# Integrating semantic NLP and logic reasoning into a unified system for fully-automated code checking

**2 authors**, including:

Jiansong Zhang
Western Michigan University

**15** PUBLICATIONS   **48** CITATIONS

1    **Integrating semantic NLP and logic reasoning into a unified system for fully-automated**
2    **code checking**

3    Jiansong Zhang[1]; and Nora M. El-Gohary[2]

4    **Abstract**

5    Existing automated compliance checking (ACC) systems are limited in their automation; they rely

6    on the use of hard-coded, proprietary rules for representing regulatory requirements, which

7    requires major manual effort in extracting regulatory information from textual regulatory

8    documents and coding these information into a rule format. To address this limitation, this paper

9    proposes a new unified ACC system that integrates: (1) semantic natural language processing

10    techniques and EXPRESS data based techniques to automatically extract and transform both

11    regulatory information (in regulatory documents) and design information [in building information

12    models (BIMs)] for automated compliance reasoning, and (2) semantic logic-based information

13    representation so that the reasoning could be fully automated. To test the proposed system, a BIM

14    test case was checked for compliance with Chapter 19 of the International Building Code 2009.

15    Comparing to a manually-developed gold standard, 98.7% recall and 87.6% precision in

16    noncompliance detection were achieved.

17    **Keywords:** Automated code checking; Automated information extraction; Automated reasoning;

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

19    Automated construction management systems.

[1] Assistant Professor, Dept. of Civil and Construction Engineering, Western Michigan University, 1903 W Michigan Ave, Kalamazoo, MI 49008.

[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.

20 **1    Introduction**

21    The manual process of regulatory compliance checking is time-consuming, costly, and error-prone

22    (Boken and Callaghan 2009). In the U.S., each building compliance review cycle usually takes

23    several weeks (State of New Jersey 2014; City of Philadelphia 2015), and a construction project

24    may be subject to multiple cycles of plan reviews due to design changes. At the city level, millions

25    of dollars are spent on manual building compliance checking each year (Department of Buildings

26    2015). Failure to comply with building regulations could further result in fines, penalties, or even

27    criminal court summons and prosecutions (Los Angeles Times 2015). Moreover, in an experiment

28    conducted by Fiatech, more disparity than agreement was found when different plan review

29    departments were asked to conduct manual code review of the same set of plans (Fiatech

30    Regulatory Streamlining Committee 2012).

31    In comparison to manual compliance checking, automated compliance checking (ACC) of

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

33    process (Eastman et al. 2009, Tan et al. 2010; Nguyen and Kim 2011; Kasim et al. 2013; Zhang

34    and El-Gohary 2013). However, the state-of-the-art ACC systems cannot achieve full automation

35    because of relying on the use of hard-coded, proprietary rules for representing regulatory

36    requirements, which requires major manual effort in extracting regulatory information from textual

37    regulatory documents and coding these information into a rule format. For example, the

38    COnstruction and Real Estate NETtwork (CORENET) project hard-coded rules in C++ programs,

39    the Solibri model checker uses a proprietary proforma-based format to hard code rules, and several

40    ACC efforts hard-coded rules for specific subdomains such as building evacuation (Choi et al.

41    2014), fall protection (Zhang et al. 2013), construction quality (Zhong et al. 2012), building safety

42    design (Qi et al. 2011), building envelope performance (Tan et al. 2010), and accessibility (Lau

2

43   and Law 2004). Such hard-coded rules could be very effective in reasoning about compliance with

44   a specific set of requirements and specific regulatory sections in a certain period of time, but such

45   rigid and static representation requires great effort in (1) adaptation to different regulatory

46   codes/sections, and (2) maintenance/update across different time periods and in response to code

47   revisions/updates. The use of hard-coded rules, thus, becomes effort-intensive and time-

48   consuming because of the large number of codes and regulations and their frequent

49   revisions/updates (Delis and Delis 1995; Dimyadi and Amor 2013).

50   In view of that, a number of researchers explored the development of generalized representations

51   for the formalization of regulatory requirements, with the aim to facilitate soft coding of rules for

52   supporting ACC. For example, Pauwels et al. (2011) proposed a semantic rule checking

53   environment, in which Notation 3 (N3) Logic is used to represent requirement rules. Hjelseth and

54   Nisbet (2011) proposed the Requirement, Applies, Select, and Exception (RASE) method to

55   represent regulatory requirements. Yurchyshyna et al. (2010; 2008) developed a conformity-

56   checking ontology that represents regulatory information, building-related knowledge, and expert

57   knowledge on checking procedures, with a representation of regulatory requirements in the form

58   of SPARQL Protocol and RDF Query Language (SPARQL) queries. Beach et al. (2013; 2015)

59   extended the RASE method for a more powerful regulatory information representation at both "the

60   block level (i.e., paragraph level) and inline (i.e., individual words or groups of words)", which

61   can be converted to Semantic Web Rule Language (SWRL) for reasoning. And, Dimyadi et al.

62   (2014) represented regulatory requirements using the Drools Rule Language (DRL).

63   These efforts have undoubtedly contributed to the improvement of flexibility and reusability of

64   regulatory representations for supporting ACC. However, they are still limited in terms of

65   automated regulatory information extraction and transformation; the state of the art in ACC still

3

66    requires major manual efforts in extracting regulatory information from textual regulatory

67    documents and transforming/encoding these information into a computer-processable rule format.

68    For example, in Pauwels et al. (2011), Hjelseth and Nisbet (2011), Yurchyshyna et al. (2010;

69    2008), Beach et al. (2013; 2015), and Dimyadi et al. (2014), the extraction of regulatory

70    information and their encoding into N3Logic, the RASE representation, the SPARQL queries, the

71    extended RASE representation, and the DRL rules, respectively, are still manually conducted. To

72    facilitate the regulatory information extraction and conversion, the SMARTcodes project led by

73    the International Code Council (ICC) developed tools to help ICC staff and building code officials

74    mark-up the ICC codes with provided tags under a predefined SMARTcodes schema. The marked

75    codes, then, can be automatically transformed into a "requirements model", which leverages the

76    IfcConstraint entities within an Industry Foundation Classes (IFC) model and therefore is

77    essentially an IFC constraint model (AEC3 2012). As the process suggests, the SMARTcodes

78    project still requires manual rule extraction and encoding efforts in the form of marking-up tasks.

79    To address these gaps of knowledge, this paper proposes a new fully-automated ACC system [the

80    authors call it semantic natural language processing (NLP)-based automated compliance checking

81    (SNACC) system] that integrates three types of algorithms in one unified computational platform:

82    (1) semantic NLP) algorithms to automatically extract the regulatory information from regulatory

83    documents (e.g., building codes) and transforms the extracted regulatory information into logic

84    rules, (2) semantic EXPRESS data processing algorithms to automatically extract the design

85    information from building information models and transform the extracted design information into

86    logic facts, and (3) semantic-based logic reasoning algorithms to automatically reason about the

87    compliance of the logic facts with the logic rules. The automated analyses are facilitated by

88    information representations that are semantic, logic-based, and generalized and flexible. This

4

89    paper presents the integration of the proposed algorithms in a unified ACC system and discusses

90    the experimental results of testing the proposed unified system using a test case.

91    **2    Proposed approach to full automation in automated compliance checking**

92    This paper proposes a fully-automated approach to ACC in construction. The approach relies on

93    the use of a set of computational techniques in an integrated manner, in one unified system. The

94    techniques include NLP, EXPRESS data processing, and logic reasoning, which are collectively

95    used for automated information processing (both design information and regulatory information)

96    and automated compliance reasoning. The automated processes are facilitated by semantic, logic-

97    based representations that are generalized and flexible.

98    2.1    Information representation

99    The choice of information representation has strong implications on information processing and is

100   of vital importance in facilitating automated processes. In ACC applications, specifically, there is

101   a need for a "standard, generalized approach for formally representing building regulations in a

102   digital format that would facilitate a variety of forms of reasoning about those codes in

103   combination with digital building information models" (Garrett et al. 2014), including automated

104   information extraction and information transformation to support complete automation of ACC.

105   The proposed representation is semantic and logic-based, in a way which is generalized and

106   flexible.

