

22 **CE Database subject headings:** Project management; Construction management; Information
23 management; Computer applications; Artificial intelligence.

24 **Author keywords:** Automated compliance checking; Automated information extraction;
25 Automated information transformation; Natural language processing; Semantic systems;
26 Automated construction management systems.

27 **Introduction**

28 Construction projects must comply with a host of regulations. The manual process of compliance
29 checking is, thus, time-consuming, costly, and error-prone (Han et al. 1998; Nguyen 2005;
30 Zhang and El-Gohary 2013c). Automated compliance checking (ACC), as an alternative to
31 manual checking, is expected to reduce the time, cost, and errors of compliance checking (CC)
32 (Tan et al. 2010; Salama and El-Gohary 2013b). In addition, ACC has many other potential
33 benefits, such as: (1) allowing earlier identification of potential non-compliance instances, which
34 could save significant time and cost caused by design modification and/or rework (Ding et al.
35 2006); (2) promoting the adoption of building information modeling (BIM) and increasing the
36 cumulative benefits of adopting BIM, since BIM would enable ACC (Pocas Martins and
37 Abrantes 2010); (3) enabling more efficient incorporation of stakeholder input into project
38 design and exploration of what-if design scenarios, since a designer would be better able to
39 experiment with different design options and check their compliance in a more time-efficient
40 manner (Niemeijer et al. 2009); and (4) reducing violations of regulations due to easier and more
41 frequent CC (Zhong et al. 2012).

42 Due to the many anticipated benefits of ACC, many efforts were undertaken in the area of ACC
43 in construction. The start of these efforts could be dated back to the 1960s, when Fenves et al.

44 (1969) formalized the American Institute of Steel Construction (AISC) specifications into
45 decision tables. These efforts took various approaches to ACC and focused on various ACC
46 purposes (or subdomains). For example, Garrett and Fenves (1987) proposed a strategy to
47 represent design standards using information networks and represent design component
48 properties using data items for ACC of structural designs; Ding et al. (2006) proposed an
49 approach to represent building codes using object-based rules and represent designs using an
50 Industry Foundation Classes (IFC)-based internal model for ACC of accessibility regulations;
51 Tan et al. (2010) proposed an approach to represent building codes and design regulations using
52 decision tables and incorporate simulation results in building information models for ACC of
53 building envelope design; the CORENET (Construction and Real Estate NETwork) project of
54 Singapore (Khemlani 2005) used an approach to represent design information using semantic
55 objects in the FORNAX library (i.e., a C++ library) and represent regulatory rules using
56 properties and functions in FORNAX objects for ACC of building control regulations, barrier
57 free access, and fire code, etc.; and the SMARTcodes project (ICC 2012) of the International
58 Code Council (ICC) used an approach to represent ICC codes in computer-processable tuple
59 format and represent designs using an IFC-based model for ACC of designs with ICC codes.
60 These efforts have all been very important in supporting ACC, and have shown the possibilities
61 of ACC through different system designs and implementations. However, despite their
62 importance, these efforts are limited in their automation capability; existing ACC efforts/systems
63 still require manual effort for the extraction of regulatory requirements from regulatory
64 documents and encoding them in a computer-processable format (Zhong et al. 2012; Zhang and
65 El-Gohary 2013c). To achieve full automation of ACC, this extraction and encoding process
66 needs to be fully automated.

67 To address this gap, the authors are proposing a new approach for automated rule extraction and
68 formalization for supporting ACC (Zhang and El-Gohary 2013a; Zhang and El-Gohary 2013b).
69 The approach utilizes semantic modeling and semantic Natural Language Processing (NLP)
70 techniques (for both information extraction and information transformation) to facilitate
71 automated textual regulatory document analysis (e.g., code analysis) and processing for
72 extracting requirements/rules from these documents and formalizing these requirements/rules in
73 a meaning-rich, computer-processable format. The approach involves developing a set of
74 algorithms and combining them into one computational platform: (1) machine-learning-based
75 algorithms for text classification (TC), (2) rule-based, semantic NLP algorithms for information
76 extraction (IE), and (3) rule-based, semantic NLP algorithms for information transformation
77 (ITr). This paper focuses on presenting the methodology and algorithms for ITr.

78 **Proposed Approach for Automated Rule Extraction and Formalization for Automated** 79 **Compliance Checking**

80 *Proposed Approach*

81 A five-phase, iterative approach for automatically extracting regulatory requirements/rules from
82 textual regulatory documents and formalizing these requirements in a logic format for further
83 automated reasoning is proposed (Figure 1). The five phases are: text classification (TC),
84 information extraction (IE), information transformation (ITr), implementation, and evaluation.
85 TC, IE, and ITr are the main processing phases.

86 

87 TC recognizes relevant sentences in a regulatory text corpus. Relevant sentences are the
88 sentences that contain the types of requirements that are relevant for an ACC scenario (e.g.,
89 environmental requirements in the scenario of environmental CC). Target information in those

90 relevant sentences are extracted and transformed in later IE and ITr processes. The TC process,
91 thus, filters out irrelevant sentences, thereby saving unnecessary processing of irrelevant
92 sentences. Such filtering also avoids unnecessary extraction and transformation errors that may
93 be caused by the processing of irrelevant sentences. The presentation of the TC algorithms and
94 results is outside the scope of this paper. For further details on the authors' work in TC, the
95 reader is referred to Salama and El-Gohary (2013a).

96 IE recognizes the words and phrases in the relevant sentences that carry target information,
97 extracts information from these words/phrases, and labels them with pre-defined information
98 tags. An information tag is a symbol/name indicating a certain type of meaning. For example, the
99 information tag 'subject' carries the semantic meaning that the information instance is a "thing"
100 (e.g., building object) that is subject to a particular regulation or norm; while the information tag
101 'JJ' carries the syntactic meaning that the information instance is an adjective that describes a
102 noun as a modifier. Target information is the information needed to check a specific type of
103 regulatory requirement. For example, for quantitative requirements, the quantified
104 values/measurements of specific properties/attributes are target information. For IE by itself, a
105 seven-phase, iterative methodology is utilized. In the IE methodology, a set of pattern-matching-
106 based IE rules are used. Both syntactic (i.e., related to syntax and grammar, such as part-of-
107 speech (POS) tags) and semantic (i.e., related to context and meaning, such as ontology concepts
108 and relations) text features are used in the IE rules. The presentation of the IE algorithms and
109 results is outside the scope of this paper. For further details on the authors' work in the area of IE,
110 the reader is referred to Zhang and El-Gohary (2013c).

111 ITr takes the extracted information instances and transforms them into logic clauses (i.e., logic
112 statements that can be further used in logic programs) using a set of pattern-matching-based rules.

113 Two types of rules are utilized for ITr: semantic mapping (SeM) rules and conflict resolution
114 (CoR) rules. Several syntactic and semantic text features are used in the rules. A bottom-up
115 method is utilized to handle complex sentence components. A “consume and generate”
116 mechanism is proposed to implement the bottom-up method and execute the SeM rules. The
117 following sections present and discuss the proposed ITr methodology in more detail. The
118 experimental implementation of the methodology in processing quantitative requirements from
119 Chapter 19 of the International Building Code (IBC) 2009 is also presented.

120 *Comparison to the State-of-the-Art*

121 In recent years, a number of research efforts, in domains such as software engineering (Breux
122 and Anton 2008; Kiyavitskaya et al. 2008) and legal compliance (Wyner and Peters 2011), have
123 been studying the extraction of regulatory rules from textual documents. Most of these efforts (1)
124 require manual annotation or mark-up of textual documents; and (2) aim at processing text at a
125 coarser granularity level, i.e., process text into text segments rather than term-level
126 concepts/relations. On the other hand, the proposed approach (1) does not require manual
127 annotation or mark-up of textual documents; and (2) aims at processing text into concepts and
128 relations at the term level (i.e., aims at performing a deeper level of NLP). To the best of the
129 authors’ knowledge, the only work that has taken a somewhat similar approach to the proposed
130 one— since it also does not require manual annotation/mark-up and aims at term-level processing,
131 in addition to utilizing a semantic and logic-based approach – is that by Wyner and Governatori
132 (2013). Wyner and Governatori (2013) have conceptually explored and analyzed the use of
133 semantic parsing and defeasible logic for regulatory rule representation. In comparison, the
134 proposed approach (1) utilizes both syntactic and semantic text features in an integrated way
135 rather than utilizing only semantic information: the use of syntactic text features in addition to

136 semantic ones allows for handling more complex expressions, (2) uses a domain ontology for
137 capturing domain-specific semantic information rather than using generic semantic information
138 produced through generic semantic parsing: capturing and using semantic text features based on
139 domain-specific meaning allows for unambiguous interpretation of concepts/relations/terms (e.g.,
140 “bridge” as an infrastructure instead of the card game) and identification of implicit semantic
141 relations (e.g., “fly ash” is a type of “cementitious material”), (3) uses first order logic (FOL)
142 rather than defeasible logic: FOL is the most widely used in automated reasoning and has been
143 extensively verified for expressivity and simplicity, and (4) has advanced to the stages of
144 implementation, testing, and evaluation: this allows for assessing the validity of the proposed
145 approach using measures of precision and recall.

