

121   into sets of cognitive synonyms (synsets) (Fellbaum 2005). In WordNet, each of the four POS

122   categories is organized into a subnet and the synsets are linked to each other using one or more

123   of the following conceptual semantic and lexical relations: synonymy, hyponymy (sub-super or

124   is-a relation), meronymy (part-whole relation), and antonymy (Fellbaum 2005). Because of the

125   abundant, explicitly-defined and well-structured conceptual semantic relations between word

126   senses in WordNet, WordNet has been widely used in semantic research, as a "lexical database"

127   (Shehata 2009; Kamps et al. 2004), a "lexical dictionary" (Varelas et al. 2005), a "semantic

128   dictionary" (Simpson and Dao 2010), or a "domain-independent background knowledge model"

129   (Suchanek et al. 2007). The lexical relations in WordNet can assist in semantic text processing.

130   The hyponymy and meronymy relations in WordNet correspond well to the is-a and part-whole

131   relations in semantic models. In addition, a synonymy relation carries an "equivalency" relation

132   between semantic classes.

133   In spite of the generally accepted assumption that semantic relations are domain dependent

134   (Orna-Montesinos 2010), WordNet, as a resource for providing semantic relations across

135   different domains, is still widely used for domain specific text/knowledge processing tasks. This

136   may be caused by a lack of domain-specific semantic relation resource with a comparative

137   coverage as that can be achieved using WordNet. For example, in OmniClass (a popular

138   classification system for the construction industry), common concepts in building codes (e.g.,

139   "cross ventilation") may be difficult to find matches. In contrast, the dictionary-level coverage in

140   WordNet contains the semantic relations for each word, and enabling semantic analysis for any

141   multi-term concept in a compositional manner. Although being few, WordNet has been used for

142   text/knowledge analysis or processing in Architecture, Engineering and Construction domain,

143   such as in (Orna-Montesinos 2010) and (Li 2010).

144   *Semantic Similarity*

145   Semantic Similarity (SS) is the conceptual/meaning distance between two entities such as

146   concepts, words, or documents (Slimani 2013). SS plays an important role in information and

147   knowledge processing tasks such as information retrieval (Rodrı´guez and Egenhofer 2003), text

148   clustering (Song et al. 2014), and ontology alignment (Jiang et al. 2014).

149   SS could be quantitatively estimated using different measures. Some popular measures are: (1)

150   Shortest Path Similarity, which utilizes the shortest path connecting two concepts in a taxonomy

151   (i.e., concept hierarchy) (Resnik 1995); (2) Leacock-Chodorow Similarity, which utilizes the

152   shortest path connecting two concepts in a taxonomy while penalizing long shortest path

153   according to the depth of the taxonomy (Resnik 1995); (3) Resnik Similarity, which utilizes the

154   information content measure from information theory to measure the information shared by two

155   concepts using the information content of the two concepts' least common subsumer (Resnik

156   1995); (4) Jiang-Conrath Similarity, which utilizes the information content of the two concepts in

157   addition to that of their least common subsumer in the taxonomy; and (5) Lin Similarity, which

158    utilizes the ratio between the information content of the least common subsumer (in the

159    taxonomy) of the two concepts and the sum of the information contents of the two concepts.

160    Shortest Path Similarity is simple and intuitive, it approximates the conceptual distance between

161    concepts by the number of edges in-between. The main limitation of Shortest Path Similarity is

162    its inability to take specificity of concepts into the measurement, which lead to the same

163    similarity results between concept pairs at shallow level of a taxonomy and concept pairs at deep

164    level of the taxonomy as long as the counts of number of edges for both concept pairs are the

165    same. This limitation is compensated for in Leacock-Chodorow similarity by taking the

166    maximum depth of the taxonomy into the measurement. Thus in Leacock-Chodorow similarity

167    concept pairs deeper in the taxonomy (i.e. more specific) would have larger similarity score than

168    concept pairs shallow in the taxonomy, with the number of edges being equal. Resnik Similarity,

169    Jiang-Conrath Similarity, and Lin Similarity are information content-based similarity measures.

170    They do not have the specificity problem of shortest path similarity because the information

171    content reflects specificity by definition. Resnik Similarity sometimes is considered coarse

172    because different concept pairs could have the same least common consumer. Jiang-Conrath

173    Similarity and Lin Similarity improve upon Resnik Similarity by taking the information content

174    of the concepts from the pair into measurement in addition to their least common consumer. The

175    main limitation for information content-based measures, however, is the need of a corpus in

176    addition to the taxonomy, which lead to different results based on different corpuses. When using

177    the measures herein to evaluate SS, WordNet was widely used as the taxonomy.

178    *Machine Learning Algorithms*

179    In any machine learning (ML) application, different ML algorithms are usually tried out and

180    tested. Some of the most commonly-used ML algorithms are summarized in Table 1.

181          Insert Table 1

182    Naïve Bayes is a simple statistical ML algorithm. It applies Bayes' rule to compute conditional

183    probabilities of predictions given evidence. It is the simplest type of algorithm among the

184    commonly-used ML algorithms. However, Naïve Bayes could outperform more complex

185    learning algorithms in some cases (Domingos 2012). Perceptron is a linear learning algorithm

186    where predictions are made based on a linear combination of feature vectors (Rosenblatt 1958).

187    Perceptron is applicable to problems that are linearly separable. The application process of

188    perceptron is iterative: a prediction vector is iteratively constructed based on each instance in the

189    training dataset (Freund and Schapire 1999). Decision Tree ML algorithms use a tree to map

190    instances into predictions. In a Decision Tree model, each non-leaf node represents one feature,

191    each branch of the tree represents a different value for a feature, and each leave node represents a

192    class of prediction. Decision Tree is a flexible algorithm that could grow with increased amount

193    of training data (Domingos 2012). K-Nearest Neigbor (k-NN) is a similarity-based ML algorithm.

194    K-NN predicts the class of an instance using the instance's k nearest instances by assigning it the

195    majority class of those k instances' classes (Cover 1967; Domingos 2012). Support Vector

196    Machines (SVM) is a kernel-based ML algorithm that has significant computational advantages

197    over standard statistical algorithms. Kernel methods is a technique for constructing nonlinear

198    features so that nonlinear functional relationships could be represented using a linear model. A

199    linear model is much simpler comparing to a nonlinear model, both theoretically and practically,

200    giving SVM its computational advantages (Cristianini and Shawe-Taylor 2000). SVM was found

201    to outperform other ML algorithms in many applications such as text classification (e.g., Salama

202    and El-Gohary 2013).

203    ML is one type of machine-based reasoning (i.e., inductive reasoning), where the various types

204    of ML algorithms induct knowledge from input data (Domingos 2012). In any machine-based

205    reasoning, successful reasoning depends on appropriate representations (Bundy 2013). What

206    features should be used to represent the data in a ML problem is, thus, an important decision.

207    **State of the Art and Knowledge Gaps**

208    Several research efforts extended the IFC schema for various purposes, such as building life

209    cycle management (Vanlande 2008), cost estimation (Ma et al. 2011a), enterprise historical

210    information retrieval (Ma et al. 2011b), parametric bridge models information exchange (Ji et al.

211    2013), and virtual construction systems data sharing (Zhang et al. 2014). However, most of the

212    extension efforts extended the schema in an arbitrary and ad-hoc manner, which lacks objectivity

213    and generality. Despite the potential of using ontology alignment and ontology mapping

214    techniques (using SSs) in developing a generalized IFC extension method, to the best of the

215    authors' knowledge there was little empirical exploration of this approach. The work of Delgado

216    et al. (2013) and Pan et al. (2008) are the closest to this approach.

217    Delgado et al. (2013) evaluated 15 ontology matching techniques in matching geospatial

218    ontologies with BIM-related ontologies (including an ontology for IFC) to discover

219    correspondences of concepts between each pair of ontologies (e.g., between CityGML ontology

220    and IFC ontology). The 15 techniques were classified into three categories: string-based

221    techniques, WordNet-based techniques, and matching systems techniques. The alignment

222    between CityGML ontology and IFC ontology is conceptually and technically similar to

223    extending IFC. The main evaluation metrics were precision, recall, and F1-measure. Precision

224    was defined as the correctly found correspondences divided by the total number of

225    correspondences found. Recall was defined as the correctly found correspondences divided by

226    the total number of correspondences that should be found. F1-measure was defined as the

227    harmonic mean of precision and recall. In their experimental results: (1) String-based techniques

228    showed the best performance (100% precision, 57.1% recall, and 23.2% F1-measure) among the

229    three tested techniques; (2) Within the WordNet-based techniques, the synonym distance

230    technique showed 41.6% precision, 14.2% recall, and 21.2% F-measure; and (3) Within the

231    matching systems techniques, the Association Rule Ontology Matching Approach (AROMA)

232    showed 40% precision, 5.7% recall, and 10% F-measure. These results show that further

233    research is needed to investigate whether the use of other semantic relations in WordNet (such as

234    hyponymy), in addition to synonymy, would result in higher levels of performance.

235    Pan et al. (2008) conducted semi-automated mapping of architectural, engineering, and

236    construction (AEC) ontologies, including an IFC ontology, using relatedness analysis techniques.

237    In their ontology mapping, three types of features were used to provide expert guidance: (1)

238    corpus-based features: co-occurrence frequencies between two concepts, (2) attribute-based

239    features: attribute value structures of ontologies, and (3) name-based features: stemmed terms of

240    the concept names to use for direct term-based matching. Further research is needed to explore

241    how different types of semantic relations among concepts could be leveraged in IFC concept

242    mapping.

243    **Proposed Method for Semi-Automated IFC Extension with Regulatory Concepts**

244    The proposed method for IFC extension with regulatory concepts is semi-automated. It utilizes

245    automation techniques in all tasks for reducing the required manual effort. The proposed method

246    includes three phases (Fig. 1): regulatory concept extraction, IFC concept selection, and

247    relationship classification.

248                             Insert Figure 1

249    *Regulatory Concept Extraction*

250    Proposed Concept Extraction Approach

251    To conduct ACC in a fully automated way, all concepts related to regulatory requirements in a

252    relevant regulatory document must be incorporated into a BIM schema (e.g., IFC schema). The

253    regulatory concept extraction phase aims to automatically extract all concepts from a selected,

254    relevant regulatory document. The proposed extraction method utilizes pattern-matching-based

255    extraction rules. After all concepts are automatically extracted from a textual regulatory

256    document, a user manually removes concepts that are not related (or irrelevant) to regulatory

257    requirements (about one and half pages of words for one chapter).