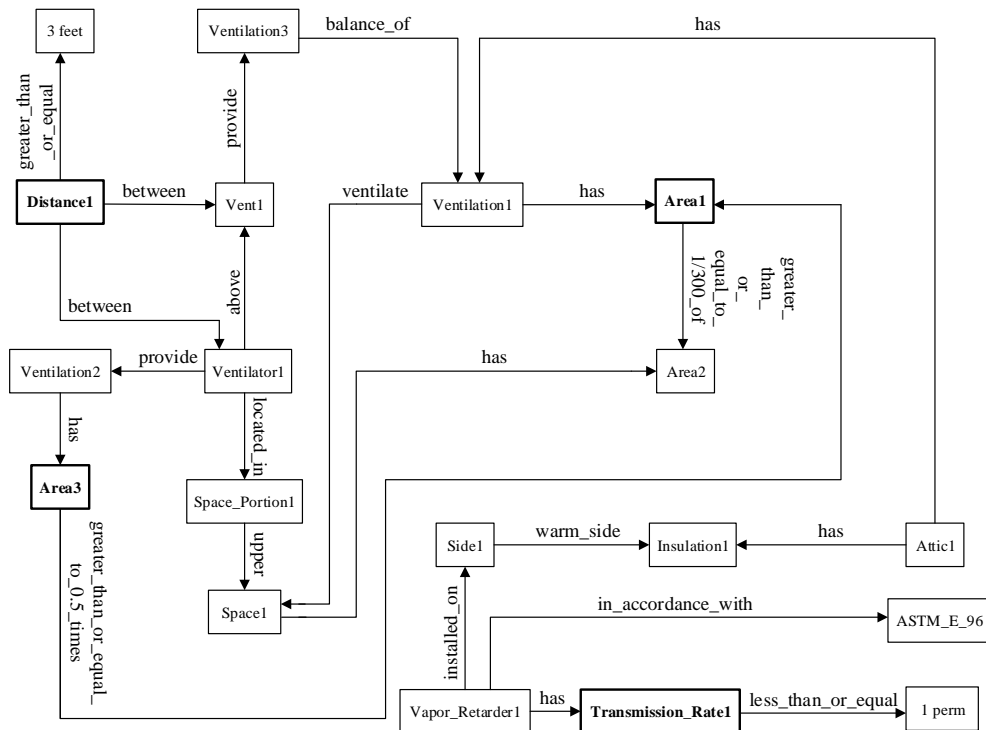107   *2.1.1   Semantic representation*

108   The representation is semantic; it uses semantic information elements and a domain ontology.

109   Semantic information elements represent the elements of a regulatory requirement, including

110   "subject," "compliance checking attribute," "deontic operator indicator," "quantitative relation,"

5

111    "comparative relation," "quantity value," "quantity unit," "quantity reference," "restriction," or

112    "exception." A building ontology is a semantic model for representing building domain knowledge

113    in the form of concept hierarchies, relationships between concepts, and axioms. The semantic

114    representation facilitates deep information processing (i.e., full-sentence analysis towards

115    capturing the entire meaning of a sentence, as opposed to shallow processing that extracts partial

116    information from a sentence). The semantic representation is also utilized to leverage domain

117    knowledge in the reasoning process, in order to handle the complex relations involved in

118    compliance reasoning and enable deep reasoning. This is important because the relations in

119    regulatory provisions could be very complex. For example, Fig. 1 shows the many relations

120    involved in one single regulatory provision in IBC 2006, leading to a very complex regulatory

121    provision. The semantic representation also facilitates human understandability and

122    interpretability of the formal representation, which is essential to facilitate usability and allow for

123    human testing and verification of the information representation and the reasoning results.

**Regulatory Provision from IBC 2006**

The minimum required net free ventilating area shall be 1/300 of the area of the space ventilated, provided a vapor retarder having a transmission rate not exceeding 1 perm in accordance with ASTM E 96 is installed on the warm side of the attic insulation and provided 50 percent of the required ventilating area provided by ventilators located in the upper portion of the space to be ventilated at least 3 feet above eave or cornice vents, with the balance of the required ventilation provided by eave or cornice vents.

**Semantic Representation of the Regulatory Provision**



**Fig. 1.** An example to illustrate the complexity of relations in provisions.

124
125

126  *2.1.2   Logic representation*

127   The representation is logic-based: regulatory information are represented as logic rules, while

128   design information are represented as logic facts. A logic-based representation was selected to take

129   advantage of the well-matured logic-based reasoning techniques. Logic-based reasoning is well-

130   suited for ACC problems because (Zhang and El-Gohary 2016b): (1) The binary nature (satisfy or

131   fail to satisfy) of logic fits the binary nature (compliance or noncompliance) of ACC; (2) Formally-

132   defined logics have sufficient expressiveness to represent concepts and relations involved in ACC;

133   (3) Once the information is properly represented in a logic format, the reasoning can be conducted

134   in a fully-automated way; and (4) Automated reasoning techniques are available in ready-to-use

135   logic reasoners.

7

136    Among the existing types of logic, FOL is the foundation of almost all work in rule representation

137    for ACC because of its expressivity; these efforts used a variety of logic

138    implementations/languages, but they all built on some restricted form of FOL. For example,

139    Pauwels and Zhang (2015) reviewed a good number of semantic rule checking applications among

140    which two main types of logic were used: N3Logic (e.g., Dimyadi et al. 2015) and SWRL (e.g.,

141    Baumgärtel et al. 2015). Both N3Logic and SWRL were created to go beyond the monotonic

142    negation limitation of FOL. N3Logic was created to avoid the paradox traps problem of FOL by

143    not using the general first-order negation but rather relying on customarily-made negated forms of

144    functions to achieve nonmonotonic negation (Berners-Lee 2005). SWRL is essentially combining

145    the Datalog Rule Markup Language (RuleML) with the Web Ontology Language (OWL), where

146    Datalog is a restricted subset of FOL using function-free Horn Clauses (HCs) (Horrocks et al.

147    2004). A HC is a restricted form of FOL that is most efficient in inference making (Saint-Dizier

148    1994). Both N3Logic and SWRL were used because of their compatibility with OWL ontologies,

149    which are the core of semantic rule checking approaches. More importantly, logic such as N3Logic

150    and SWRL need to be used to support if-then rule representation for rule checking when OWL

151    ontologies are used, because OWL is based on description logic (DL). A set of rules (if-then

152    statements) is necessary to allow for rule checking (Pauwels and Zhang 2015), and DL does not

153    allow for the representation of if-then rules. In addition to N3Logic and SWRL, Solihin and

154    Eastman (2015) took a knowledge representation approach for representing requirement rules

155    using conceptual graphs, which also has a semantic foundation in FOL and has one-to-one

156    mapping to FOL rules.

157    FOL was selected, in this paper, to support ACC not only because of its expressivity but also

158    because of its ability to represent English sentences. "A first-order sentence $\varphi$ can often be

159   translated into an English sentence which is guaranteed to be true if and only if φ is true in *I*" (i.e.,

160   the interpretation) (Hodges 2001). This property makes FOL suitable for representing regulatory

161   information to support automated compliance reasoning, because existing regulatory rules in

162   building codes and regulations are mostly coded in natural language sentences. Although FOL

163   cannot represent all provisions in building codes and regulations (Garrett et al. 2014), among those

164   provisions it can represent, FOL: (1) enables isomorphism: one-to-one mapping between an

165   English regulatory requirement and a logic clause, and (2) as a result, allows for traceability:

166   maintaining traceability is important to identify the sources of logic clauses and, thus, to facilitate

167   human verification and ensure trustworthiness of the logic clauses and the results. The scope of

168   this paper is limited to quantitative requirements – part of the regulatory requirements that are

169   representable in FOL. The representability of all possible types of regulatory requirements in FOL

170   (i.e., which requirements can be represented in FOL and which not) is an interesting topic that is

171   worth further investigation (Garrett et al. 2014), but is outside of the scope of this paper.

172   *2.1.3   Generalized and flexible representation*

173   The representation is generalized and flexible. The generalization and flexibility are achieved

174   through generalized regulatory compliance checking concepts and flexible semantic information

175   elements. Generalized regulatory compliance checking concepts (e.g., "subject" and "compliance

176   checking attribute") are used, which allows for representing regulatory provisions of any type/topic

177   (e.g., building envelope performance, facility accessibility). Flexible information elements (e.g.,

178   "subject restriction," as discussed in the following sections) are used, which allows for

179   representing all information (i.e., all concepts and relations) in a regulatory provision regardless

180   of the length and complexity of the provision (sentence). Generalization and flexibility are

9

181    important to sustain utility and robustness of the proposed system across different types of

182    regulatory documents and different types of provisions.

183    2.2    Computational techniques

184    *2.2.1    Deep natural language processing techniques*

185    It is an important impact conceived by many researchers who work on computable regulatory rule

186    representations (e.g., RASE, SMARTcodes) that the use of their representations may guide the

187    future drafting of codes and regulations (e.g., through the use of built-in annotations), so that the

188    automated extraction and transformation of regulatory information into computable rules would

189    be easily addressed. The authors also share that aspiration. But, at the same time, the authors

190    foresee that long-term goal (i.e., changing the way codes and regulations are drafted) as a big

191    challenge, potentially beyond the reach of solely the construction community, because it requires

192    harmonizing a lot of different pursuits and interests from various stakeholders (code drafters,

193    regulators, designers, etc.). Let alone that any developer of a computable regulatory rule

194    representation typically wants their own development to be adopted at a large scale, both

195    geographically and democratically – so which representation becomes a standard or becomes

196    widely adopted is another issue. On the other hand, the authors hold the ground of the status quo

197    that current building codes and regulations are mostly represented in natural language text, and

198    leverage state-of-the-art NLP techniques to develop new methods towards bridging the automation

199    gap of regulatory information extraction and transformation, under their proposed ACC framework.