146 **Background**

147 *Natural Language Processing (NLP)*

148 NLP is a subfield of artificial intelligence (AI) that aims at making natural language text or
149 speech computer-understandable, so that the text or speech could be processed by computers in a
150 human-like manner (Cherapas 1992). Examples of NLP-enabled applications include automated
151 natural language translation and automated text summarization (Marquez 2000). Examples of
152 NLP subtasks include tokenization, POS tagging, semantic role labeling (Gildea and Jurafsky
153 2002), and named entity recognition (Roth and Yih 2004). NLP tasks may take two main
154 approaches: a machine learning (ML)-based approach or a rule-based approach. A ML-based
155 approach utilizes ML algorithms for text processing (e.g., Pradhan et al. 2004), whereas a rule-
156 based approach utilizes manually-coded rules (e.g., Soysal et al. 2010). Rule-based methods
157 require more human effort for rule development, but tend to show better text processing
158 performance (Crowston et al. 2010). From another viewpoint, NLP approaches could be either

159 shallow or deep. Shallow NLP conducts partial analysis of a sentence or extracts partial, specific
160 information from a sentence (e.g., entities or concepts). Deep NLP aims at full sentence analysis
161 towards capturing the entire meaning of a sentence (Zouaq 2011). The state-of-the-art in NLP
162 has achieved reasonable performances for shallow NLP tasks, whereas it is still being challenged
163 by deep NLP tasks. Deep NLP requires elaborate knowledge representation and reasoning which
164 remains to be a challenge for AI (Tierney 2012).

165 In the construction domain, there has been a number of important research efforts that have
166 utilized NLP techniques. For example, Caldas and Soibelman (2003) have conducted ML-based
167 text classification of construction documents. For an overview of some of these efforts, the
168 reader is referred to Zhang and El-Gohary (2013c).

169 ***Rule-Based NLP using Pattern-Matching-Based Rules***

170 Pattern-matching-based rules are widely used in NLP tasks such as POS tagging (Abney 1997;
171 Yin and Fan 2013), information extraction (Califf and Mooney 2003), and text understanding
172 (Goh et al. 2006). The idea of pattern-matching-based rules is to define a set of results when the
173 matching of a pattern of a specific sequence (or structure like a tree) of elements (e.g., characters,
174 tokens, symbols, terms, concepts) occurs. Pattern-matching-based rules have a variety of
175 implementations tailored to different purposes and domains. But, they all share the same rule
176 schema of “if *pattern* then *result*” or the mapping of “from *pattern* to *result*”. For example, in the
177 proposed SeM rules, the result is the transformation of information instances into logic clause
178 elements; while in the proposed CoR rules, the result is the deletion or conversion of certain
179 information instances and/or their information tags to resolve conflicts.

180 ***Semantic Modeling and Semantic NLP***

181 A semantic model aims at capturing the meanings of a domain or topic, usually in a structured
182 manner. Ontology is a widely-used type of semantic model; it is defined as “an explicit
183 specification of a conceptualization” (Gruber 1995). An ontology is, commonly, composed of
184 concept hierarchies, relationships between concepts, and axioms. The axioms are used together
185 with the concepts and relationships to define the semantic meaning of the conceptualization. An
186 ontology is easily reusable and extendable (El-Gohary and El-Diraby 2010). The use of a
187 semantic model could help in NLP tasks. For example, semantic-based IE has been shown to
188 achieve better performance than syntactic-only IE (Soysal et al. 2010; Zhang and El-Gohary
189 2013c).

190 ***Logic-Based Information Representation and Reasoning***

191 There are several types of formally-defined logic with varying degrees of descriptive capabilities
192 (prepositional logic, predicate logic, modal logic, description logic, etc.). Among the different
193 types, FOL is the most widely-used for logic-based inference-making. A Horn Clause (HC) is
194 one of the most restricted forms of FOL. Inference-making in FOL is most efficient using HC
195 logic clauses, because of such restricted form (Saint-Dizier 1994). A HC is composed of a
196 disjunction of literals of which at most one is positive. All HCs can be represented as rules that
197 have one or more antecedents (i.e., left-hand sides (LHSs)) that are conjoined (i.e., combined
198 using ‘and’ operator), and a single positive consequent (i.e., right-hand side (RHS)). For example,
199 “*compliant(T) :- thickness(T) , exterior_basement_wall(W) , has(W,T) ,*
200 *greater_than_or_equal(T, quantity(71/2, inches))*” is a HC; where “,” is the conjunctive operator
201 (i.e., “A , B” means “A and B”) and “:-” is the implication operator (i.e., “B :- A” means “A
202 implies B”). There are three types of HCs: (1) one or more antecedents and one consequent, (2)

203 zero antecedents and one consequent, and (3) one or more antecedents and zero consequents.
204 Inference-making using HCs could be automatically and efficiently conducted, which makes it
205 suitable for supporting automated reasoning for ACC.

206 **Proposed Information Transformation Methodology**

207 The proposed ITr takes a rule-based, semantic NLP approach. It utilizes pattern-matching-based
208 rules to automatically generate logic clauses based on the extracted information instances and
209 their associated patterns of information tags. Both syntactic information tags (i.e., tags tagging
210 syntactic text features, e.g., ‘adjective’ is represented using the POS tag ‘JJ’) and semantic
211 information tags (i.e., tags tagging semantic text features, e.g., ‘compliance checking attribute’ is
212 represented using the semantic tag “a”) are used in defining the patterns. A number of NLP
213 techniques (e.g., POS tagging, term matching) are used to identify the syntactic information tags
214 of each extracted information instance, and a semantic model (an ontology that represents
215 domain knowledge) is used to identify the semantic information tags. The tagged information
216 instances are transformed into HC-type logic clauses using a set of SeM rules and CoR rules.
217 SeM rules define how to process the extracted information instances, based on their associated
218 types of information tags and the context of the information tags, so that the extracted
219 information instances could be transformed correctly into logic clauses. CoR rules resolve
220 potential conflicts that may exist in the processing of different information tags. A bottom-up
221 method is utilized to handle complex sentence components. A “consume and generate”
222 mechanism is proposed to implement the bottom-up method and execute the SeM rules.

223 The following subsections present the proposed ITr methodology (Figure 2) in more detail.

224 **Insert Figure 2**

225 *The Source: Extracted Information Instances*

226 The information source for the ITr process is the set of input information instances that were
227 obtained from the preceding IE process. Information instances have been labeled with
228 information tags during IE. The implemented changes/improvements on the authors' IE work
229 since Zhang and El-Gohary (2013c) are: (1) in addition to semantic information tags, syntactic
230 information tags and combinatorial information tags are also generated for further use in ITr; and
231 (2) instead of the top-down method for handling complex sentence components (processing
232 larger chunks of texts first, then breaking them down to process smaller chunks of texts), a
233 bottom-up method (processing smaller chunks of texts first, then aggregating them to process
234 larger chunks of texts) is adopted because – in the experiments – it has shown to achieve better
235 performance in handling complex sentence components (Zhang and El-Gohary 2013b). As such,
236 in the ITr process, the following three types of information tags (information tags will be shown
237 using single quotes hereafter) are defined and used: (1) semantic information tags, (2) syntactic
238 information tags, and (3) combinatorial information tags.

239 Semantic information tags are information tags that are related to the meaning and context of the
240 labeled information instances. Instances of semantic information tags are recognized based on
241 the concepts and relations in the domain ontology. For example, in the developed ontology, both
242 “transverse reinforcement” and “vertical reinforcement” are subconcepts of the concept ‘subject’.
243 Therefore, the appearances of “transverse reinforcement” (or “transverse reinforcements”) and
244 “vertical reinforcement” (or “vertical reinforcements”) in Chapter 19 of IBC 2009 will be
245 extracted as instances of the semantic information tag ‘subject’. The decision on which concepts
246 and relations are essential to extract and transform is based on the type of requirement (e.g.,
247 quantitative requirements) that is being checked. For example, ‘subject’ is one example of a

248 semantic information tag that is essential in the context of compliance checking of quantitative
249 requirements.

250 Syntactic information tags are information tags that are related to the grammatical role of the
251 labeled information instances. Instances of syntactic information tags are recognized based on
252 their syntactic features. Syntactic information tags carry information that is more general than
253 those carried by semantic information tags. For example, the syntactic information tag ‘noun’ is
254 describing the labeled information instance as a noun, while semantically the noun could
255 possibly belong to a ‘subject’, ‘compliance checking attribute’, or another semantic information
256 tag. In the proposed methodology, POS tags are mainly used as the syntactic features for
257 syntactic information tags. For example, ‘JJ’ is the POS tag for adjective. It is a syntactic
258 information tag for an information instance that describes properties/attributes of a noun. For
259 example, the adjective “habitable” in “habitable room” is describing the functional property of
260 “room”.

261 Combinatorial information tags are compound information tags that are composed of multiple
262 semantic and/or syntactic information tags. For example, the combination of ‘past participle verb’
263 (POS tag ‘VBN’) and ‘preposition’ (POS tag ‘IN’) is a combinatorial information tag
264 (combining two syntactic information tags) that describes a directional passive verbal relation
265 represented by bigrams like “provided by” and “located in”. The combination of ‘adjective’
266 (syntactic information tag - POS tag ‘JJ’) and ‘subject’ (semantic information tag ‘s’) is another
267 example of a combinatorial information tag (combining syntactic and semantic information tags)
268 that describes a ‘subject’ with a certain property.

269 ***The Target: Logic Clauses***

270 The target of the ITr process is the set of output logic clauses which are used to represent the
271 requirements in construction regulations. A HC format is used for such representation, in order to
272 facilitate further automated reasoning using logic programs. One single HC represents one
273 requirement. The RHS of the HC (in Prolog syntax the logical RHS appears to the left of “:-”)
274 indicates the compliance result(s). The LHS of the HC encodes the conditions for the
275 requirement using one or more predicates. Each predicate defines either a concept information
276 instance (e.g., court(C)) or a relation information instance (e.g., has(C,W)). The logic clause
277 elements in a concept predicate are called concept logic clause elements. The logic clause
278 elements in a relation predicate are called relation logic clause elements. Table 1 shows the
279 source and target for a sample sentence.