258    Concept Extraction Rules

259    The left-hand side of a rule defines the pattern to be matched and the right-hand side defines the

260    concept that should be extracted. The patterns are composed of POS features (i.e., POS tags). Fig.

261    2 shows an example concept extraction (CE) rule and its corresponding meaning. The rule could

262    extract four term concepts like "thermally isolated sunroom addition." Ten selected POS tags

263    from Penn Treebank tag set are also listed in Fig. 2. Only flattened patterns are utilized in the CE

264    rules to avoid recursive parsing. This decision is in contrast to grammar-based recursive parsing

265    where the same rule could be applied multiple times for exploring a hierarchical structure as

266    parsing results, because the purpose of the rules herein is only for identifying base noun phrases

267    instead of exploring the internal structures of the phrases. This is analogous to chunking instead

268    of finding constituents. Flattened patterns are patterns that include only terminal symbols (i.e.,

269    symbols that cannot be further broken down), which are analogous to leaf nodes in a tree-like

270    structure. For example, the following P1, P2, P3, and P4 patterns are flattened because they only

271    contain POS tags (i.e., "NN", which is the POS tag for singular noun). Non-flattened patterns, on

272    the other hand, are patterns that include non-terminal symbols (i.e., symbols that could be further

273    broken down). For example, the "NN NP" pattern in rule R2 is non-flattened because it contains

274    non-terminal symbols (i.e., "NP", which is a phrase level tag for noun phrase that could be

275    further broken down). Recursive parsing is avoided, in the proposed method, because: (1)

276    recursive parsing increases time complexities of parsing algorithms. For example in Fig. 3, to

277    match the pattern P4, the number of trials for applying rules R1 and R2 using recursive parsing

278    are minimum 4 and maximum 8, higher than the number of trials for applying rules R1, R3, R4,

279    and R5 using non-recursive parsing which are minimum 1 and maximum 4; and (2) recursive

280    parsing is less flexible. For example, if only P1 and P3 should be matched while P2 and P4

281    should not, this is easy to achieve through applying R1 and R4 using non-recursive parsing,

282    whereas it is difficult to achieve using recursive parsing.

283                                    Insert Figure 2

284                                    Insert Figure 3

285    Development of POS Pattern Set

286    The development of the set of POS patterns to use in the CE rules is conducted following the

287    algorithm shown in Fig. 4. The algorithm is executed after a gold standard of regulatory concepts

288    is created and after the POS tags for all sentences in the development text are generated. The

289    development text is a sample of regulatory text (Chapter 12 of IBC 2006, in this paper) that is

290    used to identify common POS patterns in the text for developing the POS pattern set. The

291    algorithm incrementally processes concepts in the gold standard using two levels of loops: the

292    outer loop accesses each sentence in the gold standard and the inner loop accesses each concept

293    in the sentence being accessed. In the processing of each concept, the POS pattern for the

294    concept is first tentatively collected into the POS pattern set. Then the POS pattern set is used to

295    extract concepts from all sentences. The recall and F1-measure are then calculated for the result.

296    If the recall and F1-measure increase comparing to the previous recall and F1-measure (without

297    the tentatively added POS pattern), then the addition of the POS pattern into the POS pattern set

298    is committed. This process iterates through all concepts in all sentences. The algorithm

299    iteratively improved recall and F1-measure of extraction by incorporating more POS patterns.

300                                        Insert Figure 4

301    Exclusion Word Removal

302    Exclusion words are defined, here, as words (unigram, bigram, or multigram) that match certain

303    POS patterns in the CE rules but should not be extracted as concepts. The POS tags of these

304    exclusion words usually introduce ambiguity because they carry more than one lexical or

305    functional category/meaning, which may introduce false positives (i.e., incorrectly extracted

306    concepts) in concept extraction. For example, "VBG" (POS tag for both "verb gerund" and

307    "present participle") is useful to extract "verb gerund" concepts like "opening", but it introduces

308    false positives when incorrectly extracting "present participle" words like "having" as concepts.

309    To avoid introducing false positives during concept extraction, an exclusion word list is used.

310    *IFC Concept Selection*

311    Proposed Concept Selection Approach

312    The IFC concept selection phase aims to (1) automatically find the most related concept(s) in the

313    IFC schema (called F-concept(s) hereafter) to each extracted regulatory concept (called R-

314    concept hereafter) and (2) accordingly, allow the user to select the F-concept(s) for each R-

315    concept. In this paper, the extension of the IFC schema is an incremental process; each R-

316    concept is added to the IFC schema one by one, incrementally. As a result, an R-concept that

317    gets selected (and thus added) to the IFC schema becomes part of the schema (i.e., becomes an

318    F-concept for the following automated selection step). The automated IFC concept selection

319    method includes four steps/techniques (as shown in Fig. 5): (1) Step 1: stemming, which reduces

320    words to their stems; (2) Step 2: term-based matching, which aims to find all F-concepts that

321    share term(s) with an R-concept; (3) Step 3: semantic-based matching, which aims to find all

322    semantically related F-concepts to an R-concept. Semantic-based matching is used, to add a

323    deeper level of searching, if the term-based matching fails to find candidate concepts; and (4)

324    Step 4: SS scoring and ranking, which measures the SS between each candidate F-concept (from

325    Step 2 and Step 3) and the R-concept, and accordingly ranks all candidate F-concepts related to

326    that one single R-concept for final F-concept user selection. The same process is repeated for all

327    R-concepts and their related candidate F-concepts.

328                                    Insert Figure 5

329    <u>Stemming</u>

330    Stemming is utilized in both term-based and semantic-based matching. Concepts are stemmed

331    before matching to avoid incorrect mismatching due to variant word forms (rather than variant

332    meaning). For example, with stemming applied, "foot" could be matched to "feet" (the stem of

333    "feet" is "foot").

334    <u>Term-Based Matching</u>

335    For term-based matching, three types of matching are used, based on the following two heuristic

336    rules, H1 and H2: (1) "First Term Term-Based Matching": the first term in the R-concept is

337    terminologically matched against all F-concepts to find related F-concepts, (2) "Last Term Term-

338    Based Matching": the last term in the R-concept is terminologically matched against all F-

339    concepts to find related F-concepts, and (3) "First and Last Term Term-Based Matching": the

340    first and last terms in the R-concept are terminologically matched against all F-concepts to find

341    related F-concepts. Which type of matching to use depends on two main factors: (1)  the number

342    of terms in the concept name of the R-concept, whether the concept name is unigram (i.e.,

343    concept name with only one term), bigram (i.e., concept name with two terms), or multigram (i.e.,

344    concept name with three or more terms); and (2) the types of POS patterns in the concept name

345    of the R-concept, whether the pattern is "N" (i.e., a POS pattern with only one POS tag and the

346    POS tag is a noun), "NN" (i.e., a POS pattern starting with a noun and ending with a noun), or

347    "JN" [i.e., a POS pattern starting with a prenominal modifier (e.g., adjective) and ending with a

348    noun]. The matching strategy is illustrated in Fig. 6. For example, the R-concept "interior space"

349    is a bigram, and the POS pattern ("JJ NN") matches "JN", thus "Last Term Match" is used to

350    find matching concepts that contain the term "space" (i.e., the last term in the R-concept).

351    • H1: The term that has a nominal POS tag (i.e., noun) is the primary meaning-carrying

352        term in a multi-term concept name.

353    • H2: The terms that have non-nominal POS tags (e.g., "JJ") are the secondary meaning-

354        carrying terms in a multi-term concept name, which add to or constrain the meaning as

355        modifiers.

356                          Insert Figure 6

357    Semantic-Based Matching

358    In semantic-based matching, the semantic relations of WordNet (Fellbaum 2005) are utilized to

359    find concept matches beyond term-based matching. Three types of these relations are used:

360    hypernymy, hyponymy, and synonymy. These three types were selected because they are most

361    relevant to the superclass-subclass structure of the IFC class hierarchy. Hypernymy is a semantic

362    relation where one concept is the hypernym (i.e., superclass) of the other. For example, "room"

16

363    is a hypernym of "kitchen". Hyponymy is the opposite of hypernymy where one concept is the

364    hyponym (i.e., subclass) of the other. For example, "kitchen" is a hyponym of "room".

365    Synonymy is the semantic relation between different concepts who share the same meaning. For

366    example, "gypsum board", "drywall", and "plasterboard" all share the same meaning of "a board

367    made of gypsum plaster core bonded to layers of paper or fiberboard." Three types of matching

368    are used, which semantically match the first term, the last term, or the first term and last term in

369    the R-concept, respectively, against all F-concepts to find related F-concepts: "first term

370    semantic-based matching", "last term semantic-based matching", and "first and last term

371    semantic-based matching". To conduct the semantic matching, the hypernyms, hyponyms, and

372    synonyms of the first/last term are determined, based on WordNet, and then term-matched

373    against all F-concepts to find related F-concepts. Similar to term-based matching, which type of

374    matching to use depends on (1) the number of terms in the concept name of the R-concept and (2)

375    the POS pattern types in the concept name of the R-concept. The matching strategy is illustrated

376    in Fig. 6. "Search" represents "semantic-based matching" in Figure 6.

377    Semantic Similarity Scoring and Ranking

378    The proposed SS scoring method follows heuristic rules H3, H4, and H5.

379    • H3: In a multi-term concept name, the contribution of each term's carried meaning to the

380       meaning of the whole concept decreases from right to left; the first term contributes the

381       least to the meaning of the whole concept.

382    • H4: The difference in length between two concept names (where the length is measured

383       in number of terms) is indicative of the closeness of the two concepts in a concept

384       hierarchy; the smaller the difference, the closer the two concepts are, and vice versa.

385       Sibling concepts are, thus, likely to have a small difference between their concept name

386       lengths.

387       • H5: The length of a concept name is related to its level in a concept hierarchy. The

388       shorter the length of a concept name is, the more general the concept is; and thus the

389       higher its level in a concept hierarchy. The longer the length of a concept name is, the

390       more specific the concept is; and thus the lower its level in a concept hierarchy. A

391       superconcept is, thus, likely to have a shorter concept name length than its subconcept.

392    Based on these heuristic rules, Eq. (1) and Eq. (2) are proposed as two alternative functions for

393    SS scoring, where $SS_{RF1}$ and $SS_{RF2}$ are the concept-level SS scores between an R-concept and an

394    F-concept, $SS_{RmFk}$ is the term-level SS score between the $m^{th}$ term in the R-concept and the $k^{th}$

395    term in the F-concept, $m$ is the ordinal number for the term $Rm$ in R-concept, $k$ is the ordinal

396    number for the term $Fk$ in F-concept, $L_F$ is the length of F-concept measured in number of terms,

397    and $L_R$ is the length of R-concept measured in number of terms.