200    NLP techniques are used to facilitate text analysis and processing for automatically extracting

201    regulatory information from building codes. NLP is a theoretically-based computerized approach

202    to analyzing, representing, and manipulating natural language text for the purpose of achieving

203    human-like language processing for a range of tasks or applications (Cherpas 1992). The types of

10

204    natural language analyses and techniques used highly affect the ability of NLP algorithms to

205    process complex sentences and recognize their full meaning. Full sentence understanding – of both

206    simple and complex sentences – is essential to achieve full automation in analyzing building codes

207    and extracting regulatory information. Deep NLP aims to capture the full meaning of sentences to

208    facilitate full sentence understanding by computers (Zouaq 2011). The proposed approach offers

209    a new way to achieve a deep level of text processing by integrating three types of knowledge in

210    the analysis of sentences: (1) ACC-specific knowledge: knowledge about the elements of a

211    regulatory requirement in building codes, represented in the form of semantic information

212    elements, (2) AEC domain knowledge: knowledge about the building domain, represented in the

213    form of an ontology, and (3) linguistic knowledge: knowledge about the linguistic expressions of

214    requirements in building code provisions, represented in the form of information extraction rules.

215    *2.2.2    EXPRESS data processing techniques*

216    EXPRESS data processing techniques are used for automatically extracting design information

217    from building information models. EXPRESS data processing techniques are suitable for

218    accessing information from IFC-based BIMs because the IFC schema is written in the EXPRESS

219    language. This EXPRESS language-level of processing enables the extraction and further

220    transformation of design information to be aligned with regulatory information.

221    The Java Standard Data Access Interface (JSDAI) was utilized for BIM information extraction,

222    using late binding data access methods. JSDAI is a standard data access interface (SDAI)

223    application programming interface (API) to access information from models written in EXPRESS

224    language – the ISO standard product data modeling language (ISO 2004). JSDAI provides two

225    types of data access methods: (1) early binding method, which accesses entities and attributes in

226    an EXPRESS model with specialized access methods such as "getCeilingHeight" (i.e., method to

227    get ceiling height of a floor), and (2) late binding method, which accesses entities and attributes in

228    an EXPRESS model with generalized access methods such as "getExplicitAttributes" (i.e., method

229    to get any explicit attribute). Compared to early binding, late binding allows accessing information
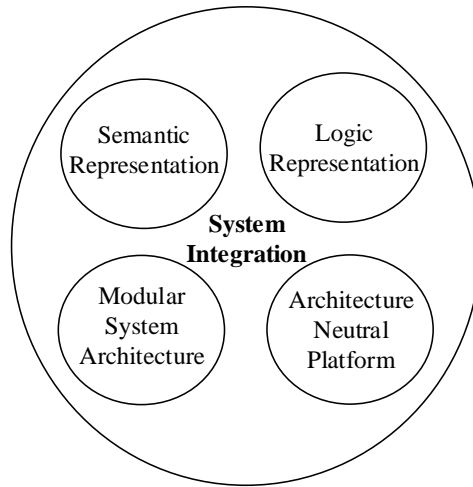
230    based on more general metadata.

231    *2.2.3   Logic reasoning and programming*

232    Logic programming is a computational programming paradigm that is based on a Horn Clause

233    (HC)-representation (Portoraro 2011). A program written in a logic programming language is

234    simply a set of logic sentences that represent facts and rules about some domain of interest. Logic

235    programming is declarative in contrast to other non-logical programming languages. For example,

236    in typical procedural programming languages like C programming language a programmer has to

237    clearly define how to solve the problem step by step, whereas in logic programming a programmer

238    only needs to define how to represent the problem in the form of facts and rules. The solution steps

239    in logic programming are already defined by a built-in reasoner through a set of organized

240    automated reasoning techniques such as search strategies and backtracking.

241    2.3   System integration

242    The proposed system offers a novel integration of natural language processing techniques,

243    EXPRESS data processing techniques, and logic reasoning into one unified computational

244    framework to allow for full automation in ACC. The integration is facilitated by the choice of (as

245    per Fig. 2): (1) a semantic representation that allows for seamless flow of information from one

246    computational paradigm to another, from one computational module to another, and from one

247    algorithm to another, (2) a logic representation, as a final representation, which allows to combine

248    partial output from two separate modules (logic rules and logic facts from module 1 and module

249    2, respectively) into one coherent representation that is ready for reasoning, (3) a modular system

250     architecture, which enables a flexible use of multiple modeling paradigms and multiple

251     programming languages, and (4) an architecture-neutral platform that can interoperate with

252     multiple programming language interfaces.



253
254                      **Fig. 2.** System integration.

255   **3    System architecture**

256     This section provides an overview of the system architecture, including an overview of: (1) the

257     system modules and how they are interlinked and integrated, and (2) the implementation of the

258     system modules and how they interact. More details on the individual system modules and

259     implementations are provided in Sections 4 and 5.

260     The system architecture is illustrated in Fig. 3. It is composed of three main modules: (1) regulatory

261     information extraction and transformation module, (2) design information extraction and

262     transformation module, and (3) compliance reasoning module. The system architecture is built on

263     top of the Java Platform (Oracle 1999). The General Architecture for Text Engineering (GATE)

264     tools (Cunningham et al. 2012), Python programming language (Python 2.7.3), B-Prolog logic

265     programming platform and reasoner (Zhou 2012), and JSDAI tools are used in the system to

266     support the computational processes in the different modules.

**Fig. 3.** System architecture of the SNACC system.

269    The regulatory information extraction and transformation module is composed of the regulatory

270    information extraction algorithm and the regulatory information transformation algorithm. The

271    information extraction algorithm aims to extract the regulatory requirements from a regulatory

272    document into a semantic information tuple representation, where each tuple contains information

273    instances for the semantic information elements (e.g., "subject," "compliance checking attribute").

274    The algorithm relies on the use of a set of pattern matching-based information extraction rules. A

275    set of syntactic and semantic features are used in the patterns of the information extraction rules.

276    The syntactic features are generated using GATE's Processing Resources (e.g., tokenizer), while

277    the semantic features are generated from the ontology using GATE's Processing Resources (e.g.,

278    gazetteer). The information extraction algorithm interacts with the Processing Resources using

14

279    GATE's API in Java. The regulatory information transformation algorithm aims to transform the

280    extracted instances of the semantic information elements in the information tuples into logic rules.

281    The algorithm relies on the use of a set of pattern matching-based semantic mapping rules and

282    conflict resolution rules, which include a set of syntactic and semantic features in their patterns.

283    The semantic features, here, are the semantic information element features (e.g., the semantic

284    feature "s" stands for "subject"). The information transformation algorithm interacts with the other

285    modules of the SNACC system (in Java) through Jython. The ontology is used to support the

286    regulatory information extraction and transformation processes by facilitating automated

287    interpretability and understandability of regulatory text based on meaning.

288    The design information extraction and transformation module is composed of the BIM information

289    extraction algorithm and the BIM information transformation algorithm. The BIM information

290    extraction algorithm aims to extract the entities and their attributes from a BIM into an information

291    tuple representation. The algorithm relies on the use of a set of entity and attribute extraction rules.

292    The data types of the entities and attributes are extracted from the BIM using late binding data

293    access methods in JSDAI. The BIM information transformation algorithm aims to transform the

294    extracted entities and attributes in the information tuples into logic facts that are aligned with the

295    logic rules. The algorithm relies on the use of initial transformation rules and semantic

296    transformation rules. The initial transformation rules transform the extracted entities and attributes

297    in the information tuples into logic facts. The semantic transformation rules further transform the

298    initially transformed logic facts into more semantic logic facts that are aligned with the predicates

299    in the logic rules. The initial transformation rules are coded in Java and the semantic transformation

300    rules are coded in B-Prolog rules. To execute the semantic transformation rules, the information

301    transformation algorithm interacts with B-Prolog's reasoner through B-Prolog's interface with

302    Java.

303    The compliance reasoning module is composed of the compliance reasoning algorithm, which

304    utilizes B-Prolog's reasoner. The compliance reasoning algorithm aims to reason about the logic

305    rules and the logic facts and generate compliance checking reports. The algorithm controls and

306    supports the reasoning about the rules and facts in B-Prolog's reasoner using a set of functional

307    built-in logic clauses. The compliance reasoning algorithm interacts with B-Prolog's reasoner

308    through B-Prolog's interface with Java. A user interacts with all the three modules through a

309    graphical user interface.

310    **4    System modules**

311    4.1    Regulatory information extraction and transformation module

312    The regulatory information extraction and transformation module is composed of four main

313    processes: preprocessing, feature generation, information extraction, and information

314    transformation.

315    Preprocessing prepares the raw natural language text of building codes for further processing. Four

316    NLP techniques are utilized: tokenization, sentence splitting, morphological analysis, and

317    dehyphenation. Tokenization divides the text into tokens (words or terms) to prepare for further

318    unit-based processing of the text. Sentence splitting recognizes the boundaries of the sentences to

319    help distinguish provisions in the building codes. Morphological analysis recognizes the different

320    forms of a word and maps them into the lexical form of that word. This helps in the recognition of

321    ontology concepts. Dehyphenation removes hyphens that indicate continuation of words between

322    lines to avoid further processing errors caused by those hyphens.

16

323   Feature generation generates a set of syntactic and semantic features that describe the text. Three

324   NLP techniques are utilized to generate the syntactic features: POS tagging, phrase structure

325   analysis, and gazetteer list analysis. POS tagging tags each word with the POS [lexical and

326   functional categories such as singular or mass noun (NN) and adjective (JJ)] of the word. Phrase

327   structure analysis tags each phrase with the phrasal tag [lexical and functional categories such as

328   noun phrase (NP) and verb phrase (VP)]. A set of application-specific phrase structure grammar

329   (PSG) rules are used to generate phrasal tags. The use of phrasal tags in addition to POS tags

330   reduces the potential number of enumerations in the patterns of the information extraction rules

331   (described in the following step). Gazetteer list analysis identifies each word that belongs to a

332   gazetteer list [a set of names based on any specific commonality possessed by those terms, e.g.,

333   "unit gazetteer list" includes inches and feet among others] and uses that information as a feature.