280 Insert Table 1

281 ***Semantic Mapping (SeM) Rules***

282 The semantic mapping (SeM) rules define how to process the extracted information instances
283 according to their semantic meaning. The semantic meaning of each information instance is
284 defined by: (1) the information tag it is associated with. For example, in Table 1, ‘subject’
285 defines the semantic meaning of “court”, i.e., it defines that “court” is the ‘subject’ of
286 compliance checking; and (2) the context of the extracted information instance, reflected by the
287 information tags of its surrounding information instances. For example, in the following sentence,
288 the semantic meaning of “not less than” (instance of ‘comparative relation’) is defined by the
289 information tag of its surrounding information instance “for each”: “The minimum net area of
290 ventilation openings shall not be less than 1 square foot for each 150 square feet of crawl space
291 area”. “For each”, here, indicates that “not less than” (relation) is not simply a relationship

292 between “net area” (instance of ‘compliance checking attribute’) and “1 square foot” (instance
293 of ‘quantity value’ + ‘quantity unit’), but it is also restricted by “150 square feet of crawl space
294 area” (instance of a ‘quantity value’ + ‘quantity reference’). The interpretation of this
295 requirement is that the quantity requirement on “minimum net area of ventilation openings” will
296 increase 1 foot for each additional “150 square feet of crawl space area”.

297 The semantic meanings of information instances are utilized in patterns on the LHS of SeM rules.
298 For the example in Table 1, the corresponding SeM rule pattern is ‘subject’ + ‘modal verb’ +
299 ‘negation’ + ‘be’ + ‘comparative relation’ + ‘quantity value’ + ‘quantity unit’ + ‘preposition’ +
300 ‘compliance checking attribute’. An SeM rule with this LHS pattern will transform the
301 information instances into the logic clause shown in the last row of Table 1. A sample action
302 defined on the RHS of this SeM rule is: “Generate predicates for the ‘subject’ information
303 instance, the ‘attribute’ information instance, and a ‘has’ information instance. The two
304 arguments of the ‘has’ information instance are from the ‘subject’ predicate and the ‘attribute’
305 predicate, respectively”. Accordingly, the following logic clause elements are generated for the
306 following statement, since “court” is recognized as a ‘subject’ information instance and “width”
307 as an ‘attribute’ information instance.

- 308 • Sentence: “Courts shall not be less than 3 feet in width”
- 309 • Logic Clause Elements: court(Court), width(Width), has(Court,Width)

310 The ITr method is intended to process each term of a sentence in a sequential manner. In general,
311 sequential processing for information transformation normally requires information instances
312 that are matched by patterns (in SeM rules) to be strictly located next to each other. Such a rigid
313 processing requirement could cause difficulty in processing sentences with different structures.

314 To avoid that, the proposed SeM rules do not follow such a rigid requirement. Instead, the SeM
315 rules allow for “look-back searching” (i.e., searching to the left of the matched words) and “look-
316 ahead searching” (i.e., searching to the right of the matched words) to find instances that match
317 certain information tags in a pattern. For example, in the following pattern, the instance of the
318 first ‘subject’ does not have to be located right next to the instance of ‘preposition’: “ ‘subject’ +
319 ‘preposition’ + ‘subject’ “. It is only required to be the ‘subject’ instance that is closest to the
320 ‘preposition’ instance from the left. “Look-back searching”, here, searches to the left of the
321 matched word for ‘preposition’ to find the closest ‘subject’ instance when the later part of the
322 pattern “ ‘preposition’ + ‘subject’ ” is matched. This allows for more flexibility in the use of
323 SeM rules to handle sentence complexities (e.g., those incurred by cases such as tail recursive
324 nested clauses). For example, an SeM rule uses the following pattern P1 to match the last three
325 information instances in InS1 (‘s’ for ‘subject’, ‘VBP’ for ‘non-3rd person singular present verb’,
326 ‘dpvr’ for ‘directional passive verbal relation’, and ‘VB’ for ‘base form verb’), finds the first
327 information instance in InS1 through “look-back searching”, and generates the logic clause
328 elements LC1 for the part of sentence S1:

- 329 • Pattern P1: ‘non-3rd person singular present verb’ ‘directional passive verbal relation’
330 ‘base form verb’
- 331 • Information Instances InS1: (‘connection’, ‘s’) ... (‘are’, ‘VBP’), (‘designed_to’, ‘dpvr’),
332 (‘yield’, ‘VB’)
- 333 • Sentence S1: “Connections that are designed to yield shall be capable of ...”
- 334 • Logic Clause Elements LC1: connection(Connection), yield(Yield),
335 designed_to(Connection, Yield)

336 In the proposed methodology, application-specific SeM rules are developed based on a
337 randomly selected sample of text (called “development text”, which is also used for text analysis
338 and further development of CoR rules). For developing a set of SeM rules for ITr, a three-step,
339 iterative methodology that shall be applied to each sentence is proposed: (1) find all relations in a
340 sentence (e.g., “of” and “not exceed” in the sentence “Spacing of transverse reinforcement shall
341 not exceed 8 inches.”), (2) for each relation, run the existing SeM rule set to check if the rule set
342 can generate the corresponding logic clause elements correctly and define the subsequent action
343 based on the following three cases: (a) if the corresponding logic clause elements are correctly
344 generated, then move to check the next relation, (b) if the corresponding logic clause elements
345 are incorrectly generated, then create a new SeM rule with a more specific pattern (i.e., a longer
346 pattern with more features) than the applied SeM rule and add it to the rule set with a higher
347 priority, and (c) if the corresponding logic clause elements are not generated, then create a new
348 SeM rule and add it to the rule set; and (3) after all relations have been checked, run the updated
349 SeM rule set on all checked sentences and check if errors have been introduced due to the added
350 SeM rules. If errors have been introduced, then identify the source(s) of errors (i.e., the rule(s)
351 that introduced the errors) and adjust those rules as necessary.

352 ***Conflict Resolution (CoR) Rules***

353 The conflict resolution (CoR) rules resolve conflicts between information tags. Two types of
354 CoR rules are used: deletion CoR rules and conversion CoR rules. Deletion CoR rules resolve
355 conflicts between information tags by deleting certain information instances. For example, the
356 following deletion CoR rule CoR1 is used to delete redundant information instances InS2 (‘cr’
357 for ‘candidate restriction’) from the set of extracted information instances InS3 (‘s’ for ‘subject’)
358 for the sentence S2:

359 • Deletion CoR Rule CoR1: “if an information instance has the tag ‘subject’ and it
360 subsumes its following information instance(s), then delete its following information
361 instance(s).”

362 • Information Instances InS2: (‘exterior’, ‘cr’), (‘basement’, ‘cr’), (‘wall’, ‘cr’)

363 • Information Instances InS3: (‘exterior basement wall’, ‘s’), (‘exterior’, ‘cr’), (‘basement’,
364 ‘cr’), (‘wall’, ‘cr’)

365 • Sentence S2: “The thickness of exterior basement walls and foundation walls shall be not
366 less than 71/2 inches.”

367 Conversion CoR rules resolve conflicts between information tags by converting information tags
368 of information instances into other types of information tags. For example, the following
369 conversion CoR rule CoR2 is used to convert information tags in information instances InS4 (‘s’
370 for ‘subject’, ‘I’ for ‘inter clause boundary relation’, and ‘a’ for ‘compliance checking attribute’)
371 to information tags in information instances InS5 (‘IN’ for ‘preposition’) for the sentence S3:

372 • Conversion CoR Rule CoR2: “if ‘with’ is directly followed by an information instance
373 that has the information tag ‘compliance checking attribute’ and ‘with’ has the
374 information tag ‘inter clause boundary relation’, then convert the information tag of ‘with’
375 to ‘preposition’.”

376 • Information Instances InS4: (‘wall segment’, ‘s’), (‘with’, ‘I’),
377 (‘horizontal_length_to_thickness_ratio’, ‘a’)

378 • Information Instances InS5: (‘wall segment’, ‘s’), (‘with’, ‘IN’),
379 (‘horizontal_length_to_thickness_ratio’, ‘a’)

380 • Sentence S3: “Wall segments with a horizontal length-to-thickness ratio less than 2.5
381 shall be designed as columns.”

382 In the proposed rule-based ITr, the CoR rules are executed before the SeM rules, after the
383 information instances have been extracted by the IE process. The development of CoR rules is
384 needed when conflicts between SeM rules cannot be resolved by adjusting SeM rule patterns and
385 actions. For developing a set of CoR rules for ITr, a five-step methodology is proposed: (1) find
386 information tags that are the sources of errors through pattern analysis of conflicting SeM rules,
387 (2) for each conflict, create a new candidate CoR rule to resolve the conflict, (3) try the candidate
388 rule and empirically analyze whether the conflict was resolved without introducing new conflicts
389 or not, (4) if the trial was successful, then add the candidate CoR rule as a new rule to the
390 existing CoR rule set, and if the trial was unsuccessful, then iterate Steps 3 and 4 until a
391 successful trial is found, and (5) after each new CoR rule is added, check all SeM rules and
392 update them as necessary according to the changes in information tags caused by the new CoR
393 rule.

394 ***Bottom-up Method for Handling Complex Sentence Components***

395 Due to the variability of natural language expressions and structures, sentences used in
396 regulatory provisions could be very complex. For example, phrases and clauses could be
397 continuously attached/nested to a sentence to constantly enrich it with more relevant information.
398 Complex sentences are difficult to process for information extraction and transformation.
399 Complex sentence components are intermediately-processed segments of text that are: (1)
400 expressed using a variety of natural language structure patterns, and (2) composed of multiple
401 concepts and relations. Complex sentence components are more likely to result in complex
402 sentence structures by embedding in or attaching more concepts and relations to a sentence.
403 Figure 3 shows a complex sentence from IBC 2006. Two methods were explored in handling
404 complex sentence components: top-down method and bottom-up method (Figure 4). The top-

405 down method starts from the top level (i.e., full sentence) and proceeds down to identify and
406 process complex sentence components. The bottom-up method starts from the lowest level (i.e.,
407 single terms/concepts/relations in a sentence) and proceeds up to identify and process complex
408 sentence components. The bottom-up method is employed in the proposed ITr approach, because
409 – based on the authors’ previous work – it has shown to achieve better performance than the top-
410 down method (Zhang and El-Gohary 2013b).