398
$$SS_{RF1} = \frac{1}{L_F} \sum_{k=1}^{k=L_F} \frac{2k}{L_F(L_F+1)} SS_{RmFk} \tag{1}$$

399
$$SS_{RF2} = \frac{1}{|L_R - L_F|+1} \sum_{k=1}^{k=L_F} \frac{2k}{L_F(L_F+1)} SS_{RmFk} \tag{2}$$

400    Any existing term pair SS measure, such as Shortest Path Similarity measure or Leacock-

401    Chodorow Similarity measure, can be used (after testing) to compute $SS_{RmFk}$. In Eq. (1) and Eq.

402    (2), each term-level SS score (i.e., $SS_{RmFk}$) is discounted using the factor $\frac{2k}{L_F(L_F+1)}$. This term-

403    level discount factor is based on heuristic rule H3. The concept-level semantic similarity score

404    between the R-concept and the F-concept (i.e., $SS_{RF}$) is determined by further discounting the

405    summation of all discounted term-level SS scores (of all term pairs formed between the matching

406    term of R-concept and each term of the F-concept). In Eq. (1), the concept-level discount factor

407    is $\frac{1}{L_F}$, which linearly discounts the summation using the length of the F-concept. This discount

408    favors concepts at higher levels in a concept hierarchy and follows heuristic rule H5 to identify

409    higher-level concepts based on the lengths of concept names. In Eq. (2), the concept-level

410    discount factor is $\frac{1}{|L_R - L_F| + 1}$, based on the absolute length difference between the concept

411    names of R-concept and F-concept. This discount favors concepts at similar levels in a concept

412    hierarchy and follows heuristic rule H4.

413    Accordingly, the proposed SS scoring method is summarized in Fig. 5. Combinations of different

414    concept-level SS scoring functions (i.e., Eq. 1 and Eq. 2) and term-level SS scoring functions

415    (i.e., existing similarity measures such as Shortest Path Similarity) should be experimentally

416    tested to select the best-performing combination. Separate testing is conducted for term-based

417    matched F-concepts (i.e., F-concepts found using term-based matching, from Step 2) and

418    semantic-based matched F-concepts (i.e., F-concepts found using semantic-based matching, from

419    Step 3). The authors' experimental testing and results are presented and discussed in the

420    Experimental Testing and Results section.

421    For SS ranking, all candidate F-concepts related to one single R-concept are ranked according to

422    their SS scores, in order of decreasing score. A threshold value or a maximum permitted value is

423    further used to filter the most related F-concept(s) among the candidate concepts. The threshold

424    is the minimum SS score below which a candidate F-concept is considered semantically not

425    related (and thus ineligible for selection for this R-concept). The maximum permitted value is a

426  natural number (default is 1) that defines at most how many number of F-concepts could be

427  selected for a single R-concept. Both, threshold value and maximum permitted value, are set by

428  the user. For example, using term-based matching, many F-concepts were found to match

429  "exterior wall" through the matching term "wall", such as "wall" and "curtain wall". Then, the

430  SS scores were computed between "exterior wall" and each of the matched F-concepts, such as

431  "wall and "wall", and "wall" and "curtain wall". The candidate F-concepts were ranked

432  according to the SS scores and the highest scored candidates were automatically selected,

433  according to the default maximum permitted value. If the maximum permitted value is set to 1

434  and Eq. (1) is used, "wall" is selected because of its highest SS score. Following a similar

435  process, but using semantic-based matching, "railing" was selected as the match to "grab bars."

436  ***Relationship Classification***

437  Proposed Classification Approach

438  The relationship classification phase aims to classify the relationship between each pair of R-

439  concept and F-concept. ML techniques are used to automatically predict the relationship between

440  a concept pair based on the concept features of the pair.

441  Types of Relationships

442  Four types of relationships are considered (Table 2): (1) equivalent concept, indicating that the

443  R-concept and the F-concept are equivalent; (2) superconcept, indicating that the R-concept is a

444  superconcept of the F-concept; (3) subconcept, indicating that the R-concept is a subconcept of

445  the F-concept; and (4) associated concept, indicating that the R-concept and the F-concept are

446  associated (bidirectional relationship).

447                                       Insert Table 2
448  Types of Features

449   The authors identified the following initial set of eight features, which includes a mix of

450   syntactic (i.e., related to syntax and grammar) and semantic (i.e., related to context and meaning)

451   features (see Table 3): (1) RTermNum: the number of terms in the concept name of the R-

452   concept, whether the concept name is unigram, bigram, or multigram; (2) RTermPOS: the type

453   of POS pattern in the concept name of the R-concept, whether the pattern is "N", "NN", or "JN";

454   (3) RMatchType: the match type of R-concept, in terms of which term in the R-concept name

455   matches a term in the F-concept name, whether it is the "first" or "last" term in the R-concept

456   name; (4) RelMatchType: the match type between R-concept and F-concept, whether it is "term-

457   based" match, "synonym"-based match (i.e., the matched term in the F-concept name is a

458   synonym of the matching term in the R-concept name), "hyponym"-based match, or

459   "hypernym"-based match; (5) FMatchType: the match type of F-concept, in terms of which term

460   in the F-concept name matches the matching term in the R-concept name, whether it is "first",

461   "middle", or "last"; (6) FTermNum: the number of terms in the concept name of the F-concept,

462   whether the concept name is unigram, bigram, or multigram; (7) FTermPOS: the type of POS

463   pattern in the concept name of the F-concept, whether the pattern is "N", "NN", or "JN"; and (8)

464   DOM: the degree of match, which is represented as a Boolean value describing if the R-concept

465   and the F-concept match term by term, with stemming applied, where one represents match and

466   zero represents no match. These features were identified based on the following heuristic rules:

467   • H5 (see above).

468   • H6: The type of POS pattern in the name of a concept affects its meaning; and since the

469     concept names are all noun phrases, the most distinguishing POS pattern is whether the

470     concept has a modifier(s), and if yes, whether the modifier(s) is/are nominal (i.e., noun or

471     noun sequences).

472    • H7: The match type, in terms of which term in each concept name is matched, affects the

473       relationship between the matched concepts.

474    • H8: The match type, in terms of the type of relationship between the matched terms in

475       both concepts, affects the relationship between the matched concepts.

476    • H9: If, in the same domain, two concept names match term by term (with stemming

477       applied), then the two concepts are likely to be equivalent.

478    Table 4 shows some example concept pairs and their features. The final set of features is

479    determined after conducting feature selection (as further discussed in the Experimental Testing

480    and Results section).

481                                 Insert Table 3

482                                 Insert Table 4

483    **Experimental Testing and Results**

484    The proposed semi-automated IFC extension method was tested on extending the IFC class

485    hierarchy (based on schema version IFC2X3_TC1) using regulatory concepts from IBC. Two

486    chapters, Chapter 12 of IBC 2006 and Chapter 19 of IBC 2009, were randomly selected. Chapter

487    12 was used for: (1) developing the set of POS patterns for use in regulatory concept extraction

488    (Phase 1), (2) selecting the best combination of SS scoring function and SS measure for IFC

489    concept selection (Phase 2), and (3) training the ML classifier for relationship classification

490    (Phase 3). Chapter 19 was used for testing and evaluating each of the following sub-

491    methods/algorithms: regulatory concept extraction, IFC concept selection, and relationship

492    classification. Each sub-method/algorithm was tested separately.

493    *Regulatory Concept Extraction*

494    Gold Standard

495   The gold standards of R-concepts for Chapter 12 of IBC 2006 and Chapter 19 of IBC 2009 were

496   manually developed by the authors. An R-concept is a concept in regulatory documents that

497   defines a "thing" (e.g., subject, object, abstract concept). The criteria for identifying R-concepts

498   in the gold standard is that the concept should be as specific as possible (i.e., including all

499   information related to the concept) without determiners and post modifiers. The longest span for

500   each noun phrase according to this criteria, thus, was manually recognized and extracted as an R-

501   concept. The longest span could be multiple terms (e.g., "continuously operated mechanical

502   operation") or one term (e.g., "ventilation") depending on its appearance in text. For example,

503   concepts in the list L1 were recognized and extracted from Sentence S1. The gold standards of

504   Chapter 12 and Chapter 19 include 368 and 821 concepts, respectively. The concepts in the gold

505   standard will be compared with concepts extracted by the algorithm for evaluating the algorithm

506   in terms of precision and recall of extracted concepts.

507   - S1: "Wall segments with a horizontal length-to-thickness ratio less than 2.5 shall be

508     designed as columns."

509   - L1: ['wall_segments', 'horizontal_length-to-thickness_ratio', 'columns']

510   Algorithm Implementation

511   The proposed regulatory concept extraction method was implemented in Python programming

512   language (v.2.7.3). The Stanford Parser (version 3.4) (Toutanova et al. 2003) was selected and

513   used to generate the POS tags for each word. The Stanford Parser used Penn Treebank tag set

514   which includes 36 tags. Ten, out of the 36 tags, were used (shown in Fig. 2).

515   Evaluation

516   Regulatory concept extraction was evaluated in terms of precision, recall, and F1-measure. The

517   definitions of these measures are similar to those in ontology matching except for the "found

518    correspondences" are replaced by "extracted concepts." A higher recall is more important than

519    precision because the overall method of IFC extension is semi-automated; precision errors could

520    be detected and eliminated by the user during user concept selection.

521    <u>Development Results and Analysis</u>

522    The development of the set of POS patterns to use in the CE rules was conducted following the

523    algorithm shown in Fig. 4. Fig. 7 shows the final set of POS patterns, which consists of 39

524    patterns. These 39 POS patterns were used as conditions for 39 CE rules, one POS pattern for

525    one CE rule. For example, the pattern "JJ" "JJ" "JJ" "NN" was used for a CE rule which extracts

526    three consecutive adjectives followed by a singular/mass noun as a concept, such as in the

527    concept "minimum net glazed area."

528    Table 5 shows the performance of extracting R-concepts from the development text (Chapter 12

529    of IBC 2006). Through error analysis two sources of errors were found: (1) POS tagging error,

530    which accounted for 38.1% of the errors. For example, "herein" was incorrectly tagged as "NN"

531    instead of the correct tag "RB", and was, thus, incorrectly extracted; and (2) ambiguity of the

532    POS tag "VBG" between gerund and present participle, which accounted for 61.9% of the errors.

533    For example, "being" was a present participle thus not representing a concept, but it was

534    extracted because the POS tag "VBG" was included in the POS patterns for representing gerund.