334   The building ontology is utilized to generate the semantic features, including terms/phrases that

335   match to the concepts and relations in the ontology. Fig. 4 shows a partial view of the ontology

336   that was used.

**Fig. 4.** Partial view of the building ontology.

337

338

339    Information extraction extracts the instances of the semantic information elements (SIEs) from the

340    building code, using a set of 146 information extraction (IE) rules. An SIE is an ontology concept,

341    an ontology relation, a "deontic operator indicator" (a term indicating an obligation, permission,

342    or prohibition), or a "restriction" (an element that places a constraint on the definition of another

343    semantic information element, where the constraint is expressed in terms of ontology concepts and

344    relations). The ten types of SIEs and their definitions are shown in Table 1. Each SIE is either a

345    "simple SIE" or a "complex SIE," and a "rigid SIE" or a "flexible SIE" (Zhang and El-Gohary

346    2013). A simple SIE is associated with a single concept/relation/indicator whereas a complex SIE

347    is expressed in terms of multiple concepts and relations. The simple SIEs are rigid [has a fixed

348    number (i.e., 1) of concepts/relations], whereas the complex SIEs are flexible [has a varying

349    number (i.e., 0 or more) of concepts/relations]. The IE rules use pattern matching; the rules extract

350    the instances of each SIE based on text patterns. The patterns consist of syntactic and semantic

351   features, which were generated during the feature generation step. For example, an IE rule for

352   extracting the instances of "subject" is shown in Fig. 5. The example IE rules use patterns that

353   consist of semantic features of "building element," "room," "space," and "quantity," and syntactic

354   features of "modal verb," "negation," "base form verb," "comparative relation," "cardinal number,"

355   "slash," and "unit." The IE rules were developed based on Chapters 12 and 23 of the International

356   Building Code (IBC) 2006 (ICC 2006). The extraction of each semantic information element is

357   separated and arranged in the following sequence because extracting all semantic information

358   elements from a sentence using a single IE rule is not efficient: "quantity value" and "quantity

359   unit/quantity reference" > "subject" > "compliance checking attribute" > "comparative relation" >

360   "quantitative relation" and "deontic operator indicator" > "subject restriction" and "quantity

361   restriction." An example illustrating the extraction is shown in Fig. 5. The text is then tagged with

362   the extracted SIEs for further information transformation.

363   Information transformation transforms the extracted information into logic rules, using a set of 9

364   conflict resolution (CR) rules, 297 semantic mapping (SM) rules, and a logic rule generator. The

365   CR rules and SM rules use pattern matching. The patterns consist of three types of information

366   tags: (1) syntactic information tags: syntactic feature tags generated during feature generation, (2)

367   semantic information tags: SIE tags generated during information extraction, and (3) combinatorial

368   information tags: compound information tags that are composed of multiple syntactic and/or

369   semantic information tags. Fig. 6 shows an example of a tagged regulatory requirement. The CR

370   rules resolve conflicts between the extracted information instances (in the form of four-element

371   tuples) based on the patterns. The SM rules transform the extracted information instances (after

372   conflict resolution) into logic components (i.e., logic predicates and logic operators) based on the

373   patterns. For example, 'n' 'c' 'v' 'u' is used as a pattern for an SM rule, which identifies a sequence

19

374  of "negation," "comparative relation," "quantity value," and "quantity unit." Fig. 7 shows an

375  example of an SM rule.

376  **Table 1.**

377  Semantic information elements (Zhang and El-Gohary 2016b; 2013).

| Semantic information element | Definition | Type |
|---|---|---|
| Subject | An ontology concept that describes a "thing" (e.g., building object, space) that is subject to a particular regulation or norm. | Simple and rigid SIE |
| Compliance checking attribute | An ontology concept that describes a specific characteristic of a "subject" by which its compliance is assessed. | Simple and rigid SIE |
| Deontic operator indicator | A term or phrase that indicates the deontic type of the requirement (i.e., whether it is an obligation, permission, or prohibition). | Simple and rigid SIE |
| Quantitative relation | A term or phrase that defines the type of relation for the quantity (e.g., "increase" is a quantitative relation). | Simple and rigid SIE |
| Comparative relation | An ontology relation that is commonly used for comparing quantitative values (i.e., comparing an existing value to a required minimum, maximum, or exact value), including "greater than or equal to," "greater than," "less than or equal to," "less than," and "equal to." | Simple and rigid SIE |
| Quantity value | A data value (or a range of values) that defines the quantified requirement. | Simple and rigid SIE |
| Quantity unit | The unit of measure for a "quantity value." | Simple and rigid SIE |
| Quantity reference | A term or phrase that refers to another quantity (which includes a value and a unit). | Simple and rigid SIE |
| Subject restriction | A term, phrase, or clause (which is composed of one or more concepts and/or relations) that places a constraint on the "subject." | Complex and flexible SIE |
| Quantity restriction | A term, phrase, or clause (which is composed of one or more concepts and/or relations) that places a constraint on the "quantity." | Complex and flexible SIE |

378

<u>Original Text:</u>

The thickness of concrete floor slabs supported directly on the ground shall not be less than 31/2 inches.

<u>Text with Features[1]:</u>

The thickness (ontology concept "quantity") of concrete floor slabs  (ontology concept "building element") supported directly on the ground shall (POS tag "MD" for modal verb) not (gazetteer list "Negation") be (POS tag "VB" for base form verb) less than (gazetteer list "Comparative relation") 31(POS tag "CD" for cardinal number)/(POS tag "Slash" for a slash)2(POS tag "CD") inches (gazetteer list "Unit").

<u>IE Rules:</u>

If "MD + Negation + VB + Comparative Relation" is matched, extract the text matched with "Negation" and the text matched with "Comparative relation" together as an instance for "comparative relation."

If ontology concept "building element" or "space" or "room" is matched, extract the matched text as an instance for "subject."

If ontology concept "quantity" is matched, extract the matched text as an instance for "compliance checking attribute."

If "CD + Slash + CD + Unit" is matched, extract the text matched with "CD + Slash + CD" as an instance of "quantity value," extract the text matched with "Unit" as an instance of "quantity unit."

<u>Extracted Instances:</u>

"thickness" as a "compliance checking attribute"

"concrete floor slab" as a "subject"

"not less than" as a "comparative relation"

"31/2" as a "quantity value"

"inches" as a "quantity unit"

1. For simplicity only features related to the IE rules below are displayed.

379
380             **Fig. 5.** Sample information extraction rules and extracted instances.

Original Text

The thickness of exterior basement walls and foundation walls shall be not less than 71/2 inches.

Information Tags

- Semantic information tags: 's' for subject, 'a' for compliance checking attribute, 'c' for comparative relation, 'v' for quantity value, 'u' for quantity unit;
- Syntactic information tags: 'CC' for conjunctive term, 'CD' for cardinal number, 'IN' for preposition, 'JJ' for adjective, 'MD' for modal verb, 'TO' for literal "to," 'VB' for base form verb, 'VBN' for past participle verb;
- Combinatorial information tags: 'dpvr' for directional passive Verbal relation, which is the combination of "past participle verb" (POS tag "VBN") and "preposition" (POS tag "IN").

Information Tuples Using Three Types of Information Tags[1]

[('thickness', 4, 9, 'a'), ('thickness', 4, 9, 'cr'), ('of', 14, 2, 'OF'), ('of', 14, 2, 'IN'), ('exterior basement walls', 17, 23, 's'), ('exterior', 17, 8, 'cr'), ('basement', 26, 8, 'cr'), ('walls', 35, 5, 'cr'), ('and', 41, 3, 'CC'), ('foundation walls', 45, 16, 's'), ('foundation', 45, 10, 'cr'), ('walls', 56, 5, 'cr'), ('shall', 62, 5, 'MD'), ('be', 68, 2, 'VB'), ('not', 71, 3, 'n'), ('less_than', 75, 9, 'c'), ('less', 75, 4, 'JJR'), ('than', 80, 4, 'IN'), ('71/2', 85, 4, 'v'), ('71/2', 85, 4, 'CD'), ('inches', 90, 6, 'u'), ('inches', 90, 6, 'cr')]

Information Tuples with Conflict Resolution Rules Applied[1]

[('thickness', 4, 9, 'a'), ('of', 14, 2, 'OF'), ('exterior basement walls', 17, 23, 's'), ('and', 41, 3, 'CC'), ('foundation walls', 45, 16, 's'), ('shall', 62, 5, 'MD'), ('be', 68, 2, 'VB'), ('not', 71, 3, 'n'), ('less_than', 75, 9, 'c'), ('71/2', 85, 4, 'v'), ('inches', 90, 6, 'u')]

Logic Components after Applying Semantic Mapping Rules[2]

thickness(Thickness),(exterior_basement_wall(Exterior_basement_wall);foundation_wall(Exterior_basement_wall)),has(Exterior_basement_wall,Thickness),not less_than(Thickness,quantity(71/2,inches))

Logic Rules Generated by Logic Rule Generator (Partial)[2]

*Primary Logic Clause*

compliance_thickness_of_Exterior_basement_wall81(Exterior_basement_wall):-
thickness(Thickness),(exterior_basement_wall(Exterior_basement_wall);foundation_wall(Exterior_basement_wall)),has(Exterior_basement_wall,Thickness),not less_than(Thickness,quantity(71/2,inches)).