411 Insert Figure 3

412 Insert Figure 4

413 In the bottom-up method, the SeM rules are used to process sentences starting from the lowest
414 level, i.e., starting from information instances (which correspond to single
415 terms/concepts/relations in a sentence). The information instances in the source text are put into
416 lists – one list for each sentence and are processed one by one until all information instances
417 have been processed. The order of the instances in the list is determined based on their order in
418 the original sentence.

419 To apply the bottom-up method, the authors propose a new “consume and generate” mechanism
420 to execute the SeM rules in a sequential manner. This mechanism follows the heuristics of the
421 “sliding window” method in computational research (i.e., a sequence of data is sequentially
422 processed, segment by segment, and each segment has a predefined fixed length (i.e., the
423 “window size”)) and the mechanism of transcription in genetics domain (i.e., a sequence of DNA
424 is sequentially transcribed, segment by segment, and each segment has a length of about 17 base-
425 pair). The “consume and generate” mechanism processes all text segments that match an SeM
426 rule pattern, where each segment matches a pattern of one SeM rule and each pattern consists of
427 information tags for a sequence of information instances. However, in comparison to the “sliding

428 window” method, the segment length in the proposed “consume and generate” mechanism is not
429 fixed across patterns to allow for flexibility in capturing complex sentence structures. The length
430 of each segment is determined according to the number of information tags in the corresponding
431 SeM rule pattern. For example, the following pattern P2 has a segment length of three and
432 matches the information instances InS6 for the part of sentence S4 to generate logic clause
433 elements LC2:

- 434 • Pattern P2: ‘compliance checking attribute’ ‘of’ ‘subject’
- 435 • Information Instances InS6: (‘area’, ‘a’), (‘of’, ‘OF’), (‘space’, ‘s’)
- 436 • Sentence S4: “The net free ventilating area shall not be less than 1/150 of the area of the
437 space ventilated ...”
- 438 • Logic Clauses Elements LC2: space(Space), area(Area), has(Space, Area)

439 The “consume and generate” mechanism allows for backward matching: if information instances
440 extracted from a segment of text match the later part of a pattern, then the information instance(s)
441 extracted from preceding text are checked for matching of the earlier part of the same pattern,
442 and corresponding logic clauses are generated if the check succeeds. For example, the following
443 information tags InT1 are associated with the five information instances from the part of
444 sentence S5. After the first three information instances InS7 are processed based on matching
445 with the pattern P3, two information instances “or” and “space” remain. These two remaining
446 information instances only match the later part (i.e., second and third information tags) of the
447 pattern P4 for ‘conjunctive subject’. Normally, this partial matching would not initiate the
448 processing of the information instances. However, under the proposed backward matching
449 mechanism, the preceding information instance “interior room” is checked for the matching of
450 the earlier part of the pattern for “conjunctive subject” (i.e., the first information tag: ‘subject’).

451 Since “interior room” matches ‘subject’, the SeM rule for “conjunctive subject” gets applied and
452 the two remaining information instances are processed to generate the logic clause elements LC3
453 (where “;” is the disjunctive operator (i.e., “A ; B” means “A or B”)).

- 454 • Information Tags InT1: ‘compliance checking attribute’, ‘of’, ‘subject’, ‘conjunctive
- 455 term’, ‘subject’
- 456 • Sentence S5: “...the floor area of the interior room or space...”
- 457 • Information Instances InS7: “floor area”, “of”, “interior room”
- 458 • Pattern P3: ‘compliance checking attribute’ + ‘of’ + ‘subject’
- 459 • Pattern P4: ‘subject’ + ‘conjunctive term’ + ‘subject’
- 460 • Logic Clause elements LC3: interior_room(Interior_room); space(Interior_room)

461 ***Validation***

462 Results are evaluated in terms of precision, recall, and F1 measure. Precision is the number of
463 correctly generated logic clause elements divided by the total number of generated logic clause
464 elements. Recall is the number of correctly generated logic clause elements divided by the total
465 number of logic clause elements that should be generated. F1 measure is the harmonic mean of
466 precision and recall, assigning equal weights to precision and recall. Ideally, both 100% recall
467 and precision are desired. However, given the inherent trade-off between the two measures, it is
468 difficult to achieve such a result. The ultimate goal for ACC is, therefore, to achieve 100% recall
469 of non-compliance instances – with high precision.

470 **Experimental Implementation and Validation**

471 For testing and validation, the proposed ITr methodology was empirically implemented in
472 transforming information instances of quantitative requirements, which were automatically
473 extracted from the IBC 2009, into logic clauses.

474 *Source Text Selection*

475 The proposed ACC approach and ITr methodology are intended to process information from a
476 variety of construction-related textual regulatory documents (e.g., building codes, environmental
477 regulations, safety regulations and standards). Since building codes are the primary sets of
478 regulations governing the design, construction, operation, and maintenance of residential and
479 commercial buildings, they were chosen for testing the proposed ITr methodology. In the U.S.,
480 almost all state authorities (except for Delaware, Massachusetts, Mississippi, and Missouri)
481 adopt versions of the IBC by ICC. Thus, IBC was selected as the source text corpus. More
482 specifically, IBC 2006 and IBC 2009 were selected because of their availability and easiness for
483 comparison (with the authors' previous NLP work in which IBC 2006 and IBC 2009 were used
484 for testing and validation) (Zhang and El-Gohary 2013c).

485 The SeM and CoR rules were developed based on Chapters 12 and 23 of IBC 2006, and the
486 proposed ITr algorithms were tested in processing information instances of "quantitative
487 requirements" that were extracted from Chapter 19 of IBC 2009. A quantitative requirement is a
488 requirement which defines the relationship between an attribute of a certain building
489 element/part and a specific quantity value (or quantity range). For example, the following
490 sentence, states that the width (attribute) of court (building element/part) should be greater than
491 or equal to 3' (quantity value): "Courts shall not be less than 3 feet in width". The authors decided
492 to The experiment on the extraction of quantitative requirements because: (1) IBC 2006 and IBC

493 2009 describe many quantitative requirements (e.g., on average, quantitative requirements
494 represent 41% of the requirements in Chapters 12 and 23 of IBC 2006 and Chapter 19 of IBC
495 2009), which ensures a sufficient amount of relevant sentences for development and testing; and
496 (2) sentences describing quantitative requirements appear to be more complex than those
497 describing other types of requirements (e.g., existential requirements, which requires the
498 existence of a certain building element/part), which implies that they are more difficult to
499 process. This makes quantitative requirements good candidates for testing.

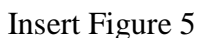
500 ***Tool Selection***

501 The proposed TC, IE, and ITr algorithms were combined into one computational platform. The
502 representation of Prolog was selected for logic clause representation, in order to facilitate future
503 CR. Prolog is an approximate realization of the logic programming computational model on a
504 sequential machine (Sterling and Shapiro 1986). It is the most popular logic programming
505 language with a reasoner. The syntax of B-Prolog was used. B-Prolog is a Prolog system with
506 extensions for programming concurrency, constraints, and interactive graphics. It has bi-
507 directional interface with C and Java (Zhou 2012). To facilitate quantitative reasoning, a set of
508 built-in rules were developed to perform arithmetic and comparative operations on the proposed
509 quantitative representation. The TC and IE algorithms were implemented using the General
510 Architecture for Text Engineering (GATE) tools (Univ. of Sheffield 2013). GATE has a variety
511 of built-in tools for a variety of text processing functions (e.g., tokenization, sentence splitting,
512 POS tagging, gazetteer compiling, and morphological analysis). For ITr, the SeM rules and CoR
513 rules were implemented using Python programming language (v3.3.2). The “re” module (i.e.,
514 regular expression module) in Python was used for pattern matching, so that each extracted
515 information instance could be used for subsequent processing steps based on their information

516 tags (example tags are shown in Figure 3). A domain ontology was developed and used to
517 facilitate semantic IE and ITr. In developing the ontology, the ontology development
518 methodology in El-Gohary and El-Diraby (2010) was followed. The GATES' built-in ontology
519 editor was used for ontology building and editing.

520 ***Information Representation***

521 Two types of logic statements in B-Prolog syntax were utilized: facts and rules. A rule has the
522 form: "H :- B1, B2, ..., Bn. (n>0)". H, B1, ..., Bn are atomic formulas. H is called the head, and
523 the RHS of ':-' is called the body of the rule. A fact is a special kind of rule whose body is
524 always true (Zhou 2012). Each requirement rule in IBC 2006 and IBC 2009 is represented as one
525 single B-Prolog rule. Instances of concepts are represented using unary predicates. For example,
526 the information instance "floor" is represented by the predicate "floor(F)", with "floor" being the
527 predicate name and the variable "F" (all variables in B-Prolog start with capitalized letter) being
528 the argument for the predicate. Instances of relations are represented using binary or n-ary
529 predicates. For example, "provided with" is a relation which is represented as the predicate
530 "provided_with(A,B)", while the variables "A" and "B" could be defined in the predicates
531 interior_space(A) and space_heating_system(B). Each design fact, on the other hand, is
532 represented using one B-Prolog fact. The B-Prolog reasoner can then automatically reason about
533 the facts and rules and, accordingly, determine the compliance checking result(s). An example is
534 shown in Figure 5.

535 

536 *Information Tags*

537 A total of 40 information tags were developed for use in the SeM rules and CoR rules for ITr. A
538 total of 17, 22, and 1 semantic information tags, syntactic information tags, and combinatorial
539 information tags were used, respectively.