535    While addressing error type (1) depends on improvement of POS taggers, error source (2) was

536    addressed by adding the false positive present participle terms (e.g., "having," "being,"

537    "involving") to the exclusion word list. Membership in the exclusion word list prevents a

538    word/phrase from being extracted in spite of matching a POS pattern in the set. The performance

539    of regulatory concept extraction using the exclusion word list is shown in Table 5. Precision

540    increased, from 93.4% to 97.1%, without decreasing recall.

541                                           Insert Figure 7

542                                           Insert Table 5

543    <u>Testing Results and Analysis</u>

544    The regulatory concept extraction algorithm was tested on Chapter 19 of IBC 2009. The

545    precision, recall, and F1-measure are 89.4%, 94.2%, and 91.7%, and 88.7%, 94.2%, and 91.4%,

546    with and without the use of exclusion word list, respectively. Table 6 shows the performance

547    results. Through error analysis, when using the exclusion word list, four sources of errors were

548    found: (1) POS tagging errors, which accounted for 20.0% of the errors. For example,

549    "corresponding" was incorrectly tagged as "NN" (as opposed to "VBG"); and, thus,

550    "force_level_corresponding" was incorrectly extracted as a concept; (2) ambiguity of POS tag

551    "VBG" between gerund and present participle, which accounted for 7.9% of the errors. For

552    example, "excluding" was incorrectly extracted as a concept because the POS tag for present

553    participle was "VBG" (although it does not represent a meaningful nominal concept); (3)

554    word continuation using hyphen, which accounted for 27.9% of the errors. For example,

555    "pro_vide" was incorrectly extracted as a concept because the word continuation in "pro-vide"

556    led to "pro" and "vide" be tagged as two words with the tags "JJ" and "NN"; and (4) missing

557    POS patterns, which accounted for 44.3% of the errors. For example,

558    "concrete_breakout_strength" and "breakout_strength_requirements" were incorrectly extracted

559    as two concepts (instead of one concept, "concrete_breakout_strength_requirements") because

560    the POS pattern *"JJ" "JJ" "JJ" "NN" "NNS"* was missing.

561    Preventing errors from source (1) requires improvement of POS taggers. Preventing errors from

562    source (3) requires a better word continuation representation manner instead of using hyphen, in

563    order to avoid confusion with hyphens used for conjoining noun modifiers. Preventing errors

564     from sources (2) and (4) could be partially prevented by further developing the exclusion word

565     list and POS pattern set, respectively. The use of the developed exclusion word list (to prevent

566     errors from source (2)) prevented 6 instances of false positives and increased precision from 89.0%

567     to 89.6%. More terms could be added, iteratively, to the exclusion word list to further enhance

568     performance. Similarly, errors from source (4) could be prevented by adding more patterns to the

569     POS pattern set until all possible POS patterns are included. While theoretically this POS pattern

570     set is infinite (e.g., infinite number of "JJ" before a "NN"), in practice this POS pattern set is

571     quite limited [e.g., words with more than 7 prenominal modifiers (e.g., white thin high strong

572     stone north exterior ancient wall) are seldom (if not never) seen].

573     To test the effect of iterative development of the exclusion word list and POS pattern set, three

574     more experiments were conducted to: (1) add the false positive present participle terms

575     (identified as a result of initial testing) to the exclusion word list and use it in further testing; (2)

576     add the missing POS patterns (identified as a result of initial testing) to the pattern set and use it

577     in further testing; and (3) use both, the extended exclusion word list and the extended POS

578     pattern set, in further testing. Table 7 shows the performance results of the three experiments.

579     The results show that the use of the extended exclusion word list and the POS pattern set both

580     improve the performance of concept extraction, with the latter showing a larger improvement.

581                                    Insert Table 6

582                                    Insert Table 7

583     *IFC Concept Selection*

584     Gold Standard

585     The gold standards of F-concepts for Chapter 12 of IBC 2006 and Chapter 19 of IBC 2009 were

586     manually developed by the authors. The F-concepts were initially identified using the matching

587    and ranking algorithms and then manually filtered. The gold standards of Chapter 12 and

588    Chapter 19 include 343 and 588 F-concepts, respectively.

589    Algorithm Implementation

590    The proposed IFC concept selection method and algorithms were implemented in Python

591    programming language (v.2.7.3). The Porter Stemmer (Porter 1980) was used for stemming. The

592    "re" (regular expression) module in python was utilized to support the matching algorithms. The

593    hypernymy, hyponymy, and synonymy relations in WordNet were utilized through the Natural

594    Language Toolkit (NLTK) (Bird et al. 2009) WordNet interface in python.

595    Evaluation

596    IFC concept selection was evaluated in terms of adoption rate. Adoption rate is defined as the

597    number of automatically selected F-concepts that were adopted divided by the total number of

598    automatically selected F-concepts.

599    Development Results and Analysis

600    For term-based matched F-concepts, Table 8 shows the results of testing combinations of

601    different concept-level SS scoring functions and term-level SS scoring functions. Table 9 shows

602    some example concepts that were extracted and matched using the different combinations. For

603    concept-level SS scoring, Eq.1 and Eq. 2 were tested. As shown in Table 8, Eq. (1) consistently

604    outperformed Eq. (2). Eq. (1) prefers shorter F-concepts and, thus, tends to select F-concepts that

605    are higher in the concept hierarchy (most likely a superclass). In comparison, Eq. (2) prefers F-

606    concepts with similar length to the R-concept and, thus, tends to select F-concepts that are at a

607    similar level in the concept hierarchy to the R-concept. However, an F-concept located at a

608    similar level to the R-concept may deviate a lot in meaning because concepts at similar level in a

609    concept hierarchy could belong to different branches of the hierarchy. A matched higher-level F-

610    concept, thus, usually has higher relatedness to the R-concept than a matched similar-level F-

611    concept. For example, using Shortest Path Similarity (for term-level SS scoring), Eq. (1) resulted

612    in matching of "net_free_ventilating_area" and "quantity_area", whereas Eq. (2) resulted in the

613    matching of "net_free_ventilating_area" and "annotation_fill_area_occurrence". "Quantity_area"

614    was correctly a superconcept of "net_free_ventilating_area" and was adopted. On the other hand,

615    the    meaning    of    "annotation_fill_area_occurrence"    was    far    from    that    of

616    "net_free_ventilating_area" despite being at a similar level in the concept hierarchy. Based on

617    these experimental results, Eq. (1) was selected for concept-level SS scoring for term-based

618    matched F-concepts.

619    For term-level SS scoring, the following five existing SS measures were tested: Shortest Path

620    Similarity, Jiang-Conrath Similarity, Leacock-Chodorow Similarity, Resnik Similarity, and Lin

621    Similarity (see Background section). The Shortest Path Similarity is the simplest among the five

622    tested measures, and achieved the best adoption rate of 86.5%. The Shortest Path Similarity and

623    Leacock-Chodorow Similarity are based on shortest path between two concepts in a taxonomy.

624    The other three SS measures are based on information content of the two concepts' least

625    common subsume (i.e., the lowest-level concept that is a superconcept of both concepts). The

626    performance drop from the Shortest Path Similarity to the other similarity measures (except for

627    Leacock_Chodorow Similarity) shows the advantage of a shortest path measure in comparison to

628    an information content of the least common subsumer measure. Empirically, this is because the

629    length of path between two concepts is more distinctive than the information content of their

630    least common subsumer. For example, in the concept hierarchy of Fig. 8, the shortest paths

631    between C2 and C5, C4 and C5, C7 and C5 are different, but their least common subsumers are

632    all the same (i.e., C1). For shortest path measures, the Leacock-Chodorow Similarity takes the

633    depth of the taxonomy into consideration, in addition to the use of shortest path. The

634    performance drop from the Shortest Path Similarity to the Leacock-Chodorow Similarity

635    indicates that the absolute taxonomy depth is not a distinctive feature in the context of concept

636    matching. Based on these experimental results, the Shortest Path Similarity was selected for

637    term-level SS scoring for term-based matched F-concepts.

638                                    Insert Figure 8

639                                    Insert Table 8

640                                    Insert Table 9

641    For semantic-based matched F-concepts, Table 10 shows the results of testing combinations of

642    different concept-level SS scoring functions and term-level SS scoring functions. Table 11 shows

643    some examples of concepts that were extracted and matched using the different combinations. As

644    shown in Table 10, for concept-level SS scoring, Eq. (1) and Eq. (2) did not show any variability

645    in performance. Since both functions performed equally, for consistency with term-based

646    matching of F-concepts, Eq. 1 was selected for concept-level SS scoring for semantic-based

647    matched F-concepts.

648    For term-level SS scoring, the Shortest Path Similarity outperformed all other SS measures. This

649    is consistent with the results obtained for term-based matched F-concepts. Based on the

650    experimental results, the Shortest Path Similarity was selected for term-level SS scoring for

651    semantic-based matched F-concepts.

652    Thus, the same term-level SS scoring function (Shortest Path Similarity) and concept-level SS

653    scoring function (Eq. 1) were selected for both term-based matching and semantic-based

654    matching algorithms. This shows consistency of performance across both types of matching.

655     Insert Table 10

656     Insert Table 11

657     Testing Results and Analysis

658     The proposed IFC concept selection method and algorithms [using Eq. (1) and Shortest Path

659     Similarity] were tested in automatically selecting F-concepts for the extracted R-concepts (from

660     Phase I). The testing results are summarized in Table 12. The total adoption rate is 84.5%. The

661     adoption rates for term-based and semantic-based matched F-concepts are 84.8% and 82.7%,

662     respectively, both which are close to the training performance (87.1% and 82.5%, respectively).

663     This shows initial stability in the performance of the proposed IFC concept selection method.

664     Insert Table 12

665     ***Relationship Classification***

666     Gold Standard

667     The aim of the classifier is to predict the relationship between each pair of R-concept and F-

668     concept. Two gold standards, one for training and one for testing, were manually developed by

669     the authors and three other graduate students in construction management. The training and

670     testing gold standards included pairs of concepts from Chapter 12 of IBC 2006 and Chapter 19 of

671     IBC 2009, respectively. The training data set was used for feature selection, ML algorithm

672     selection, and classifier training, and the testing data set for evaluating the classifier's

673     performance. In each gold standard, the relationship between each R-concept and F-concept was

674     defined. Four types of relationships were defined, as per Table 2.