*Activation Condition Logic Clause*

...thickness(Thickness),(exterior_basement_wall(X);foundation_wall(X)),has(X,Thickness)->
check_thickness_of_Exterior_basement_wall81(X);true,...

*Compliance Checking Consequence Logic Clause*

check_thickness_of_Exterior_basement_wall81(X):-(compliance_thickness_of_Exterior_basement_wall81(X)->
writeln((X,is,compliant,with,section,1909-6-1,rule81));writeln((X,is,noncompliant,with,section,1909-6-1,thickness,should,be,not,less_than,71/2,inches,rule82))).

1. Each tuple includes four elements: the information instance, its location (the starting point in the sentence), its length (in number of letters), and its information tag.
2. In this logic syntax, comma represents conjunction, semicolon represents disjunction, "not" represents negation, ":-" represents implication, predicate takes the form of pred(arg1,arg2,...), rule takes the form of predh(arg1,arg2,...) :- pred1(arg1,arg2,...), pred2(arg1,arg2,...)..., predn(arg1,arg2,...).

381
382 **Fig. 6.** An Example to illustrate regulatory information transformation.

SM Rule Pattern

Logic Clause

'subject' (s) 'negation' 'comparative relation' (cr) 'value' (v) 'unit' (u)

s(S),not cr(S,quantity(v,u))

*Note: An upper case represents a variable.*

383
384 **Fig.7.** An example of a semantic mapping rule.

22

385 The logic rule generator generates three types of logic rules based on the logic components:

386 primary logic clauses, activation condition logic clauses, and compliance checking consequence

387 logic clauses. A primary logic clause is the main representation of a requirement; the premise of

388 the rule represents the conditions of a requirement and the conclusion of the rule represents the

389 consequent result (i.e., the compliance with the requirement). For example, in the primary logic

390 clause in Fig. 6, the logic components to the right of ":-" represent the conditions of the wall

391 thickness requirement (for exterior basement walls and foundation walls) and the logic

392 components to the left of ":-" represent the conclusion of that requirement. An activation condition

393 logic clause represents the conditions that activate the checking of a requirement, which are the

394 existence of the corresponding information in the BIM (e.g., the existence of exterior basement

395 wall or foundation wall and thickness information for the example in Fig. 6). Activation conditions

396 are used to help prevent missing information from leading to false positives because missing

397 information would lead to failure in activation. A compliance checking consequence logic clause

398 represents the consequences of the compliance checking result (compliance or noncompliance).

399 For example, if the result is noncompliant, a corrective suggestion is provided (e.g., "thickness

400 should be not less than 71/2 inches," as per Fig. 6).


401 4.2 Design information extraction and transformation module

402 The design information extraction and transformation module is composed of two main processes:

403 BIM information extraction and BIM information transformation.


404 BIM information extraction utilizes EXPRESS data processing techniques in a BIM information

405 extraction algorithm to extract all entities and their attributes in an IFC file into information tuples

406 based on their metadata, in a recursive and exhaustive manner. The information tuples store

407 information for each entity, the attributes of the entity, and the values of the attributes of the entity,

408    to prepare for the following transformation process. The BIM IE algorithm exhaustively extracts

409    the values (e.g., '2O2Fr$t4X7Zf8NOew3FNld') for each attribute (e.g., global ID) of each entity

410    (e.g., wall). Recursion is used in two ways (as illustrated in Fig. 8): (1) when an entity is being

411    extracted, not only the explicit attributes of the entity are extracted, but all explicit attributes that

412    belong to the supertype of that entity and supertype of supertype (until no supertype can be found)

413    of that entity are extracted too. For example, when a "door" entity is being extracted, not only the

414    explicit attributes "overall height" and "overall width" are extracted, but all the following explicit

415    attributes that belong to the supertypes of "door" are extracted too: "global ID," "owner history,"

416    "name," "description," "object type," "object placement," "representation," and "tag;" and (2) if

417    an attribute is of an aggregation data type (i.e., aggregation of multiple attributes), then the member

418    attributes of the aggregation are recursively accessed for extracting their values. For example,

419    because the attribute "related objects" of a "rel associates material" is of an aggregation data type

420    (i.e., set data type in this case), when a "related objects" instance is being processed, each of its

421    member objects is accessed recursively for extracting their values. The late binding data access

422    method in JSDAI is used to support the entity and attribute extraction in the BIM IE algorithm.

423    Late binding accesses each entity and attribute using standard access methods in Java.

```
                         ┌─────────┐
                         │  Start  │
                         └─────────┘
                              │
                              ▼
                   ╱─────────────────────╲
                  ╱   Any unprocessed      ╲ ◄──────────────────┐
                  ╲   entity in the IFC file? ╱                 │
                   ╲─────────────────────╱                      │
                              │ Yes                             │
                              ▼                                 │
              ┌──────────────────────────────┐   Subroutine S1  │
              │ Read an entity in an IFC file │                  │
              │ into E;                       │                  │
              └──────────────────────────────┘                  │
                              │                                  │
                              ▼                                  │
                   ╱─────────────────────╲                      │
                  ╱   Is the entity an     ╲──Yes──┐             │
                  ╲   aggregation type?    ╱       │             │
                   ╲─────────────────────╱         ▼             │
                              │ No       ┌───────────────────┐   │
                              ▼          │ Iterate through    │  │
              ┌──────────────────────┐   │ each sub-entity SE │  │
              │ Get explicit         │   │ in the aggregate,  │  │
              │ attributes (names and│   │ process each SE    │  │
              │ values) of E into    │   │ using the          │  │
              │ attribute set AS;    │   │ subroutine S1;     │  │
              └──────────────────────┘   └───────────────────┘   │
                              │                    │             │
                              ▼                    │             │
          No      ╱─────────────────────╲          │             │
        ◄─────────╲   Entity E has a     ◄──────────┘             │
                  ╱   supertype?         ╱                        │
                   ╲─────────────────────╱                       │
                              │ Yes                              │
                              ▼                                  │
              ┌──────────────────────────────┐                   │
              │ Get explicit attributes      │                   │
              │ (names and values) of        │───────────────────┘
              │ supertype of E into AS;      │
              │ Assign the supertype of E    │
              │ to E;                        │
              └──────────────────────────────┘
                              │ No
                              ▼
                         ┌─────────┐
                         │   End   │
                         └─────────┘
```

**Fig. 8.** The BIM IE Algorithm based on two recursive processes.

BIM information transformation transforms the extracted BIM information in the information tuples into logic facts (concept facts and relation facts) in two steps: initial transformation and alignment transformation. Initial transformation transforms the extracted entities and their attributes into concept facts and relation facts using three main initial transformation rules. These rules transform elements in the entities, attributes, and values into predicate names or arguments based on their metadata. For example, the first initial transformation rule in Fig. 9 converts a line

25

432    in IFC data with referenced attribute values into logic facts. After initial transformation, alignment

433    transformation further transforms the generated logic facts into a logic fact representation that is

434    aligned with the predicates in the logic rules (that represent the corresponding regulatory

435    requirements). A set of semantic transformation (ST) rules are used in the alignment

436    transformation step. For example, Fig. 9 shows a set of logic facts after initial transformation and

437    after alignment transformation using two ST rules. Compared to the logic facts before alignment

438    transformation, the logic facts after alignment transformation are more easily understandable and

439    aligned with the logic rules.

---

IFC Data

#39592=IFCHEIGHT();
#39594=IFCWALL($,$,$,$,$,$,$,$);
#39595=IFCACCBIRELATION($,$,$,$,'has',#39594,#39592,$);

Initial Transformation Rules

(1) an entity is transformed into a concept fact (i.e., a predicate) by using
the name of the entity as the name of the predicate, and using the name
of the entity concatenated with the ID of the entity as the argument (i.e.,
an entity constant) of the predicate.
(2) an attribute of an entity is transformed into a relation fact (i.e., a
predicate), using the name of the attribute preceded by "has_" as the
name of the predicate, using the corresponding entity constant as the first
argument of the predicate, and using the value of the attribute as the
second argument of the predicate (if the value is not a reference to
another entity).
(3) if the value of an attribute is a reference to another entity, then the
referred entity constant is used as the second argument of the predicate.

Logic Facts After Initial Transformation

acc_bi_relation(acc_bi_relation39595).
has_type_name(acc_bi_relation39595,has).
has_relating_element(acc_bi_relation39595,wall39594).
has_related_element(acc_bi_relation39595,height39592).

Alignment Transformation Rules

acc_bi_relation(X), has_type_name(X,Name),
has_relating_element(X,Y),has_related_element(X,Z) → relation(Name, Y, Z).
relation(Name, Y, Z) → Name(Y,Z).

Logic Facts After Alignment Transformation

has(wall39594,height39592).

---

440
441         **Fig. 9.** An example to illustrate BIM information extraction and transformation.