540 Two main types of semantic information tags were defined (as per Figure 6): essential
541 information tags and secondary information tags. Essential information tags are tags for
542 information that must be defined for this specific type of requirement. Six main types of essential
543 information tags were defined for quantitative requirements: subject, compliance checking
544 attribute, comparative relation, quantity value, quantity unit, and quantity reference. A 'subject'
545 is an ontology concept; it is a "thing" (e.g., building object, space) that is subject to a particular
546 regulation or norm. A 'compliance checking attribute' is an ontology concept; it is a specific
547 characteristic of a 'subject' by which its compliance is assessed. A 'comparative relation' is an
548 ontology relation which is commonly-used for comparing quantitative values (i.e., comparing an
549 existing value to a required minimum or maximum value). Five subtypes of comparative
550 relations were further defined: 'greater than or equal to', 'greater than', 'less than or equal to',
551 'less than', and 'equal to'. A 'quantity value' is a value, or a range of values, which defines the
552 quantified requirement. A 'quantity unit' is the unit of measure for the 'quantity value'. A
553 'quantity reference' is a reference to another quantity (which includes a value and a unit).

554 Secondary information tags are tags for information that are not necessary for this specific type
555 of requirement, but may exist in defining the requirement. Two main types of secondary
556 information tags were defined for quantitative requirements: 'restriction' and 'exception'. A
557 'restriction' is a concept that places a constraint on the 'subject', 'compliance checking attribute',
558 'comparative relation', pair of 'quantity value' and 'quantity unit', pair of 'quantity value' and

559 'quantity reference', or the full requirement. A 'subject restriction' is a concept that places a
560 constraint on the 'subject'. Two subtypes of 'subject restriction' were further defined: 'possessive
561 subject restriction' and 'nonpossessive subject restriction'. A 'possessive subject restriction' places
562 a possessive constraint on the 'subject', thereby restricting the 'subject' to possess certain
563 building parts or properties. For example, in the following requirement sentence, "having
564 windows opening on opposite sides" is a 'possessive subject restriction' on "court": "Courts
565 having windows opening on opposite sides shall not be less than 6 feet in width". A
566 'nonpossessive subject restriction' places a nonpossessive constraint on the 'subject', thereby
567 restricting the 'subject' not to possess certain building parts or properties. A 'compliance
568 checking attribute restriction' places a constraint on the 'compliance checking attribute', thereby
569 restricting the 'compliance checking attribute' to a more specific type. For example, in the
570 following requirement sentence, "to the outdoors" is a 'compliance checking attribute restriction'
571 on "minimum openable area": "The minimum openable area to the outdoors shall be 4 percent of
572 the floor area being ventilated". A 'comparative relation restriction' places a constraint on the
573 'comparative relation', thereby restricting the 'comparative relation' using new conditions. For
574 example, in the following requirement sentence, "for each 150 square feet of crawl space area" is
575 a 'comparative relation restriction' on "not less than": "The minimum net area of ventilation
576 openings shall not be less than 1 square foot for each 150 square feet of crawl space area". A
577 'quantity restriction' places a constraint on the 'quantity value' + 'quantity unit'/'quantity
578 reference' pair, thereby specifying the properties (e.g., range) of the pair. A 'full requirement
579 restriction' places a constraint on the whole quantitative requirement, thereby restricting the
580 quantitative requirement with new preconditions. An 'exception' defines a condition where the
581 described requirement does not apply.

582 For syntactic information tags, the Hepple POS Tagger was used to generate POS tag features.
583 Some additional syntactic features that were not in the Hepple POS Tagger (e.g., the preposition
584 “of”) were also defined. Each selected POS type and defined syntactic feature represents a
585 syntactic information tag such as adjective (POS tag ‘JJ’) and preposition “of” (the literal “OF”).
586 One combinatorial information tag was defined for use in this implementation and was called
587 ‘directional passive verbal relation’, which is the combination of ‘past participle verb’ (POS tag
588 ‘VBN’) and ‘preposition’ (POS tag ‘IN’). Combinatorial information tags are expressive and
589 flexible. Thus, more combinatorial information tags may be defined and used if more complex
590 information tags are needed to capture complex meanings or patterns.

591 Insert Figure 6

592 ***Gold Standard***

593 The gold standard for Chapter 19 of IBC 2009 was developed semi-automatically. In the authors’
594 previous work, all sentences that include a number (both appearances of digits and words forms
595 of a number) were automatically extracted to ensure a 100% recall of sentences describing
596 quantitative requirements. Then, one of the authors manually deleted false positive sentences.
597 After that, one of the authors manually coded the logic clauses based on the extracted
598 information instances from each sentence. The gold standard was reviewed by two other
599 researchers to verify its correctness. Because of the unambiguous nature of quantitative
600 requirements, along with the well-defined information representation that is used in the proposed
601 methodology, there was an agreement in formulating the gold standard. For Chapter 19, 62
602 sentences containing quantitative requirements were recognized. Correspondingly, 62 logic
603 clauses were coded. In these 62 logic clauses, 1901 logic clause elements were identified,

604 including 568 logic clause elements for describing concepts and 1333 logic clause elements for
605 describing relations between concepts.

606 *Algorithm Implementation*

607 The proposed ITr methodology was implemented using Python programming language. The
608 processing steps of an example sentence and the pseudo codes for the main algorithm and the
609 “consume and generate” mechanism are shown in Figure 7, Figure 8, and Figure 9, respectively.

610 Insert Figure 7

611 Insert Figure 8

612 Insert Figure 9

613 As shown in Figure 7, the IE process tags the original sentence with information tags (from Part I
614 to Part II). The main ITr algorithm then represents each information instance in the tagged
615 sentence into a four-tuple (from Part II to Part III). The CoR rules in the main algorithm then
616 process the information instance tuple list to resolve conflicts between tuples (from Part III to
617 Part IV). The “consume and generate” code then executes the set of SeM rules to process each
618 tuple in the list and generate logic clause elements based on matching of SeM rule patterns (from
619 Part IV to Part V). For each information instance, the four-tuple is used to store: (1) the
620 information instance itself, (2) the location of the information instance in the corresponding
621 sentence (represented by the starting point of the information instance in the sentence), (3) the
622 length of the information instance in terms of number of letters, and (4) the information tag of
623 the information instance (e.g., ‘Interior’, 0, 15, and ‘s’ for the first information instance in Part
624 III of Figure 7).

625 In the main algorithm (Figure 8), the CoR rules are executed through the function “resolve
626 conflicts”. Then, the SeM rules are executed using the “consume and generate” code to process

627 the conflict-free information instances for each sentence of the source text file (in the format of a
628 list of four tuples) to generate and display the corresponding logic clause. As shown in Figure 9,
629 the “consume and generate” code checks through the patterns for each SeM rule (*PATTERN1*,
630 *PATTERN2*, *PATTERN3*...) and generates logic clauses as a result of matching to SeM rules. In
631 case of no matching, the default negative step length enables backward matching.

632 **Experimental Results and Discussion**

633 The proposed ITr algorithms were tested in transforming information instances of quantitative
634 requirements, which were automatically extracted from Chapter 19 of IBC 2009, into logic
635 clauses. The following two experiments were conducted for comparing the performances of two
636 methods of information representation: (1) using essential semantic information tags only, and (2)
637 using essential, as well as secondary, semantic information tags.

638 In Experiment #1, only the essential semantic information tags were used: ‘subject’, ‘compliance
639 checking attribute’, ‘comparative relation’, ‘quantity value’, ‘quantity unit’, and ‘quantity
640 reference’. A subset of the gold standard (including logic clause elements corresponding to the
641 essential semantic information instances) was used as the gold standard for Experiment #1. A
642 total of 53 and 11 SeM and CoR rules, respectively, were developed.

643 In Experiment #2, both essential and secondary information tags were used. Figure 3 shows
644 examples of some of the information tags that were used. A total of 297 and 9 SeM and CoR
645 rules, respectively, were encoded. The gold standard of Experiment #2 (the full gold standard set)
646 contains 177% more logic clause elements than those in the gold standard of Experiment #1.
647 This shows that for quantitative requirements, the source text contains much secondary
648 information instances.

649 The SeM rules that were developed in the experiments are classified into four main types: simple
650 SeM rules, multiple action SeM rules, multiple condition SeM rules, and complex SeM rules. A
651 simple SeM rule is the simplest type where a strict SeM pattern directly maps to a logic clause.
652 For multiple action SeM rules, other actions (called “supportive actions”) such as “look-ahead
653 searching” and “look-back searching” are involved in addition to mapping SeM patterns to logic
654 clauses. For multiple condition SeM rules, the mapping from SeM patterns to logic clauses are
655 encoded in subrules to handle subtly different cases in rule conditions such as the existence/non-
656 existence status of certain information instances. A complex SeM rule is a combination of the
657 first three types of rules; it utilizes both supportive actions and subrules to support mappings
658 from SeM patterns to logic clauses.

659 The logic clauses generated from the SeM rules are classified into three main types: single
660 predicate logic clauses, multiple predicate logic clauses, and compound predicate logic clauses.
661 A single predicate logic clause includes only one single predicate (e.g., “space(Space)”). A
662 multiple predicate logic clause includes more than one predicate (e.g., “space(Space), area(Area),
663 has(Space, Area)”). A compound predicate logic clause has predicate(s) that embed other
664 predicate(s) as argument(s) (e.g., “greater_than_or_equal(T, quantity(71/2, inches))”).

665
666 Table 2 shows the patterns of the most applied SeM rules (i.e., rules applied at least three times)
667 in the experiments. The patterns of the rest of the applied SeM rules are shown in Table 3.

668 Insert Table 2

669 Insert Table 3

670 The overall performance results of Experiment #1 and Experiment #2 are summarized in Table 4
671 and Table 5, respectively.