675     Algorithm Implementation

676     The proposed relationship classification algorithms were developed and tested in Waikato

677     Environment for Knowledge Analysis (Weka) data mining software system (Hall et al. 2009). A

30

678    program for generating the ML features was developed using Python programming language

679    (v.2.7.3). The following ML algorithms were tested: (1) weka.classifiers.bayes.NaiveBayes for

680    Naïve Bayes; (2) weka.classifiers.trees.J48 for Decision Tree; (3) weka.classifiers.lazy.IBk for k-

681    NN; and (4) weka.classifiers.functions.SMO for SVM. Tenfold cross-validation was applied to

682    each training experiment, which randomly split the data to training subset and testing subset ten

683    times and averaged the results from the ten trials of training and testing.

684    Evaluation

685    Relationship classification was evaluated in two ways: (1) the performance across all

686    relationships was evaluated, together, in terms of precision, and (2) the performance for each

687    type of relationship was evaluated, separately, in terms of precision, recall, and F1-measure. In

688    the first case, precision is defined as the number of correctly classified concept pairs divided by

689    the total number of classified concept pairs. In the second case, precision is defined as the

690    number of correctly classified concept pairs in a relationship type divided by the total number of

691    concept pairs that are classified into that relationship type. Recall is defined as the correctly

692    classified concept pairs in a relationship type divided by the total number of concept pairs that

693    should be classified into that relationship type. F1-measure is the harmonic mean of precision

694    and recall.

695    ML Algorithm Selection, Feature Selection, and Classifier Training

696    The training data set was used for feature selection and classifier training. The results of testing

697    the four ML algorithms are summarized in Table 13. While three out of the four ML algorithms

698    achieved a precision greater than 85%, k-NN achieved the best precision of 90.98%. Decision

699    Tree ranked second in performance, with 86.07% precision.

31

700 A "leave-one-out" feature analysis was used for feature selection. Feature selection, in this paper,

701 aims at selecting – based on performance – a subset (or the full set) of the complete/initial

702 feature set (the eight features, see Table 3) for use in representing the concepts. The "leave-one-

703 out" feature analysis is a method to analyze the contribution of each feature by comparing the

704 performance with and without that feature. The analysis was conducted using the top-two

705 performing ML algorithms (Decision Tree and k-NN). The feature analysis results are

706 summarized in Table 14. The bold highlighted values indicate the precision values that

707 outperformed the baseline precision (underlined, where all eight features were used). The results

708 show that four out of the eight features (RTermNum, RTermPOS, RelMatchType, FTermNum)

709 were not discriminating when using Decision Tree, and one out of the eight features (FTermNum)

710 was not discriminating when using k-NN. Using only the discriminating features (i.e., Features

711 RMatchType, FMatchType, FTermPOS, and DOM for Decision Tree, and Features RTermNum,

712 RTermPOS, RMatchType, RelMatchType, FMatchType, FTermPOS, and DOM for k-NN),

713 Decision Tree achieved a precision of 87.43% and k-NN achieved a precision of 91.26%. This

714 difference shows that, in comparison to Decision Tree, k-NN was able to achieve a higher

715 performance with a larger feature size. This may indicate that the additional features used by k-

716 NN provided better discriminating ability to the classifier. As such, based on the experimental

717 results, the above-mentioned seven discriminating features and the k-NN algorithm were selected

718 for training the classifier.

719 The results also show that the following four features were discriminating for both algorithms:

720 RMatchType, FMatchType, FTermPOS, and DOM. DOM was discriminating because a term-

721 by-term match could provide a strong indication of concept equivalency. The fact that

722 RMatchType and FMatchType were discriminating shows that the arrangement of terms could

723 affect the meanings of concepts and that the locations of the matching terms in a concept pair

724 could affect the relationship between the two concepts in the pair. In addition to these four

725 features, the following three features were discriminating for k-NN: RTermNum, RTermPOS,

726 and RelMatchType. The fact that these features were discriminating in k-NN but not in Decision

727 Tree may attributed to the different types of ML algorithms. More importantly, the fact that the

728 RelMatchType is discriminating shows that the semantic features could benefit the task of

729 concept relationship classification and result in further improvement of precision.

730 Insert Table 13

731 Insert Table 14

732 Testing Results and Analysis

733 The testing data set was used for testing and evaluating the performance of the classifier. The

734 testing results are summarized in Table 15. The overall precision across all relationships is

735 87.94%. This is close to the overall training precision (90.98%), which shows the initial stability

736 in the performance of the relationship classifier. The subconcept relationship type achieved the

737 best precision of 93.4% and best recall of 93.4%. The analysis of the results shows that in many

738 cases the R-concept was a bigram or multigram (e.g., "structural concrete") whose last term

739 matched with the only term in a unigram F-concept (e.g., "concrete"). This pattern has a strong

740 predictive effect. Comparing to the subconcept relationship type, the superconcept relationship

741 type shares a similar pattern but did not achieve a performance as high. The precision and recall

742 for the superconcept relationship type were 88.5% and 75.4%, respectively. One observation was

743 that the classifier tends to prefer subconcept relationship types over superconcept relationship

744 types, when both the R-concept and the F-concept were bigram or multigram. For example, there

745 were six cases where a superconcept relationship was incorrectly classified as a subconcept

746   relationship, but zero cases where a subconcept relationship was incorrectly classified as a

747   superconcept relationship.  This could be due to the fact that there were only two instances of

748   bigram/multigram concept pairs with superconcept relationship in the training data set. The

749   equivalent relationship type achieved a precision of 91.9% and recall of 86.1%. The associated

750   relationship type achieved a precision of 62.5% and recall of 80.0%, which is the lowest among

751   the four types of relationships. This is probably because: (1) the size of the training data was

752   limited for this relationship type, and (2) the associated relationship includes more semantic

753   types than the other types of relationships and has more variability in the expression of concepts.

754   Thus, while the data set might provide enough variability for concepts related to the other

755   relationship types, the associated relationship may require more data. Overall, the precision is

756   87.94%, which is considered a good performance [within the range of 80% to 90% (Spiliopoulos

757   et al. 2010)].

758                                         Insert Table 15

759   **Limitations and Future Work**

760   Two limitations of this work are acknowledged, which the authors plan to address as part of their

761   ongoing/future research. First, due to the large amount of manual effort required in developing

762   the gold standard for each phase, the proposed method was only tested on one Chapter of IBC

763   2009. Similar good performance is expected on other chapters of IBC and other regulatory

764   documents. However, different performance results might be obtained due to the possible

765   variability of contents across different chapters of IBC 2009 or across different types of

766   regulatory documents. As such, in their future work, the authors plan to test the proposed method

767   on more chapters of IBC 2009 and on other types of regulatory documents (e.g., EPA

768   regulations). Second, only unigram (single terms) semantic-based matching was used for finding

34

769    semantically related F-concepts to an R-concepts. While the combinatorial nature of term

770    meanings [i.e., the meanings of single terms (e.g., "exterior" and "door") in a concept name are

771    combined to form the overall meaning of the whole concept (e.g., "exterior door")] renders this

772    unigram method effective, there may be cases where bigram (pairs of terms) or multigram

773    (groups of three or more terms) matching could be effective. As such, in future work, the authors

774    plan to extend the semantic-based matching method to incorporate semantic relations between

775    bigram and multigram to test whether such bigram or mutligram considerations could further

776    improve the performance of concept matching.

777    **Contributions to the Body of Knowledge**

778    This study contributes to the body of knowledge in four main ways. First, this research offers a

779    method for automated concept extraction that utilizes POS-pattern-matching-based rules to

780    extract regulatory concepts from natural language regulatory documents. The set of POS patterns

781    that was developed captures natural language knowledge, which allows for the recognition of

782    concepts based on the lexical and functional categories of their terms. The pattern set includes

783    only flattened patterns to avoid recursive parsing, which allows for efficient computation. The

784    set of POS patterns are also generalized, and thus can be used to extract concepts in other

785    domains.  Second, this research offers a matching-based method for identifying and selecting the

786    most related IFC concepts to the extracted regulatory concepts. The proposed method leverages

787    both syntactic and semantic knowledge, which allows for the recognition of related concept pairs

788    based on the syntactic and semantic similarities of their terms. As part of this method, two new

789    concept-level semantic similarity (SS) scoring functions are offered. In the context of schema

790    extension, existing SS scoring functions allow for measuring SS at the term-level. These

791    proposed two functions further allow for measuring SS at the concept-level. Third, this research

792    offers an automated machine learning classification method for classifying the relationships

793    between the extracted regulatory concepts and their most related IFC concepts. The classification

794    results show that semantic features could benefit the task of relationship classification and result

795    in further improvement of precision. The proposed method is also generalized and can be used to

796    classify the relationships between any two concepts, based on eight syntactic and semantic

797    features of their terms, into the following four types: "equivalent concept", "superconcept",

798    "subconcept", and "associated concept". Fourth, the experimental results show that the three

799    proposed methods could be effectively combined in a sequential way for extending the IFC

800    schema with regulatory concepts from regulatory documents. This offers a new method for

801    objectively extending the IFC schema with domain-specific concepts that are extracted from

802    natural language documents. The proposed combined method is also generalized and can be used

803    to extend the IFC schema with other types of concepts (e.g., environmental concepts) from other

804    types of documents (e.g., environmental documents) or to extend other types of class hierarchies

805    (e.g., of an ontology) in the construction domain or in other domains.

806    **Conclusions**

807    This paper presented a new method for extending the IFC schema with regulatory concepts from

808    relevant regulatory documents for supporting automated compliance checking. The proposed

809    method utilizes semantic natural language processing (NLP) techniques and machine learning

810    techniques, and is composed of three primary methods that are combined into one computational

811    platform: (1) a method for concept extraction that utilizes POS-pattern-matching-based rules to

812    extract regulatory concepts from regulatory documents, (2) a method for identifying and

813    selecting the most related IFC concepts to the extracted regulatory concept, which utilizes term-

814    based and semantic-based matching algorithms to find candidate related IFC concepts and a

815  semantic similarity (SS) scoring and ranking algorithm to measure the SS between each

816  candidate IFC concept a regulatory concept, and (3) a machine learning classification method for

817  predicting the relationship between the extracted regulatory concepts and their most related IFC

818  concepts based on the syntactic and semantic features of their terms. The proposed IFC extension

819  method was evaluated on extending the IFC schema with regulatory concepts from Chapter 19 of

820  IBC 2009. Each of the three methods were evaluated separately, and achieved 91.7%, 84.5%,

821  and 87.94% F1-measure, adoption rate, and precision, respectively. The performance results

822  indicate that the proposed IFC extension method is potentially effective. The results also show

823  that semantic features of the concept terms and their interrelationships are helpful in IFC

824  extension and result in performance improvement. In their future work, the authors plan to: (1)

825  test the proposed method on other chapters of IBC 2009 and other construction regulatory

826  documents (e.g., EPA regulations); and (2) extend the semantic-based matching method to

827  incorporate semantic relations between bigram and multi-term to test whether such bigram or

828  multigram considerations could further improve the performance of concept matching.