442    4.3    Compliance reasoning module

443    The compliance reasoning module utilizes B-Prolog's reasoner to reason about the logic rules and

444    the logic facts and generate compliance checking reports. A set of functional built-in logic clauses

445    were developed and embedded into the system to provide basic arithmetic functions (e.g., unit

446    conversion) and define the sequence of execution/checking. For execution, the user specifies the

447    list of subjects (e.g., walls and doors) or subjects and attributes (e.g., walls and their heights) to

448    check, and accordingly the subjects in the specified list are sequentially checked one by one. By

449    default, a "select all" option is used.


450    **5    System implementation**

451    The proposed SNACC system was implemented in a proof-of-concept prototype. The main

452    platform of the prototype was built using Java programming language (Java Standard Edition

453    Development Kit 6u45). The regulatory information extraction algorithm was implemented using

454    GATE's Processing Resources and Java programs. The following Processing Resources were used:

455    (1) the English Tokenizer, Sentence Splitter, POS Tagger, and Gazetteer in the A Nearly-New

456    Information Extraction (ANNIE) system for tokenization, sentence splitting, POS tagging, and

457    gazetteer compiling, (2) the Morphological Analyzer for morphological analysis, (3) the Flexible

458    Gazetteer for generating semantic features based on the ontology, and (4) the Java Annotation

459    Patterns Engine (JAPE) rules for encoding the IE rules. The information extraction algorithm

460    interacts with the Processing Resources using GATE's API 7.0.

461    The regulatory information transformation algorithm was implemented using Python

462    programming language (Python 2.7.3). The SM rules and CR rules were coded as Python

463    conditional statements. The "re" module (i.e., regular expression module) in Python was used for

464    both extracting the syntactic and semantic features from the information tuples and conducting
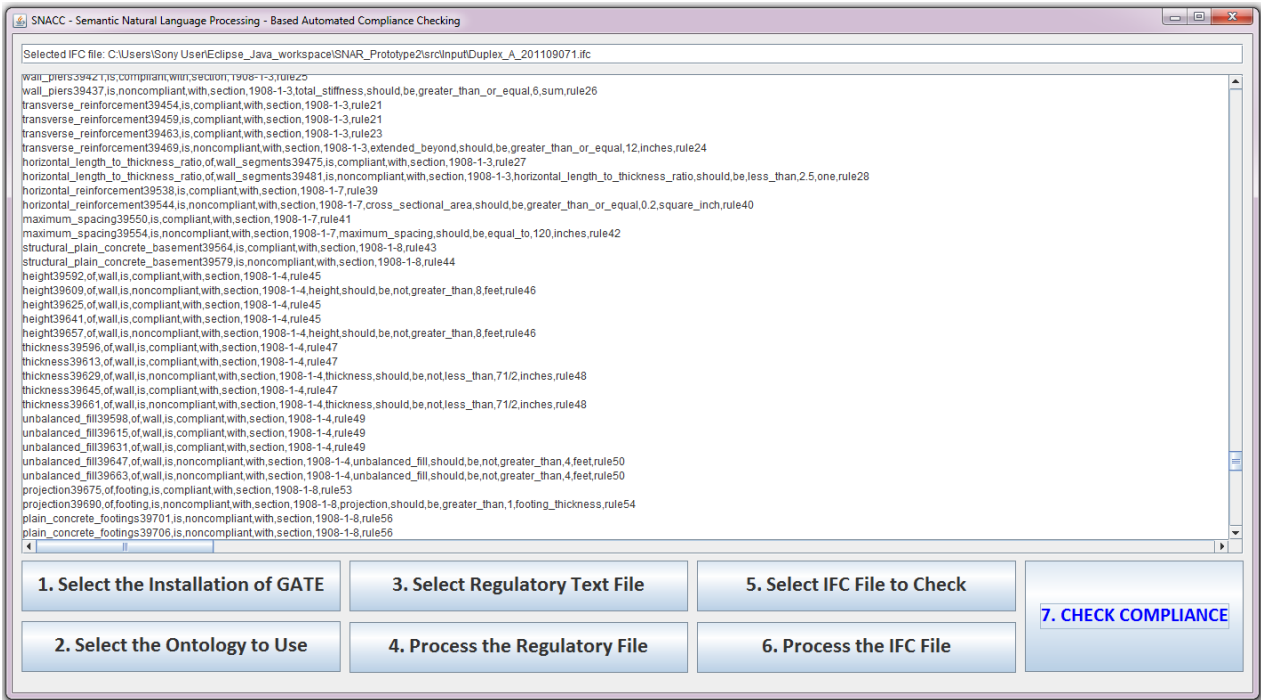
465    pattern matching. The information transformation algorithm interacts with the other modules of

466    the SNACC system (in Java) through Jython 2.2.1.

467    The BIM information extraction and transformation algorithms were implemented in Java

468    programs and B-Prolog rules, respectively. The JSDAI runtime (JSDAI 4.3.0) was used to access

469    the information in IFC-based BIMs (IFC files) for entity and attribute extraction. String processing

470    methods in Java were used for initial transformation. Static rules and dynamic rules in B-Prolog

471    were used for alignment transformation. Static rules are rules that only use static predicates.

472    Dynamic rules are rules that use at least one dynamic predicate. A static predicate is a predicate

473    that cannot be updated during execution whereas a dynamic predicate is a predicate that can be

474    updated during execution. The rules for entity extraction, attribute extraction, and initial

475    transformation were coded as Java conditional statements. The rules for alignment transformation

476    (i.e., ST rules) were coded as B-Prolog rules.

477    The logic-based automated reasoning algorithm was implemented in Java. The functional built-in

478    logic clauses were encoded in B-Prolog. The automated reasoning algorithm interacts with the

479    logic clauses and logic reasoner through B-Prolog's bi-directional interface 7.8 with Java

480    programming language.

481    The graphical user interface of the SNACC system is shown in Fig. 10. As shown in Fig. 10, the

482    SNACC system requires the download of the GATE tool and the availability of a building ontology

483    to execute the regulatory information extraction and transformation algorithms. A user could then

484    select the regulatory document (.txt file) and the BIM (.ifc file) for automated compliance checking.

485    The information extraction and information transformation algorithms for regulatory information

486    and design information could be executed in parallel. After all information have been extracted

487    and transformed, pressing the "check compliance" button activates the automated reasoning

488    process using B-Prolog. The compliance checking results are then automatically displayed to the

489    user in the text field of the graphical user interface (as shown in Fig. 10).



490
491                    **Fig. 10.** Graphical user interface of the SNACC system.

## 6    System testing

493    The SNACC system was tested in checking the compliance of a BIM test case with Chapter 19 of

494    IBC 2009. IBC was selected because it is predominantly adopted in the United States. Chapter 19

495    was then randomly selected. For the test case, it was developed based on the Duplex Apartment

496    Project from buildingSMARTalliance of the National Institute of Building Sciences (East 2013).

497    Design information were added in the BIM model, based on an extended version of the

498    IFC_2X3_TC1 schema (BuildingSmart 2014) (Zhang and El-Gohary 2016a). The test case

499    included design information for each provision in Chapter 19 of IBC 2009. The design information

500    included both compliant and noncompliant design information. If a provision had more than one

501    requirement, then compliant and noncompliant design information for each requirement was

502    included. For example, the following regulatory provision (*RP1)* is a complex provision that

29

503    contains three quantitative requirements: "*In dwellings assigned to Seismic Design Category D or*

504    *E, the height of the wall shall not exceed 8 feet (2438 mm), the thickness shall not be less than 71/2*

505    *inches (190 mm), and the wall shall retain no more than 4 feet (1219 mm) of unbalanced fill.*"

506    Thus, five information sets were created for *RP1* which correspond to the scenarios that (1) only

507    height is noncompliant, (2) only thickness is noncompliant, (3) only unbalanced fill is

508    noncompliant, (4) all three attributes are noncompliant, and (5) no attributes are noncompliant.

509    *7    Results and discussion*

510    The ACC prototype system was evaluated using precision, recall, and F1-measure of

511    noncompliance detection. Precision is defined as the number of correctly-detected noncompliance

512    instances divided by the total number of noncompliance instances detected. Recall is defined as

513    the number of correctly-detected noncompliance instances divided by the total number of

514    noncompliance instances that should be detected. F1-measure is the harmonic mean of precision

515    and recall. A manually-developed gold standard was used for the evaluation. A gold standard refers

516    to a benchmark against which testing results are compared for evaluation. The gold standard

517    includes the ground truth of compliant and noncompliant instances.

518    The testing results are summarized in Table 2. As shown in Table 2, the recall, precision, and F1-

519    measure of noncompliance detection is 98.7%, 87.6%, and 92.8%, respectively. The relevant

520    provision numbers and rule numbers for the compliant and noncompliant instances were also

521    correctly reported. For each noncompliance instance, a suggestion on how to fix the

522    noncompliance case was also correctly reported (partially shown in Fig. 10).

523    **Table 2.**
524    Noncompliance detection testing results.

| Parameter/measure | Result |
|---|---|
| Number of noncompliance instances in gold standard | 79 |

| | |
|---|---|
| Number of noncompliance instances detected | 89 |
| Number of noncompliance instances correctly detected | 78 |
| Recall of noncompliance detection | 98.7% |
| Precision of noncompliance detection | 87.6% |
| F1-measure of noncompliance detection | 92.8% |

525   These high performance results show that the proposed ACC system is promising. In addition, the

526   fact that the proposed ACC system achieved higher recall (98.7%) than precision (87.6%) shows

527   its suitability for the ACC application; in noncompliance detection, recall is more important than

528   precision. Recall errors are more critical because they might result in missing noncompliance

529   instances, whereas precision errors could be easily double-checked and filtered out by the user.