672 Insert Table 4

673 Insert Table 5

674 A comparison between the results of Experiment #1 and those of Experiment #2 is summarized
675 in Table 6. The number of information tags in Experiment #2 increased 400% from that used in
676 Experiment #1. The increase in the number of SeM rules was of similar magnitude (460%).
677 Through analysis, the causes of this increase in the number of SeM rules were found to be: (1)
678 the use of more information tags increases the length of patterns in SeM rules, which in turn
679 increases the specificity of each pattern; and (2) the use of more information tags increases the
680 complexity of patterns in SeM rules, which in turn increases the possible number of patterns. In
681 contrast to SeM rules, the number of CoR rules decreased from Experiment #1 to Experiment #2.
682 This results from the use of more information tags, which leads to better distinguishable
683 information instances, and in turn leads to less conflicts between information instances.

684 The algorithms achieved 92.5% and 98.2%, 95.1% and 99.1%, and 93.8% and 98.6% overall
685 precision, recall, and F1 measure for Experiment #1 and Experiment #2, respectively. Both
686 precision and recall improve in Experiment #2, because the use of more information tags could:
687 (1) better distinguish and capture the variations in expressions; and (2) help define SeM rules
688 with more specificity in patterns. Based on the comparative analysis, the following conclusion
689 can be drawn: the use of more information tags helps in improving the performance of
690 information transformation.

691 Insert Table 6

692 The precisions of relation logic clause elements are lower than other precision and recall values
693 across Experiment #1 and Experiment #2. Through analysis, four main causes for this relatively
694 lower performance of precision (89.8% and 97.5% for Experiment #1 and Experiment #2,
695 respectively) of relation logic clause elements are recognized: (1) Structural ambiguity caused by
696 conjunctive terms: For example, in the following part of sentence, there are two possible
697 syntactic uses of “and” – either linking “wall piers” and “such segments” or linking the
698 preceding clause and the following clause: “...shear wall segments provide lateral support to the
699 wall piers and such segments have a total stiffness...”. The ability of the SeM rules to handle
700 structural ambiguity is limited by the development text, which may lead to errors; (2) Incorrect
701 tagging during IE: For example, “professional” (in “registered design professional”) was
702 incorrectly tagged as an adjective instead of noun. This is due to the imperfection of state-of-the-
703 art POS tagging methods; (3) Errors due to morphological analysis (MA): MA was used for
704 improving the recall of semantic information instances by finding all forms of a term based on its
705 lexical form. However, while useful in this regard, MA also introduced false positive instances.
706 For example, as a result of MA, “supported” was stemmed into “support”, matched with the
707 concept “support” in the ontology, and as a result incorrectly recognized as an instance of
708 ‘subject’; and (4) Errors caused by certain SeM rules: For example, an SeM rule selects the
709 immediate left neighbor of a preposition as the first argument of that preposition. In cases where
710 the immediate left neighbor of a preposition is not its real first argument, this SeM rule causes
711 errors. For example, in the following part of sentence, “gypsum concrete” was mistakenly
712 identified as the first argument rather than “clear span”: “clear span of the gypsum concrete
713 between supports”.

714 Analyzing other errors (other than those influencing precision of relation logic clause elements),
715 two additional causes of errors are recognized: (1) Missing tags in IE: For example, based on the
716 concepts in the ontology, “connection” should have been semantically-tagged as ‘subject’.
717 However, in a few instances, it was missing the ‘subject’ information tag. This is due to the
718 inherent errors in the NLP tools that were used (no existing NLP tool can achieve 100%
719 performance); and (2) Error in processing sentences with uncommon syntactic expression
720 structures: For example, in the part of sentence “...which have been water soaked for at least 24
721 hours...”, “soaked” (‘compliance checking attribute’) was not recognized because: (a) “soaked”
722 was not semantically-recognized because the ontology did not cover this concept, and (b) the
723 syntactic feature of “soaked” (i.e., past participle) was not a common syntactic expression for
724 ‘compliance checking attribute’ (in contrast, noun is a common expression for ‘compliance
725 checking attribute’).

726 **Limitations and Future Work**

727 The experimental results show that the proposed approach is promising in automatically
728 transforming the extracted information instances into logic clauses for further compliance
729 reasoning. In spite of the high performance that was achieved (98.2%, 99.1%, and 98.6% for
730 precision, recall, and F1 measure, respectively), three main limitations of this work are
731 acknowledged, which the authors plan to address as part of their ongoing/future research. First,
732 the methodology was only tested on processing quantitative requirements. The types of semantic
733 patterns and conflicts in other types of requirements (e.g., existential requirements) may vary and,
734 thus, may lead to different performance results. Although the processing of other types of
735 requirements is expected to be less or equally complex than that of quantitative requirements –
736 and thus is expected to have similar or better performance, in future work, the authors plan to test

737 the proposed methodology on other types of requirements (e.g., existential requirements) for
738 validation. Second, due to the large amount of manual effort required in developing a gold
739 standard, the proposed ITr algorithms were tested only on one chapter of IBC 2009. Similar high
740 performance is expected when testing on other chapters of IBC and on other regulatory
741 documents, since all regulatory documents share similarities in expressions. However, different
742 performance results might be obtained due to the possible variability of text across different
743 chapters or different regulatory documents. As such, in future work, the authors plan to test the
744 proposed ITr methodology on more chapters of IBC 2009 and on other types of regulatory
745 documents (e.g., environmental regulations). Third, the validation of the proposed ITr algorithms
746 was focused on precision and recall. At this stage, the computational efficiency of the proposed
747 algorithms was not evaluated, although it was taken into consideration when developing the
748 algorithms. For example, the more efficient and stable merge sort (rather than quick sort) was
749 used when a sorting algorithm was needed. In future work, the authors plan to perform
750 algorithm optimization to improve the computational efficiency of the proposed algorithms, if/as
751 necessary.

752 **Contribution to the Body of Knowledge**

753 This research contributes to the body of knowledge in four main ways. First, domain-specific,
754 semantic NLP-based information processing methods that can achieve full sentence processing
755 and information extraction (i.e., all terms of a sentence are processed), as opposed to partial
756 sentence processing and information extraction (i.e., only specific terms/concepts are
757 processed/extracted) are offered. Domain-specific semantics allow for analyzing complex
758 sentence structures that would otherwise be too complex and ambiguous for automated IE and
759 ITr, recognizing domain-specific text meaning, and in turn allowing for

760 processing/understandability of full sentences. Full sentence processing/understandability allows
761 for a deeper level of NLP, namely natural language understanding. Second, this research shows
762 that a hybrid approach that combines rule-based NLP methods and semantic NLP methods could
763 achieve high performance for the combination of IE and ITr from/of regulatory text, in spite of
764 the complexity inherent in natural language text. Domain-specific expert NLP knowledge
765 (encoded in the form of rules), along with domain knowledge (represented in the form of an
766 ontology), facilitates deep text processing/understandability. Previous work (Zhang and El-
767 Gohary 2013c) showed high performance for rule-based, semantic IE. This paper further shows
768 high performance for rule-based, semantic ITr. Third, a new context-aware and flexible way of
769 utilizing pattern-matching-rule-based methods through the use of context-aware semantic
770 mapping rules is offered. This way of utilizing pattern-matching-based rules captures the details
771 (in terms of the expression, language structure, etc.) of complex sentence components, in a
772 context-aware manner, and through flexible pattern lengths. Fourth, a new mechanism
773 ("consume and generate" mechanism) for processing and transforming complex regulatory text
774 into logic clauses is offered. The proposed mechanism follows the bottom-up method, which has
775 shown based on the experimental results to outperform the top-down method in ITr. The high
776 performance that the mechanism achieved verifies that the bottom-up method is suitable for such
777 ITr tasks.

778 From a practical perspective, this work is expected to have significant impacts on four main
779 levels. First, this work facilitates ACC in the construction domain. ACC could bring down the
780 time, cost, and errors of the checking process; promote compliance of construction projects to
781 various regulations (due to easier and more frequent checking); and encourage the adoption of
782 BIM in the AEC industry. Second, the novel IE and ITr methods and algorithms proposed in this

783 work could be adopted/applied to automate a variety of other tasks in the construction domain,
784 such as contract document analysis and construction accident record analysis. Third, the
785 proposed ITr methodology could be adopted/applied outside of the construction domain, which
786 would contribute to the general domain of natural language processing/understanding. Fourth,
787 the results of this research could ultimately lead to defining principles for the drafting of future
788 regulations in a manner to support ACC. For example, the use of uncommon expressions that
789 tend to cause processing errors could be avoided when drafting future regulations.

790 **Conclusions**

791 This paper presented a rule-based, semantic NLP methodology for automated information
792 transformation (ITr) of information instances, which were automatically extracted from
793 construction regulatory documents, into logic clauses. A set of semantic mapping (SeM) rules
794 and conflict resolution rules (CoR) are used in ITr. CoR rules resolve conflicts between
795 information instances, while SeM rules transform the information instances into logic clause
796 elements. The SeM rules use context-aware and flexible information patterns. Both syntactic and
797 semantic information tags are utilized in the patterns. Syntactic information tags (e.g., POS tags)
798 are generated using NLP techniques. A semantic model helps recognize the semantic information
799 tags of each extracted information instance. A “consume and generate” mechanism is proposed
800 to handle complex sentence components and execute the SeM rules. The ITr method, thus,
801 processes almost all terms of a sentence. Such full sentence processing enables deep NLP
802 towards natural language understanding.

803 The proposed ITr algorithms were tested in transforming information instances of quantitative
804 requirements, which were automatically extracted from Chapter 19 of IBC 2009, into logic
805 clauses. The transformation results were compared with a manually-developed gold-standard.