829  **Acknowledgement**

834  **References**

835  AISC. (2014). "Technology integration." <http://www.aisc.org/content.aspx?id=26044>. (Aug

836      12, 2014).

837 Bird, S., Klein, E., and Loper, E. (2009). "Natural language processing with Python." O'Reilly

838      Media Inc., Sebastopol, CA.

839 Böhms, M., Bonsma, P., Bourdeau, M., and Kazi, A.S. (2009). "Semantic product modelling and

840      configuration: challenges and opportunities." *J. Inf. Techno. Constr.*, 14, 507-525.

841 BuildingSmart. (2014). "Industry Foundation Classes (IFC) data model." <

842      http://www.buildingsmart.org/standards/ifc/model-industry-foundation-classes-ifc> (Aug

843      12, 2014).

844 Bundy, A. (2013). "The interaction of representation and reasoning." *Proc., R. Soc. A*, 469(2157),

845      1-18.

846 Cover, T.M. (1967). "Nearest neighbor pattern classification." *IEEE Transactions on

847      Information Theory*, 13(1), 21-27.

848 Cristianini, N., and Shawe-Taylor, J. (2000). "An introduction to support vector machines and

849      other kernel-based learning methods." Cambridge University Press, 1st edition,

850      Cambridge, U.K.

851 Dietrich, S.W., and Urban, S.D. (2011). "Fundamentals of object databases: object-oriented and

852      object-relational design." *Synthesis Lectures on Data Management*, Morgan & Claypool

853      Publishers, San Rafael, California.

854 Domingos, P. (2012). "A few useful things to know about machine learning." *Communications

855      of the ACM*, 55(10), 78-87.

856 Delgado, F., Martínez-González, M.M., and Finat, J. (2013). "An evaluation of ontology

857      matching techniques on geospatial ontologies." *Int. J. Geogr. Inf. Sci.*, 27(12), 2279-2301.

858 Eastman, C., Lee, J., Jeong, Y., and Lee, J. (2009). "Automatic rule-based checking of building

859      designs." *Autom. Constr.*, 18(8), 1011-1033.

860    El-Gohary, N.M., and El-Diraby, T.E. (2010). "Domain ontology for processes in infrastructure

861         and construction." *J. Constr. Eng. Manage.*, 136(7), 730–744.

862    Fellbaum, C. (2005). "WordNet and wordnets." In: Brown, Keith et al. (eds.), *Encyclopedia of*

863         *Language and Linguistics*, Second Edition, Oxford: Elsevier, 665-670.

864    Freund, Y., and Schapire, R. (1999). "Large margin classification using the perceptron

865         algorithm." *Mach. Learn.*, 37(3), 277-296.

866    Fritz, D. (2006). "The Semantic Model: A basis for understanding and implementing data

867         warehouse requirements." <http://www.tdan.com/view-articles/4044> (Aug 12, 2014).

868    Gruber, T.R. (1995). "Toward principles for the design of ontologies used for knowledge

869         sharing." *Int. J. Hum.-Comput. St.*, 43, 907-928.

870    Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., and Witten, I.H. (2009). "The

871         WEKA data mining software: an update." *SIGKDD Explor.*, 11(1), 10-18.

872    Hamil, S. (2012). "Building information modelling and interoperability." <

873         http://www.thenbs.com/topics/bim/articles/bimAndInteroperability.asp> (Aug 12, 2014).

874    Hanis, T., and Noller, D. (2011). "The role of semantic models in smarter industrial operations."

875         *Oper.*, IBM Corporation, Armonk, New York.

876    Isikdag, U., Aouad, G., Underwood, J, and Wu, S. (2007). "Building information models: a

877         review on storage and exchange mechanisms." *Proc., 24th W78 Conf. & 5th ITCEDU*

878         *Workshop & 14th EG-ICE Workshop, Bringing ITC Knowledge to Work*, International

879         Council for Research and Innovation in Building and Construction (CIB), Rotterdam,

880         Netherlands, 135-144.

881    Ji, Y., Borrmann, A., Beetz, J., and Obergrießer, M. (2013). "Exchange of parametric bridge

882         models using a neutral data format." *J. Comput. Civ. Eng.*, 27(6), 593–606.

883    Jiang, Y., Wang, X, and Zheng, H. (2014). "A semantic similarity measure based on information

884        distance for ontology alignment." *Inf.. Sci.*, 278(2014), 76-87.

885    Kasim, T., Li, H., Rezgui, Y., and Beach, T. (2013). "Automated sustainability compliance

886        checking process: proof of concept." *Proc., 13th Int. Conf. Constr. App. Vir. Real.*,

887        Teesside University, Tees Valley, UK, 11-21.

888    Kamps, J., Marx, M., Mokken, R.J., and Rijke, M. (2004). "Using WordNet to measure semantic

889        orientations of adjectives." *Proc., LREC-04,* European Language Resources Association

890        (ELRA), Paris, France, 1115-1118.

891    Khemlani, L. (2005). "CORENET e-PlanCheck: Singapore's automated code checking system."

892        AECbytes        "Building        the        Future"        Article,

893        <http://www.novacitynets.com/pdf/aecbytes_20052610.pdf> (Aug 12, 2014).

894    Klas,    W.,    and    Schrefl,    M.    (1995).    "Metaclasses and their application -

895        Data model tailoring and database integration." *Lect. Notes. Comput. Sc.*, 943, Springer-

896        Verlag, Berlin Heidelberg, Germany, 1-7.

897    Li,    T.    (2010).    "Practice    and    exploration    of    ontology    creation    algorithms."    <

898        http://www.cs.ubc.ca/~carenini/TEACHING/CPSC503-14/FINAL-REPORTS-

899        10/CPSC503_Project_Report_Tianyu.pdf > (Mar. 21st, 2015).

900    Ma, Z., Lu, N., and Song, W. (2011a). "Identification and representation of information

901        resources for construction firms." *Adv. Eng. Inform..*, 25(4), 2011, 612-624.

902    Ma, Z., Wei, Z., Song, W, Lou, Z. (2011b). "Application and extension of the IFC standard in

903        construction cost estimating for tendering in China." *Autom. Constr.*, 20(2), 196-204.

904    Martins, J.P., and Monteiro, A. (2013). "LicA: A BIM based automated code-checking

905        application for water distribution systems." *Autom. Constr.*, 29(2013), 12–23.

906    National Institute of Building Sciences. (2014). "National BIM standard – United States Version

907        2." < http://www.nationalbimstandard.org/faq.php#faq1> (Aug 12, 2014).

908    Nawari, N.O. (2011). "Automating codes conformance in structural domain." *Proc., Comput.*

909        *Civ. Eng.*, ASCE, Reston, VA, 569-577.

910    Nguyen,T., and Kim, J. (2011). "Building code compliance checking using BIM

911        technology." *Proc., 2011 Winter Simulation Conference (WSC)*, IEEE, New York, NY,

912        3395-3400.

913    Niemeijer, R.A., Vries, B. D., and Beetz, J. (2009). "Check-mate: automatic constraint checking

914        of IFC models." In A Dikbas, E Ergen & H Giritli (Eds.), *Manag. IT in Constr. Manag.*

915        *Constr. for Tomorrow*, CRC Press, London, UK, 479-486.

916    Orna-Montesinos, C. (2010). "Hyponymy relations in construction textbooks: a corpus-based

917        analysis." *Linguistic and Translation Studies in Scientific Communication, Linguistic*

918        *Insights*, 86, 96-114.

919    Pan, J., Cheng, C.J., Lau, G.T., and Law, K.H. (2008). "Utilizing statistical semantic similarity

920        techniques for ontology mapping - with applications to AEC standard models." *Tsinghua*

921        *Sci. and Technol.*, 13(S1), 217-222.

922    Porter, M. (1980). "An algorithm for suffix stripping." *Program (Autom. Libr. and Inf. Syst.)*,

923        14(3), 130-137.

924    Resnik, P. (1995). "Using information content to evaluate semantic similarity in a taxonomy."

925        *Proc., IJCAI'95*, IJCAI, Inc., Somerset, NJ, 448-453.

926    Rodrı´guez, M.A., and Egenhofer, M.J. (2003). "Determining semantic similarity among entity

927        classes from different ontologies." *IEEE T. Knowl. Data En.,* 15(2), 442-456.

928    Rosenblatt, F. (1958). "The perceptron: a probabilistic model for information storage and

929          organization in the brain." *Psychol. Rev.*, 65(6), 1958.

930    Salama, D. and El-Gohary, N. (2013). "Semantic text classification for supporting automated

931          compliance checking in construction." *J. Comput. Civ. Eng.*, 10.1061/(ASCE)CP.1943-

932          5487.0000301 (Feb. 23, 2013).

933    Shehata, S. (2009). "A WordNet-based semantic model for enhancing text clustering." *2009*

934          *IEEE Int. Conf. Data Mining. Workshops*, IEEE, Piscataway, NJ, 477-482.

935    Simpson, T., and Dao, T. (2010). "WordNet-based semantic similarity measurement." <

936          http://www.codeproject.com/Articles/11835/WordNet-based-semantic-similarity-

937          measurement> (Aug 14, 2014).

938    Sinha, S., Sawhney, A., Borrmann, A., and Ritter, F. (2013). "Extracting information from

939          building information models for energy code compliance of building envelope." *COBRA*

940          *2013 Conf.*, International Council for Research and Innovation in Building and

941          Construction (CIB), Rotterdam, Netherlands.

942    Song, W., Liang, J.Z., and Park, S.C. (2014). "Fuzzy control GA with a novel hybrid semantic

943          similarity strategy for text clustering." *Inform. Sciences*, 273(2014), 156-170.

944    Soysal, E., Cicekli, I., and Baykal, N. (2010). "Design and evaluation of an ontology based

945          information extraction system for radiological reports." *Comput. in Biology and Med.*,

946          40(11-12), 900-911.