530   An error analysis was also conducted to identify the sources of the errors in noncompliance

531   detection. The noncompliance detection errors originated from errors in regulatory information

532   extraction and regulatory information transformation; there were no errors in BIM information

533   extraction, BIM information transformation, or compliance reasoning. The errors were attributed

534   to errors made by GATE's processing resources, limitations of rules used in regulatory information

535   extraction and information transformation, and limitations of the state-of-the-art NLP techniques

536   [e.g., state-of-the-art Part-of-Speech (POS) tagging has an accuracy of around 97% (Manning

537   2011)]. For example, "concrete floor slab" was not successfully extracted as the subject (i.e., a

538   false negative) for the following requirement because of errors made by GATE's processing

539   resources: "The thickness of concrete floor slabs supported directly on the ground shall not be less

540   than 31/2 inches (89 mm)" (Provision 1910.1 of IBC 2009).

541   **8   Contribution to the body of knowledge**

542   This research contributes to the body of knowledge in three main ways. First, this research offers

543   a novel system for fully-automated checking of building information models for compliance with

544   building codes. The proposed system goes beyond the current state-of-the-art of ACC by allowing

545    fully-automated (1) extraction of both regulatory and design information from regulatory

546    documents and IFC-based BIM models, respectively, and (2) alignment of the representations of

547    these two sets of information, so that they can be interpreted together in one system. Second, this

548    research offers integrated NLP and first order logic methods for automatically extracting

549    regulatory information from regulatory documents and automatically representing the extracted

550    information in an ACC first order logic-based representation that is used in automated ACC logic

551    reasoning. The proposed methods/algorithms offer a novel way for, both, deep information

552    extraction (i.e., full-sentence analysis to capture the entire meaning of a provision) and generalized

553    and flexible ACC representation; both – together – enable the extraction and representation of

554    information even in long and complex provisions, which is important to sustain utility and

555    robustness of ACC system performance across different types of regulatory documents and

556    different types of provisions. Third, this research offers a novel combination of NLP techniques

557    with both semantic analysis and logic-based reasoning into one computational framework. In this

558    research, a set of information extraction, information transformation, and automated reasoning

559    algorithms are effectively implemented into one proof-of-concept ACC system. The combined

560    performance of all algorithms, into the system, shows high automated noncompliance detection

561    performance (98.7%, 87.6%, and 92.8% recall, precision, and F1-measure, respectively).

562    **9    Conclusions**

563    This paper presented a unified system that integrates a set of techniques and algorithms for

564    automatically checking the compliance of BIM-based building designs with building codes. The

565    proposed system offers a fully-automated approach to ACC in construction. The approach relies

566    on the use of a set of computational techniques in an integrated manner, in one unified system.

567    The techniques include NLP, EXPRESS data processing, and logic reasoning, which are

32

568    collectively used for automated information extraction, automated information transformation, and

569    automated compliance reasoning. The automation is facilitated by semantic, logic-based

570    representations that are generalized and flexible.

571    The system is composed of three main modules: (1) a regulatory information extraction and

572    transformation module, which utilizes semantic natural language processing algorithms to

573    automatically extract regulatory information from building codes and transform the extracted

574    information into logic rules, (2) design information extraction and transformation module, which

575    utilizes EXPRESS data processing-based  algorithms to automatically extract design information

576    from building information models and transform the extracted information into logic facts, and (3)

577    compliance reasoning module, which utilizes semantic-based logic reasoning algorithms to

578    automatically reason about the compliance of the logic facts with the logic rules. The algorithms

579    were implemented in different programming languages and integrated into one proof-of-concept

580    prototype system (the SNACC system). The integration is facilitated by the choice of a semantic

581    representation, a logic representation, a modular system architecture, and an architecture-neutral

582    platform.

583    The SNACC system was tested in checking the compliance of a BIM test case with Chapter 19 of

584    IBC 2009. A recall of 98.7%, a precision of 87.6%, and an F1-measure of 92.8% in noncompliance

585    detection were achieved. The high performance results, of all algorithms when combined into one

586    unified system, show that the proposed ACC system is promising. In addition, the higher recall

587    shows the suitability of the proposed system for ACC, because recall is more critical than precision

588    for noncompliance detection.

## 10  Limitations and future work

589

590    As mentioned above, at this point, the system proposed in this paper focused on quantitative

591    requirements. It could be extended to support the checking of other types of requirements such as

592    existential requirements (i.e., rules that require the existence of certain building elements, etc.),

593    but it cannot go beyond the limitations of machine intelligence or represent and reason with rules

594    that require human judgement by nature.

595    Also, in spite of the authors' firm belief in automation and early evidence of low consistency in

596    manual noncompliance checking (Fiatech 2014), how the automated information extraction and

597    transformation approach proposed in this paper compares to the state-of-the-art semi-automated

598    information extraction and transformation approaches (e.g., such as RASE-based or

599    SMARTcodes-based, which rely on manual annotation) in terms of accuracy and efficiency

600    requires further investigation.

601    As part of their future/ongoing research work, the authors will test the proposed ACC system on

602    more building code chapters and more BIM test cases. In addition, other types of requirements

603    (e.g., existential requirements) will be tested, and different ways of handling information

604    incompleteness cases during ACC will be proposed and tested.

605    In future research – by the authors or the larger research community, the proposed information

606    extraction and transformation algorithms could also be applied to other logic-based representations

607    such as SWRL and N3Logic. In this case the JSDAI-based BIM information processing can be

608    partially replaced by existing conversion methods such as those in Pauwels and Terkaj (2016) and

609    Beetz et al. (2009). However, in this case, further semantic transformation of BIM information

610    would still be needed to align the concept representations of the design information to those of the

34

611   regulatory information. Similarly, further research could be conducted to study how to best link

612   the proposed algorithms with OWL representations and other semantic modeling approaches and

613   assess the advantages and limitations of the proposed methods in this context. The authors expect

614   that the proposed information extraction, information transformation, and automated reasoning

615   methods would lend themselves well to such integrative efforts. However, further research is

616   needed to study practicality, benefits, and limitations.

617   **Acknowledgements**

622   **References**

623   AEC3. (2012). "International Code Council." ⟨http://www.aec3.com/en/5/5_013_ICC.htm⟩ (Jun.

624       30, 2016).

625   Baumgärtel, K., Kadolsky, M. and Scherer, R.J. (2015). "An ontology framework for improving

626       building energy performance by utilizing energy saving regulations." Proc., 10th European

627       Conference on Product and Process Modelling (ECPPM), ECPPM, 519-526.

628   Beach, T.H., Kasim, T., Li, H., Nisbet, N., and Rezgui, Y. (2013). "Towards automated

629       compliance checking in the construction industry." *H. Decker et al. (Eds.): DEXA 2013, Part*

630       *I, LNCS 8055*, 366-380.

631   Beach, T.H., Rezgui, Y., Li, H., and Kasim, T. (2015). "A rule-based semantic approach for

632         automated regulatory compliance in the construction sector." *Expert Systems with*

633         *Applications*, 42(12), 5219-5231.

634   Beetz, J., van Leeuwen, J., and de Vries, B. (2009). "IfcOWL: A case of transforming EXPRESS

635         schemas into ontologies." *Artificial Intelligence for Engineering Design, Analysis and*

636         *Manufacturing*, 23(SP01), 89-101.

637   Berners-Lee, T. (2005). "Status: An early draft of a semi-formal semantics of the N3 logical

638         properties." <https://www.w3.org/DesignIssues/N3Logic> (Jun. 9, 2016).

639   Boken, P., and Callaghan, G. (2009). "Confronting the challenges of manual journal entries."

640         *Protiviti*, Alexandria, VA, 1-4.

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

642         <http://www.buildingsmart-tech.org/specifications/ifc-overview> (Jan 19, 2015).

643   Cherpas, C. (1992). "Natural language processing, pragmatics, and verbal behavior." *Anal. Verbal*

644         *Behav.*, 10, 135–147.

645   Choi, J., Choi, J., and Kim, I. (2014). "Development of BIM-based evacuation regulation checking

646         system for high-rise and complex buildings." *Autom. Constr.*, 46, 38-49.

647   City   of   Philadelphia.   (2015).   "Licenses   and   inspections:   building   permits."

648         <https://business.phila.gov/licenses-and-inspections-building-permits/> (Sept. 4, 2015).

649   Cunningham, H., et al. (2012). "Developing language processing components with gate version 7

650         (a user guide)." Univ. of Sheffield, Dept. of Computer Science, Sheffield, U.K.

651   Delis, E.A., and Delis, A. (1995) "Automatic fire-code checking using expert-system technology."

652         *J. Comput. Civ. Eng.*, 9(2), 141-156.

653     Department of Buildings. (2015). "Hearing on the fiscal 2016 preliminary budget and the fiscal

654           2015          preliminary          mayor's          management          report."

655           <ttp://council.nyc.gov/html/budget/2016/Pre/dob.pdf> (Sept. 4, 2015).