806 The results showed 98.2%, 99.1%, and 98.6% precision, recall, and F1 measure, respectively.
807 This high performance shows that the proposed ITr methodology is promising. Through error
808 analysis, the following six causes of errors were recognized: (1) missing tags in IE; (2) incorrect
809 tagging during IE; (3) errors in processing sentences with uncommon expression structures; (4)
810 errors due to morphological analysis; (5) errors caused by certain SeM rules; and (6) structural
811 ambiguity. In future work, the authors plan to further refine the proposed methodology to avoid
812 those causes of errors – as much as possible, in an effort to further enhance the performance of
813 the ITr algorithms. Also, as part of the authors’ ongoing/future research, the proposed ITr
814 methodology will be tested on more chapters of building codes and on other types of
815 construction regulatory documents (e.g., environmental regulations). Similar high performance is
816 expected. However, variability in performance is possible due to differences in the characteristics
817 of the text across different chapters or documents.

818 **Acknowledgement**

819 The authors would like to thank the National Science Foundation (NSF). This material is based
820 upon work supported by NSF under Grant No. 1201170. Any opinions, findings, and conclusions
821 or recommendations expressed in this material are those of the authors and do not necessarily
822 reflect the views of NSF.

823 **References**

- 824 Abney, S. (1997). "Part-of-speech tagging and partial parsing." *Text, Speech and Language*
825 *Technology*, 2(1997), 118-136.
- 826 Avolve Software Corporation. (2013). *Electronic plan review for building and planning*
827 *departments*. <<http://www.avolvesoftware.com/index.php/solutions/building-departments/>>
828 (Oct 3, 2013).

- 829 Breaux, T.D., and Anton, A.I. (2008). "Analyzing regulatory rules for privacy and security
830 requirements." *IEEE Transactions on Software Eng.*, 34(1), 5-20.
- 831 Caldas, C.H., and Soibelman, L. (2003). "Automating hierarchical document classification for
832 construction management information systems." *Autom. Constr.*, 12(2003), 395-406.
- 833 Califf, M. E., and Mooney, R. J. (2003). "Bottom-up relational learning of pattern matching rules
834 for information extraction." *J. Machine Learning Research*, 4(2003), 177-210.
- 835 Cherpas, C. (1992). "Natural language processing, pragmatics, and verbal behavior." *The
836 Analysis of Verbal Behavior*, 10, 135-147.
- 837 Crowston, K., Liu, X., Allen, E., and Heckman, R. (2010). "Machine learning and rule-based
838 automated coding of qualitative data." *Proc., 73rd ASIS&T Annual Meeting: Navigating
839 Streams in an Information Ecosystem*, Association for Information Science and
840 Technology, Silver Spring, Maryland, 1-2.
- 841 Ding, L., Drogemuller, R., Rosenman, M., Marchant, D., and Gero, J. (2006). "Automating code
842 checking for building designs – DesignCheck." *Clients Driving Innovation: Moving Ideas
843 into Practice*, CRC for Construction Innovation, Brisbane, Australia, 1-16.
- 844 El-Gohary, N.M., and El-Diraby, T.E. (2010). "Domain ontology for processes in infrastructure
845 and construction." *J. Constr. Eng. Manage.*, 136(7), 730–744.
- 846 Fenves, S.J., Gaylord, E.H., and Goel, S.K. (1969). "Decision table formulation of the 1969
847 AISC specification." *Civ. Eng. Studies: Structural Research Series*, 347, University of
848 Illinois, Urbana, IL, 1-167
- 849 Garrett, J.H., Jr., and Fenves, S.J. (1987). "A knowledge-based standard processor for structural
850 component design." *Eng. with Comput.*, 2(4), 219-238.

- 851 Gildea, D., and Jurafsky, D. (2002). "Automatic labeling of semantic roles." *J. Comput. Linguist.*,
852 28(3), 245-288.
- 853 Goh, O. S., Depickere, A., Fung, C.C., and Wong, K. W. (2006). "Topdown natural language
854 query approach for embodied conversational agent." *Proc., Intl. MultiConf. Eng. and*
855 *Comput. Sci. (IMECS 2006)*, The International Association of Engineers (IAENG), Hong
856 Kong, China.
- 857 Gruber, T.R. (1995). "Toward principles for the design of ontologies used for knowledge
858 sharing." *Intl. J. Human-Computer Studies*, 43, 907-928.
- 859 Han, C.S., Kunz, J.C., and Law, K.H. (1998). "Client/server framework for online building code
860 checking." *J. Comput. Civ. Eng.*, 12(4), 181-194.
- 861 International Code Council (ICC). (2012). "International Code Council." *AEC3*,
862 <http://www.aec3.com/en/5/5_013_ICC.htm> (Oct. 26, 2013).
- 863 Khemlani, L. (2005). "CORENET e-PlanCheck: Singapore's automated code checking system."
864 *AECbytes* "Building the Future" Article,
865 <<http://www.aecbytes.com/buildingthefuture/2005/CORENETePlanCheck.html>> (Oct 26,
866 2013).
- 867 Kiyavitskaya, N., Zeni, N., Breaux, T.D., Anton, A.I., Cordy, J.R., Mich, L., and Mylopoulos, J.
868 (2008). "Automating the extraction of rights and obligations for regulatory compliance."
869 *Lecture Notes in Comput. Sci.*, 5231(2008), 154-168.
- 870 Marquez, L. (2000). "Machine learning and natural language processing." *Proc., "Aprendizaje*
871 *automatico aplicado al procesamiento del lenguaje natural"*.
- 872 Nguyen, T. (2005). "Integrating building code compliance checking into a 3D CAD system."
873 *Proc., Intl. Conf. Comput. Civ. Eng.*, ASCE, Reston, VA, 1-12.

- 874 Niemeijer, R.A., Vries, B. de, and Beetz, J. (2009). "Check-mate: automatic constraint checking
875 of IFC models." *Managing IT in Construction/Managing Construction for Tomorrow*, A
876 Dikbas, E Ergen & H Giritli (Eds.), CRC Press, London, 479-486.
- 877 Pocas Martins, J.P., and Abrantes, V. (2010). "Automated code-checking as a driver of BIM
878 adoption." *Intl. J. Housing Sci.*, 34(4), 286-294.
- 879 Pradhan, S., Ward, W., Hacıoglu, K., Martin, J.H., and Jurafsky, D. (2004). "Shallow semantic
880 parsing using support vector machines." *Proc, NAACL-HLT*, The Association for
881 Computational Linguistics, East Stroudsburg, PA, 233-240.
- 882 Roth, D., and Yih, W. (2004). "A linear programming formulation for global inference in natural
883 language tasks." *Proc., 2004 Conf. Comput. Natural Language Learning (CoNLL-2004)*,
884 SIGNLL, Boston, MA, 1-8.
- 885 Saint-Dizier, P. (1994). "Advanced logic programming for language processing." Academic
886 Press, San Diego, CA.
- 887 Salama, D., and El-Gohary, N. (2013a). "Semantic text classification for supporting automated
888 compliance checking in construction". *J. Comput. Civ. Eng.*, Accepted and published online
889 ahead of print.
- 890 Salama, D., and El-Gohary, N. (2013b). "Automated compliance checking of construction
891 operation plans using a deontology for the construction domain." *J. Comput. Civ. Eng.*,
892 27(6), 681-698.
- 893 Soysal, E., Cicekli, I., and Baykal, N. (2010). "Design and evaluation of an ontology based
894 information extraction system for radiological reports." *Comput. in Biology and Med.*,
895 40(11-12), 900-911.

- 896 Sterling, L., and Shapiro, E. (1986). "The art of Prolog: advanced programming techniques."
897 MIT Press, Cambridge, Massachusetts, London, England.
- 898 Tan, X., Hammad, A., and Fazio, P. (2010). "Automated code compliance checking for building
899 envelope design." *J. Comput. Civ. Eng.*, 24(2), 203-211.
- 900 Tierney, P.J. (2012). "A qualitative analysis framework using natural language processing and
901 graph theory." *The Intl. Review of Research in Open and Distance Learning*, 13(5).
902 University of Sheffield. (2013). "General architecture for text engineering." <<http://gate.ac.uk/>>
903 (Oct. 13, 2013).
- 904 Wyner, A., and Governatori, G. (2013). "A study on translating regulatory rules from natural
905 language to defeasible logic." *Proc., RuleML 2013: The 7th Intl. Web Rule Symposium*,
906 Springer-Verlag, Berlin Heidelberg, Germany.
- 907 Wyner, A., and Peters, W. (2011). "On rule extraction from regulations." *Proc., JURIX 2011:*
908 *The 24th Intl. Conf. Legal Knowledge and Info. Systems*, IOS Press, Amsterdam, The
909 Netherlands, 113-122.
- 910 Yin, S., and Fan, G. (2013). "Research of POS tagging rules mining algorithm." *Applied*
911 *Mechanics and Materials*, 347 – 350(2013), 2836-2840.
- 912 Zhang, J., and El-Gohary, N.M. (2013a). "Information transformation and automated reasoning
913 for automated compliance checking in construction." *Proc., 2013 ASCE Intl. Workshop*
914 *Comput. in Civ. Eng.*, ASCE, Reston, VA, 701-708.
- 915 Zhang, J., and El-Gohary, N.M. (2013b). "Handling sentence complexity in information
916 extraction for automated compliance checking in construction." *Proc., CIB W78 2013*,
917 Conseil International du Bâtiment (CIB), Rotterdam, The Netherlands.

918 Zhang, J., and El-Gohary, N. (2013c). "Semantic NLP-based information extraction from
919 construction regulatory documents for automated compliance checking." *J. Comput. Civ.*
920 *Eng.*, Accepted and published online ahead of print.

921 Zhong, B.T., Ding, L.Y., Luo, H.B., Zhou, Y., Hu, Y.Z., and Hu, H.M. (2012). "Ontology-based
922 semantic modeling of regulation constraint for automated construction quality compliance
923 checking." *Autom. Constr.*, 28, 58-70.