947    Spiliopoulos, V., Vouros, G.A., and Karkaletsis, V. (2010). "On the discovery of subsumption

948          relations for the alignment of ontologies." *Web Semant.*, 182, 1-20.Slimani, T. (2013).

949          "Description and evaluation of semantic similarity measures approaches." *Int. J. Comput.*

950          *Appl.*, 80(10), 25-33.

951    Suchanek, M., Kasneci, G., and Weikum, G. (2007). "YAGO: A core of semantic knowledge

952           unifying WordNet and Wikipedia." *Proc., WWW 2007*, Association for Computing

953           Machinery, New York, NY, 697-706.

954    Tan, X., Hammad, A., and Fazio, P. (2010). "Automated code compliance checking for building

955           envelope    design."    *J.    Comput.    Civ.    Eng.*,    10.1061/1192    (ASCE)0887-

956           3801(2010)24:2(203), 203–211.

957    Toutanova, K., Klein, D., Manning, C., and Singer, Y. (2003). "Feature-rich part-of-speech

958           tagging with a cyclic dependency network." *Proc., HLT-NAACL 2003*, 252-259.

959    Vanlande, R., Nicolle, C., and Cruz, C. (2008). "IFC and building lifecycle management." *Autom.*

960           *Constr.,* 18(1), 70-78.

961    Varelas, G., Voutsakis, E., and Raftopoulou, P. (2005). "Semantic similarity methods in

962           WordNet and their application to information retrieval on the web." *Proc., 7th annual*

963           *ACM intl. workshop on Web inform. and data manage. (WIDM '05)*, Association for

964           Computing Machinery, New York, NY, 10-16.

965    Young, N.W., Jr., Jones, S.A., Bernstein, H.M., and Gudgel, J.E. (2009). "The business value of

966           BIM: getting building information modeling to the bottom line." The McGraw-Hill

967           Companies, New York, NY.

968    Zhang, J., and El-Gohary, N. (2013). "Semantic NLP-based information extraction from

969           construction regulatory documents for automated compliance checking." *J. Comput. Civ.*

970           *Eng.*, Accepted and published online ahead of print.

971    Zhang, J., Yu, F., Li, D., and Hu, Z. (2014). "Development and implementation of an Industry

972           Foundation Classes-based graphic information model for virtual construction." *Comput.-*

973           *Aided Civ. Inf.*, 29(2014), 60-74.

974    Zhou, P., and El-Gohary, N. (2014). "Ontology-based multi-label text classification for enhanced

975        information retrieval for supporting automated environmental compliance checking."

976        *Proc., 2014 ASCE Constr. Res. Congress (CRC)*, ASCE, Reston, VA, 2238-2245.

977
978
979
980
981
982
983
984
985

986    **Tables**

987    Table 1. Commonly-used Machine Learning Algorithms

|  | Machine learning algorithm | | | | |
|---|---|---|---|---|---|
|  | Naïve Bayes | Perceptron | Decision Tree | k-NN | SVM |
| Key feature | simple but effective | linear | flexible | similarity-based | kernel-based |

988
989    Table 2. Types of Relationships Considered

| Relationship type | Relationship interpretation |
|---|---|
| Equivalent concept | $R^1$ is equivalent to $F^2$ |
| Superconcept | R is superconcept of F |
| Subconcept | R is subconcept of F |
| Associated concept | R and with F are associated |

990    [1] R means R-concept
991    [2] F means F-concept

992
993    Table 3. The Syntactic and Semantic Features used for the Relationship Classifier

| Feature name | RTermNum | RTermPOS | RMatch Type | RelMatch Type | FMatchType | FTermNum | FTerm POS | DOM |
|---|---|---|---|---|---|---|---|---|
| Possible values | unigram, bigram, multigram | N, NN, JN | first, last | synonym, hypernym, hyponym, term-based | first, middle, last | unigram, bigram, multigram | N, NN, JN | 1, 0 |

994
995    Table 4. Example R-concepts, Matched F-concepts, and Their Feature Values

| R-concept | F-concept | RTerm Num | RTerm POS | RMatch Type | RelMatch Type | FMatch Type | FTermNum | FTerm POS | DOM |
|---|---|---|---|---|---|---|---|---|---|
| construction | construction resource | unigram | N | first | term-based | first | bigram | NN | 0 |
| floor joist | beam | bigram | N | last | synonym | first | unigram | NN | 0 |
| preconstruction | preconstruction | bigram | NN | first | term- | first | bigram | NN | 1 |

| testing | test | | | | based | | | | |
|---|---|---|---|---|---|---|---|---|---|
| skylight | window | unigram | N | first | synonym | first | unigram | N | 0 |
| water-proof joint | structural connection | bigram | NN | last | hypernym | last | bigram | JN | 0 |

996

997 Table 5. Performance of Extracting Regulatory Concepts from Development Text (Chapter 12 of
998 IBC 2006)

| Method | Number in gold standard | Number extracted | Number correctly extracted | Precision | Recall | F1-measure |
|---|---|---|---|---|---|---|
| Without exclusion word list | 368 | 391 | 365 | 93.4% | 99.2% | 96.2% |
| With exclusion word list | 368 | 376 | 365 | 97.1% | 99.2% | 98.1% |

999 Table 6. Performance of Extracting Regulatory Concepts from Testing Text (Chapter 19 of IBC
1000 2009)

| Method | Number in gold standard | Number extracted | Number correctly extracted | Precision | Recall | F1-measure |
|---|---|---|---|---|---|---|
| Without exclusion word list | 821 | 871 | 773 | 88.7% | 94.2% | 91.4% |
| With exclusion word list | 821 | 865 | 773 | **89.4%** | 94.2% | **91.7%** |

1001

1002 Table 7. Performance of Regulatory Concept Extraction after Improvements

| Method | Number in gold standard | Number extracted | Number correctly extracted | Precision | Recall | F1-measure |
|---|---|---|---|---|---|---|
| Baseline Condition (from Table 6) | 821 | 865 | 773 | 89.4% | 94.2% | 91.7% |
| With extended exclusion word list | 821 | 856 | 774 | 90.4% | 94.3% | 92.3% |
| With extended POS pattern set | 821 | 860 | 784 | 91.2% | 95.5% | 93.3% |
| With both extended exclusion word list and extended POS pattern set | 821 | 851 | 785 | **92.2%** | **95.6%** | **94.0%** |

1003

1004 Table 8. Performances of Different SS Scoring Methods for Term-Based Matched F-Concepts

| Proposed concept-level SS scoring function | Term-level SS scoring function | Number of related F-concepts found | Number of related F-concepts adopted | Adoption rate |
|---|---|---|---|---|
| Eq. (1) | Shortest Path Similarity | 286 | 249 | **87.1%** |
| Eq. (2) | Shortest Path Similarity | 286 | 225 | 78.7% |
| Eq. (1) | Jiang-Conrath Similarity | 286 | 244 | 85.3% |
| Eq. (2) | Jiang-Conrath Similarity | 286 | 224 | 78.3% |
| Eq. (1) | Leacock-Chodorow Similarity | 286 | 237 | 82.9% |
| Eq. (2) | Leacock-Chodorow Similarity | 286 | 202 | 70.6% |
| Eq. (1) | Resnik Similarity | 286 | 246 | 86.0% |
| Eq. (2) | Resnik Similarity | 286 | 228 | 79.7% |
| Eq. (1) | Lin Similarity | 286 | 246 | 86.0% |
| Eq. (2) | Lin Similarity | 286 | 224 | 78.3% |

1005

1006   Table 9. Examples of Matched R-Concepts and F-Concepts Using Different SS Scoring Methods
1007   for Term-Based Matched F-Concepts

| Extracted R-concept | Proposed concept-level SS scoring function | Matched F-concept using Shortest Path Similarity[1] | Matched F-concept using Leacock-Chodorow Similarity[1] | Matched F-concept using Jiang-Conrath Similarity[1] |
|---|---|---|---|---|
| adjacent dwelling unit | Eq. (1) | dwelling unit | *unit assignment* | *derived unit* |
|  | Eq. (2) | *context dependent unit* | *context dependent unit* | *context dependent unit* |
| square foot | Eq. (1) | feet | *footing* | feet |
|  | Eq. (2) | feet | *footing* | feet |
| international energy conservation code | Eq. (1) | code | code | code |
|  | Eq. (2) | international mechanical code | international mechanical code | international mechanical code |
| required ventilating area | Eq. (1) | area | area | area |
|  | Eq. (2) | *annotation fill area* | *annotation fill area* | *annotation fill area* |
| exterior walls | Eq. (1) | wall | wall | wall |
|  | Eq. (2) | curtain wall | curtain wall | curtain wall |

1008   [1] italicized concepts were not adopted

1009   Table 10. Performances of Different SS Scoring Methods for Semantic-Based Matched F-
1010   Concepts

| Proposed concept-level SS scoring function | Term-level SS scoring function | Number of related F-concepts found | Number of related F-concepts adopted | Adoption rate |
|---|---|---|---|---|
| Eq. (1) | Shortest Path Similarity | 114 | 94 | **82.5%** |
| Eq. (2) | Shortest Path Similarity | 114 | 94 | **82.5%** |
| Eq. (1) | Jiang-Conrath Similarity | 114 | 92 | 80.7% |
| Eq. (2) | Jiang-Conrath Similarity | 114 | 92 | 80.7% |
| Eq. (1) | Leacock-Chodorow Similarity | 114 | 93 | 81.6% |
| Eq. (2) | Leacock-Chodorow Similarity | 114 | 93 | 81.6% |
| Eq. (1) | Resnik Similarity | 114 | 93 | 81.6% |
| Eq. (2) | Resnik Similarity | 114 | 93 | 81.6% |
| Eq. (1) | Lin Similarity | 114 | 93 | 81.6% |
| Eq. (2) | Lin Similarity | 114 | 93 | 81.6% |

1011

1012   Table 11. Examples of Matched R-Concepts and F-Concepts Using Different SS Scoring
1013   Methods for Semantic-Based Matched F-Concepts

| Extracted R-concept | Proposed concept-level SS scoring Function | Matched F-concept using Shortest Path Similarity[1] | Matched F-Concept using Leacock-Chodorow Similarity[1] | Matched F-Concept using Jiang-Conrath Similarity[1] |
|---|---|---|---|---|
| corrosion-resistant wire cloth screening | Eq. (1) | hardware cloth | hardware cloth | hardware cloth |
|  | Eq. (2) | hardware cloth | hardware cloth | hardware cloth |
| grab bars | Eq. (1) | railing | railing | railing |
|  | Eq. (2) | railing | railing | railing |
| outdoors | Eq. (1) | outside horizontal clear space | outside horizontal clear space | outside horizontal clear space |
|  | Eq. (2) | outside horizontal clear space | outside horizontal clear space | outside horizontal clear space |
| installed | Eq. (1) | *contaminant sources* | *contaminant source* | *light source* |

| shower heads | Eq. (2) | *contaminant sources* | *contaminant source* | *light source* |