656     Dimyadi, J., and Amor, R. (2013). "Automated building code compliance checking - where is it

657           at?" *Proc. 19th Int. CIB World Build. Congress*, Brisbane, Australia.

658     Dimyadi, J., Clifton, C., Spearpoint, M., and Amor, R. (2014). "Regulatory knowledge encoding

659           guidelinens for automated compliance audit of building engineering design." *Comput. Civ.*

660           *Build. Eng.* (2014), ASCE, Reston, VA, 536-543.

661     Dimyadi, J., Pauwels, P., Spearpoint., M., Clifton, C., and Amor, R.W. (2015). "Querying a

662           regulatory model for compliant building design audit." *Proc., CIB W78 2015*, Conseil

663           International du Bâtiment (CIB), Rotterdam, The Netherlands, 139-148.

664     East, E.W. (2013). "Common building information model files and tools." <

665           http://www.nibs.org/?page=bsa_commonbimfiles&hhSearchTerms=%22common+and+BI

666           M+and+file%22> (Jun. 27, 2014).

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

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

669     Fiatech Regulatory Streamlining Committee, (2012). "AutoCodes project: phase 1, proof-of-

670           concept:                    final                    report."

671           <http://www.fiatech.org/images/stories/techprojects/project_deliverables/Updated_project_d

672           eliverables/AutoCodesPOCFINALREPORT.pdf> (Dec. 24, 2013).

673     Fiatech. (2014). "Automated code plan checking tool-proof-of-concept (phase 2)."

674           <http://www.fiatech.org/images/stories/projects/FiatechAutoCodesPh2-Report-Sept2015.pdf>

675           (Jun. 16, 2016).

676    Garrett, J.H.Jr., and Palmer, M.E. (2014). "Delivering the infrastructure for digital building

677         regulations." *J. Comput. Civ. Eng.*, 2014(28), 167-169.

678    Hjelseth, E. and Nisbet, N. (2011). "Capturing normative constraints by use of the semantic mark-

679         up RASE methodology." *Proc., CIB W78 2011*, Conseil International du Bâtiment (CIB),

680         Rotterdam, The Netherlands.

681    Hodges, W. (2001). "Classical logic I - first-order logic." *Goble,* 2001, 9-32.

682         <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.137.4783&rep=rep1&type=pdf>

683         (Dec. 26, 2013).

684    Horrocks, I., Patel-Schneider, P.F., Boley, H., Tabet, S., Grosof, B., and Dean, M. (2004). "SWRL:

685         A   Semantic   Web   Rule   Language   Combining   OWL   and   RuleML."   <

686         https://www.w3.org/Submission/SWRL/> (Jun. 9, 2016).

687    International Code Council (ICC). (2006). "2006 international building code." 2006 Int. Codes, 〈

688         http://publiccodes.cyberregs.com/icod/ibc/2006f2/〉 (Oct. 25, 2015).

689    ISO. (2004). "ISO 10303-11:2004 - Part 11: Description methods: The EXPRESS language

690         reference                                                                            manual."

691         <http://www.iso.org/iso/iso_catalogue/catalogue_tc/catalogue_detail.htm?csnumber=38047>

692         (Dec. 05, 2014).

693    Java Standard Edition Development Kit 6u45. [Computer Software]. Redwood Shores, CA, Oracle.

694    JSDAI 4.3.0. [Computer software]. Kuenzell, Germany, LKSoftWare GmbH.

695    Jython 2.2.1. [Computer software]. <http://www.jython.org/> (May 02, 2015).

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

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

698         University, Tees Valley, UK, 11-21.

699    Lau, G. T., and Law, K. (2004). "An information infrastructure for comparing accessibility

700         regulations and related information from multiple sources." *Proc., 10th Int. Conf. on*

701         *Computational Civil and Building Engineering (ICCCBE)*, ISCCBE, Hong Kong, China.

702    Los        Angeles        Times.        (2015).        "Public        &        Legal        notices."

703         <http://classifieds.latimes.com/classifieds?category=public_notice> (Sept. 4, 2015).

704    Manning, C.D. (2011). "Part-of-Speech tagging from 97% to 100%: is it time for some linguistics?"

705         *Proc., 12th International Conference on Intelligent Text Processing and Computational*

706         *Linguistics*, CICLing, Mexico.

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

708         *Proc., 2011 Winter Simulation Conference*, Association for Computing Machinery, New York,

709         3400 - 3405.

710    Oracle.        (1999).        "Essentials        of        the        Java        programming        language,        Part        1."

711         <http://www.oracle.com/technetwork/java/index-138747.html> (Jun. 25, 2015).

712    Pauwels, P., and Terkaj, W. (2016). "EXPRESS to OWL for construction industry: Towards a

713         recommendable and usable ifcOWL ontology." *Autom. Constr.*, 63(Mar. 2016), 100-133.

714    Pauwels, P., and Zhang, S. (2015). "Semantic rule-checking for regulation compliance checking:

715         an overview of strategies and approaches." *Proc., CIB W78 2015,* Conseil International du

716         Bâtiment (CIB), Rotterdam, The Netherlands, 619-628.

717    Pauwels, P., Van Deursenc, D., Verstraetena, R., De Rooc, J., De Meyera, R., Van de Wallec, R.,

718         Van Campenhoutb, J. (2011). "A semantic rule checking environment for building

719         performance checking." *Autom. Constr.*, 20(5), 506-518.

720    Portoraro, F. (2011). "Automated reasoning." *The Stanford encyclopedia of philosophy (Summer*

721         *2011        Edition)*,        Edward        N.        Zalta        (ed.)

722        <http://plato.stanford.edu/archives/sum2011/entries/reasoning-automated/>    (Dec.    26,

723            2014).

724    Python v2.7.3 [Computer software]. Beaverton, OR, Python Software Foundation.

725    Qi, J., Issa, R., Hinze, J., and Olbina, S. (2011). "Integration of safety in design through the use of

726        building information modeling." *Int. Workshop on Computing in Civil Engineering 2011*,

727        ASCE, Reston, VA, 698-705.

728    Saint-Dizier, P. (1994). "Advanced logic programming for language processing." Academic Press,

729        San Diego, CA.

730    Solihin, W., and Eastman, C. (2015). "A knowledge representation approach to capturing BIM

731        based rule checking requirements using conceptual graph." *Proc., CIB W78 2015*, Conseil

732        International du Bâtiment (CIB), Rotterdam, The Netherlands, 686-695.

733    State        of        New        Jersey.        (2014).        "Plan        review        instructions."        <

734        http://www.state.nj.us/dca/divisions/codes/forms/pdf_bcpr/pr_app_guide.pdf>    (Sept.    4,

735        2015).

736    Tan, X., Hammad, A., and Fazio, P. (2010). "Automated code compliance checking for building

737        envelope design." *J. Comput. Civ. Eng.,* 10.1061/ 1195 (ASCE)0887-3801(2010)24:2(203),

738        203-211.

739    Yurchyshyna, A., Faron-Zucker, C., Thanh, N.L., and Zarli, A. (2010). "Adaptation of the domain

740        ontology for different user profiles: application to conformity checking in construction." *Web*

741        *Information Systems and Technologies, Lecture Notes in Business Information Processing*,

742        45, 128-141.

743    Yurchyshyna, A., Faron-Zucker, C., Thanh, N.L., and Zarli, A. (2008). "Towards an ontology-

744        enabled approach for modeling the process of conformity checking in construction." *Proc.,*

745        *CAiSE'08 Forum 20th Intl. Conf. Adv. Info. Sys. Eng.*, dblp team, Germany, 21-24.

746    Zhang, J., and El-Gohary, N. (2013). "Semantic NLP-based information extraction from

747        construction regulatory documents for automated compliance checking." *J. Comput. Civ. Eng.*,

748        10.1061/(ASCE)CP.1943-5487.0000346, 04015014.

749    Zhang, J., and El-Gohary, N.M. (2016a). "Extending building information models semi-

750        automatically using natural language processing techniques." *J. Comput. in Civ. Eng.*,

751        10.1061/(ASCE)CP.1943-5487.0000536, C4016004.

752    Zhang, J., and El-Gohary, N.M. (2016b). "Semantic-based logic representation and reasoning for

753        automated    regulatory    compliance    checking."    *J.    Comput.    in    Civ.    Eng.*,

754        10.1061/(ASCE)CP.1943-5487.0000583, 04016037. Zhang, S., Teizer, J., Lee, J., Eastman,

755        C.M., and Venugopal, M. (2013). "Building information modeling (BIM) and safety:

756        automatic safety checking of construction models and schedules." *Autom. Constr.*, 29(2013),

757        183-195.

758    Zhong, B., Ding, L., Luo, H., Zhou, Y., Hu, Y., and Hu, H. (2012). "Ontology-based semantic

759        modeling of regulation constraint for automated construction quality compliance checking."

760        *Autom. Constr.*, 28, 58-70.

761    Zhou, N. (2012). "B-Prolog user's manual (version 7.8): Prolog, agent, and constraint

762        programming." Afany Software. <http://www.probp.com/manual/manual.html> (Dec. 28,

763        2013).

764     Zouaq, A. (2011). "An overview of shallow and deep natural language processing for ontology

765         learning." *Ontology learning and knowledge discovery using the web: Challenges and recent*

766         *advances*, IGI Global, Hershey, PA, 16-38.