1014 [1] italicized concepts were not adopted

1015

1016 Table 12. Testing Results of IFC Concept Selection Method

| Concept matching type | Concept-level SS scoring function | Term-level SS scoring function | Number of related F-concepts found | Number of related F-concepts adopted | Adoption rate |
|---|---|---|---|---|---|
| Term-based matching | Eq. (1) | Shortest Path Similarity | 598 | 507 | 84.8% |
| Semantic-based matching | | | 98 | 81 | 82.7% |
| Total | | | 696 | 588 | 84.5% |

1017

1018 Table 13. Results of Testing Different ML Algorithms

| Metric | ML algorithm | | | |
|---|---|---|---|---|
| | Naïve Bayes | Decision Tree | k-NN | SVM |
| Total number of relationship instances | 366 | 366 | 366 | 366 |
| Number of correctly classified relationship instances | 279 | 315 | 333 | 314 |
| Precision | 76.23% | 86.07% | **90.98%** | 85.79% |

1019

1020 Table 14. Leave-One-Out Feature Analysis Precision Results

| ML algorithm | Excluded feature | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | None | RTermNum | RTermPOS | RMatchType | RelMatchType | FMatchType | FTermNum | FTermPOS | DOM |
| Decision Tree | 86.07% | **86.89%** | **86.61%** | 81.15% | **86.61%** | 84.70% | **86.89%** | 86.07% | 81.98% |
| k-NN | 90.98% | 89.07% | 89.89% | 86.34% | 88.80% | 86.89% | **91.26%** | 90.44% | 88.25% |

1021

1022 Table 15. Relationship Classifier Testing Results

| Relationship type | Number of relationship instances in gold standard | Number of classified relationship instances | Number of correctly classified relationship instances | Precision | Recall | F1-Measure |
|---|---|---|---|---|---|---|
| Equivalent concept | 79 | 74 | 68 | 91.9% | 86.1% | 88.9% |
| Subconcept | 241 | 241 | 225 | 93.4% | 93.4% | 93.4% |
| Superconcept | 61 | 52 | 46 | 88.5% | 75.4% | 81.4% |
| Associated concept | 50 | 64 | 40 | 62.5% | 80.0% | 70.2% |
| Total | 431 | 431 | 379 | 87.94% | 87.94% | 87.94% |

1023

1024

1025

1026 Figure 1. Proposed IFC extension method

1027

1028
1029
1030
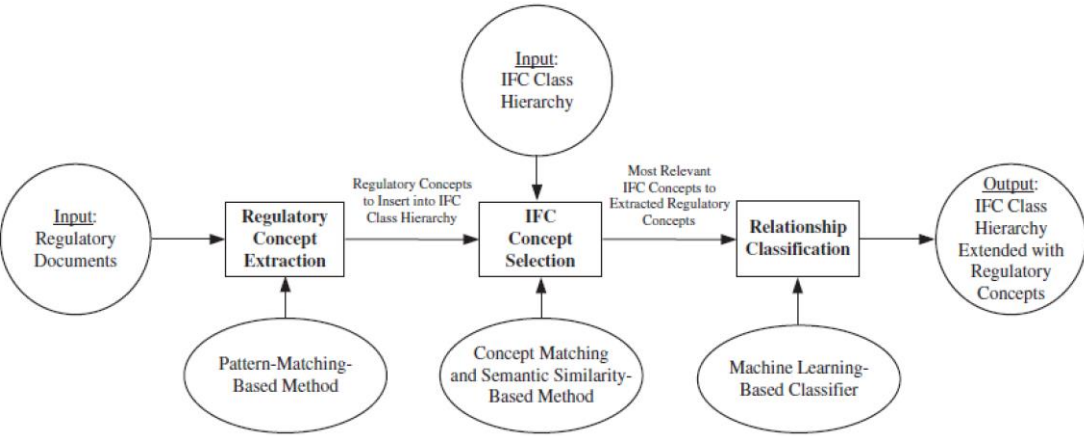1031     Figure 2. Example concept extraction rule and its meaning
1032

**Concept Extraction Rule**: RB VBN NN NN -> Extract the four matched terms

**Meaning**: If four consecutive terms are adverb, past participle verb, singular/mass noun, and singular/mass noun, then these four terms should be extracted as a concept.

| POS Tag | Meaning |
|---------|---------|
| NN | Singular or mass noun |
| NNS | Plural noun |
| NNP | Singular proper noun |
| NNPS | Plural proper noun |
| JJ | Adjective |
| RB | Adverb |
| VBN | Past participle verb |
| VBP | Non-3rd person singular present verb |
| VBD | Past tense verb |
| VBG | Gerund or present participle verb |

1033
1034
1035     **Figure 3.** Sample of Patterns and Concept Extraction Rules

- R1: NP -> NN NN
- R2: NP -> NN NP (non-flattened)
- R3: NP -> NN NN NN
- R4: NP -> NN NN NN NN
- R5: NP -> NN NN NN NN NN
- P1: NN NN
- P2: NN NN NN
- P3: NN NN NN NN
- P4: NN NN NN NN NN

1036
1037
1038
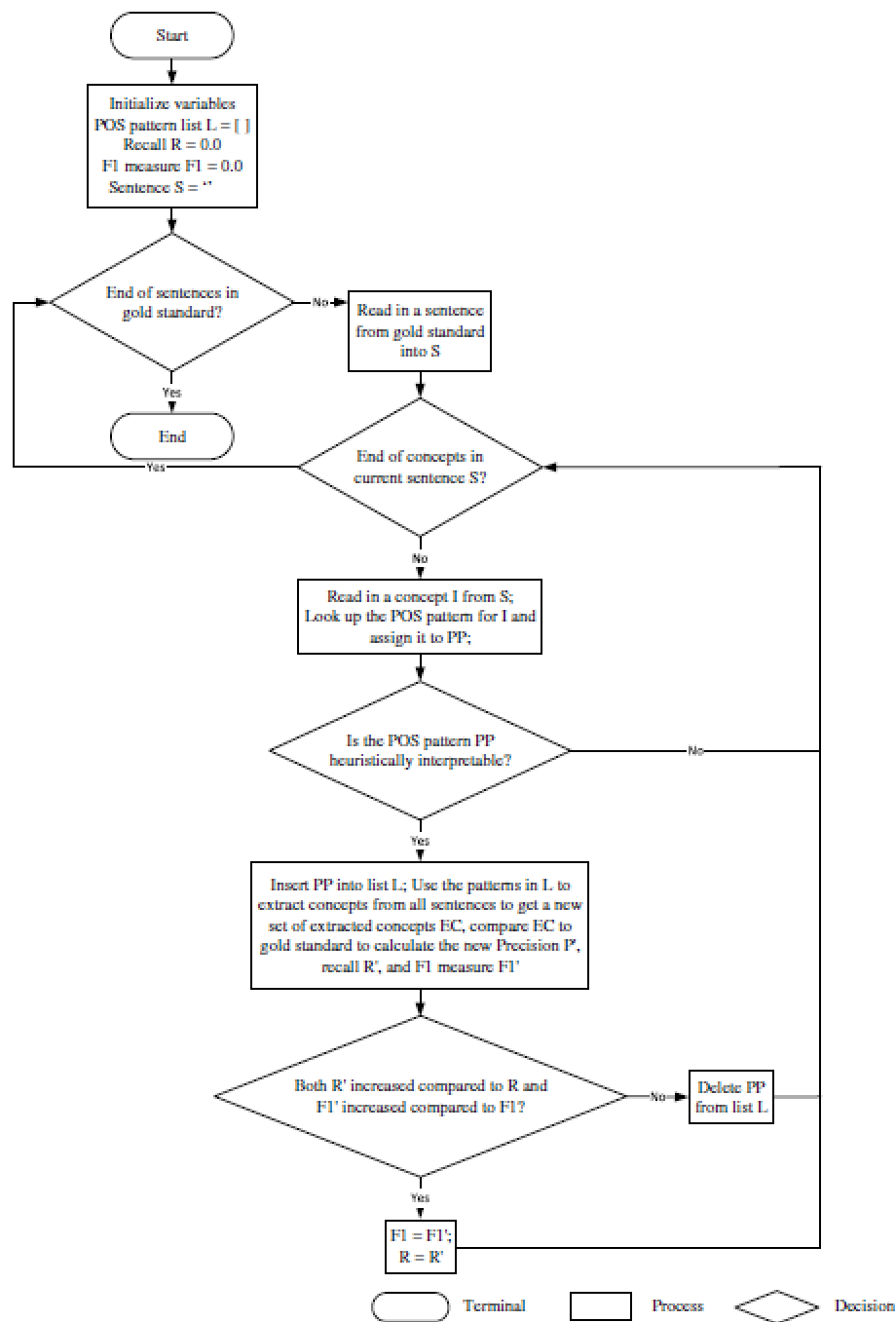1039     Figure 4. Flow chart of the POS pattern set development algorithm
1040

49

Figure 5. IFC concept selection method

50

1045
1046
1047        Figure 6. Term-based and semantic-based matching strategy
1048



1049
1050
1051        Figure 7. Set of POS patterns developed
1052

| NN | VBN NN | VBN NN NN |
|---|---|---|
| NNP | VBN NNS | VBN NN NNS |
| NNS | JJ JJ NN | JJ JJ JJ NN |
| VBG | JJ JJ NNS | JJ JJ NN NN |
| JJ NN | JJ NN NN | JJ NN NN NN |
| JJS NN | JJ NN NNS | NN NN NN NNS |
| JJ NNS | NN NN NN | NNP NN NN NN |
| NN NN | NN NN NNS | NNP NN NN NNS |
| NN NNS | NNP NN NN | NNP NNP NNP NNP |
| NNP NNP | NNP NNP NNP | RB VBN JJ NN |
| NNP NNS | NNP VBD NNS | RB VBN NN NN |
| VBG NN | NN VBG NN | RB VBN NNP NN |
| VBG NNS | RB JJ NNS | JJ VBN JJ JJ VBG NN |

1053
1054

1055 **Figure 8.** A sample concept hierarchy structure

1056



1057