924 Zhou, N. (2012). "B-Prolog user's manual (version 7.7): Prolog, agent, and constraint
925 programming." Afany Software. <<http://www.probp.com/manual/manual.html>> (Nov. 19,
926 2012).

927 Zouaq, A. (2011). "An overview of shallow and deep natural language processing for ontology
928 learning." *Ontology Learning and Knowledge Discovery Using the Web: Challenges and*
929 *Recent Advances*, IGI Global., Hershey, PA, 16-38.

930

931

932

933

934

935

936

937

938

939

940

941

942 **Tables**

943 **Table 1: A Transformation Example**

Requirement Sentence	Courts shall not be less than 3 feet in width.					
Source – Information Tag	Subject	Compliance Checking Attribute	Comparative Relation	Quantity Value	Quantity Unit	Quantity Reference
Source – Information Instance	court	width	not less than	3	feet	NA
Target – Logic Clause	compliant_width_of_court(Court) :- width(Width), court(Court), has(Court,Width), greater_than_or_equal(Width,quantity(3,feet)).					

944

945 **Table 2: Patterns of the Most Applied SeM Rules in the Experiments**

SeM Rule Pattern	Action	Condition Case	Logic Clause Generated	SeM Rule Type
['a' 's' 'cr'] (a) 'OF' (b) ['a' 's' 'cr'] (c)			a(A),c(C),has(C,A)	Simple
'dpvr' (a) ['s' 'cr'] (b)	look-back search for attribute or subject (s); look-back search for negation (n)	n exists	s(S),b(B),not a(S,B)	Complex
		n not exists	s(S),b(B),a(S,B)	
'c' (a) 'v' (b)	look-back search for attribute or subject (s); look-ahead search for unit or reference (u); look-back search for negation (n)	n exists	not a(S, quantity(b,u))	Complex
		n not exists	a(S, quantity(b,u))	
'I' 's'	skip			Multiple action
'c' (a) 'v' (b) 'u' (c) 'IN' (d) 's' (e)	look-back search for attribute or subject (s)		distance(Distance),s(S),e(E), d(S,E,Distance),a(Distance, quantity(b,c))	Multiple action
['a' 's' 'cr'] (a) 'CC' (b) ['a' 's' 'cr'] (c)			(a(A);c(A))	Simple
['VB' ^ 'be'] (a) 'IN' (b) ['cr' 'a' 's'] (c)	look-back search for subject or attribute (s)		s(S),c(C),b(S,C)	Multiple action
['a' 's' 'cr'] (a) 'IN'			a(A),c(C),b(A,C)	Simple

(b) ['a' 's' 'cr'] (c)				
'Except'	mark the beginning of exception			Multiple action
'n' (a) 'c' (b) 'v' (c) 'u' (d)	look-back search for attribute or subject (s)		s(S),not b(S,quantity(c,d))	Multiple action
['a' 's'] (a) 'OF' (b) 'v' (c) ['u' 'a'] (d)		pattern preceded by ['a' 's' 'cr'] (e) ['Has' 'NoHas' 'IN' 'OF' '^ 'between'] (f)	a(A),e(E),equal_to(E, quantity(c,d))	Multiple condition
		otherwise	a(A),equal_to(A, quantity(c,d))	
'VBP' (a) 'VBN' (b)	look-back search for attribute or subject (s)		b(S)	Multiple action
I' 'CC'	skip			Multiple action
's' (a) 'MD' (b) 'Has' (c) 'a' (d)	look-back search for attribute or subject (s)	pattern preceded by 'IN'	s(S),d(D),has(S,D)	Complex
		otherwise	a(A),d(D),has(A,D)	
'TO' (a) 'VB' (b) ['s' 'cr' 'a'] (c)	look-back search for attribute or subject (s)	s not exists	c(C),a_b(C)	Complex

- 946 (1) ': A pair of single quotes encloses information tags
947 (2) '^: A caret separates optional information tags from exceptions
948 (3) (a) , (b) , (c) , etc., show the mapping of components (in SeM patterns) to logic clause
949 elements (in generated logic clauses), where an upper case represents a variable
950 (4) Contents in the "logic clause generated" column are case-sensitive
951
952
953 Table 3: Patterns of the Rest of the SeM Rules Applied in the Experiments

SeM Rule Pattern	
['a' 's' 'cr'] 'MD' 'n' 'VB' 'c' 'v' 'u'	'VBP' 'dpvr' 'VB'
's' 'JJ' 'n' 'c' 'v' 'u'	'n' 'c' 's'
'IN' 'ea' ['v' 'CD'] 'u' 'OF' 's'	['s' 'cr'] 'VBD' ['cr' 's']
'I' 'CC' 'n' 'C' 'v' 'u'	'IN' 'VBG' ['cr' 's']
'JJ' 'IN' 'c' 'v' ['u' 'cr']	['s' 'cr'] 'VBP' ['VBN' 'JJ']
'VB' 'IN' 'c' 'v' ['cr' 's']	'dpvr' 'v' 'u'
's' 'MD' 'VB' 'dpvr' ['VBZ' 'cr' 'VB']	'RB' 'TO' ['s' 'cr']
'CC' 'v' 'u' 'IN' 'a'	'MD' 'VB' 'VBN'
TO' ['s' 'cr']	'a' 'OF' 'v' 'u' 'by' 'v' 'u'

's' 'MD' 'n' 'VB' 'dpvr'	['cr' 's' 'a'] ['OF' 'IN' 'Has' 'NoHas' '^' 'for'] 's' 'IN' 's'
['s' 'a' 'cr'] 'I' 'VBG' ['cr' 'a' 's'] 'I'	'MD' 'VB' ['a' 's' 'cr']
'JJ' 'CC' 'JJR' 's'	'n' 'c' 'v'
's' 'WDT' 'VBP' 'cr'	'n' 'c' 'CD'
'VBG' 'cr' 'VBP' 'VBN'	'v' ['s' 'cr']
'MD' 'VB' 'v' 'u'	's' 'VBN'
'c' 'v' 'ea' ['cr' 's']	'JJR' 'IN'
'IN' 'JJ' 'CC' 's'	'TO' ['s' 'cs']
['s' 'cr'] 'with' 'a'	'Except' 'IN'
'n' 'c' 'v' ['cr' 's']	'rv' ['a']
'JJR' 'IN' 'v' 'u'	'VBZ' 'dpvr'
's' 'Has' 'a' 'OF' 'c' 'v' 'u'	'VB' ['cr' 'a' 's']
's' 'MD' 'VB' 'OF'	'IN' ['cr' 'a' 's']
'MD' 'VB' 'dpvr' 's'	['u' 'JJR'] ['^' 'stories']
['cr' 'a' 's'] 'MD' 'VB' ['cr' 'a' 's']	'I' 'a'
's' 'MD' 'Has' 's'	'I' 'VBD'
'cs' 'MD' 'Has' 's'	'I' 'JJ'
'v' 'u' 'CC' 'JJR'	'VBD' 'I'
's' 'MD' 'VB' 'dpvr'	

954
955

Table 4: Experimental Results Using Essential Information Tags Only

	Concepts	Relations	Total
Number of logic clause elements in gold standard	334	749	1083
Total number of logic clause elements generated	328	786	1114
Number of logic clause elements correctly generated	324	706	1030
Precision	0.988	0.898	0.925
Recall	0.970	0.943	0.951
F1 measure	0.979	0.920	0.938

956
957

Table 5: Experimental Results Using Both Essential and Secondary Information Tags

	Concepts	Relations	Total
Number of logic clause elements in gold standard	570	1349	1919
Total number of logic clause elements generated	569	1367	1936
Number of logic clause elements correctly generated	568	1333	1901
Precision	0.998	0.975	0.982
Recall	0.996	0.988	0.991
F1 measure	0.997	0.982	0.986

958
959

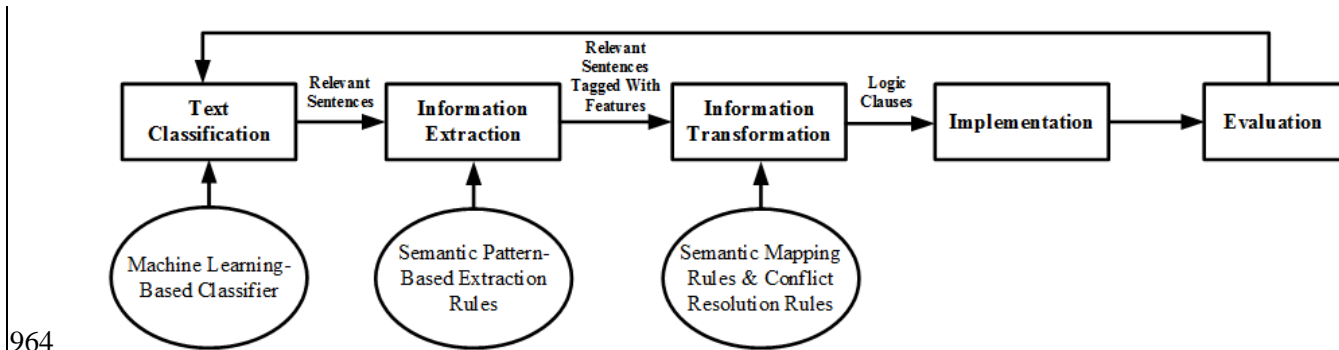
Table 6: Comparative Summary of Experiment #1 and Experiment #2

	Experiment #1	Experiment #2	Increase
Number of information tags used	8	40	+ 400%
Number of semantic mapping rules used	53	297	+ 460%
Number of conflict resolution rules used	11	9	- 18%
Number of logic clause elements built	1114	1936	174%
Precision	0.925	0.982	6%
Recall	0.951	0.991	4%
F1 Measure	0.938	0.986	5%

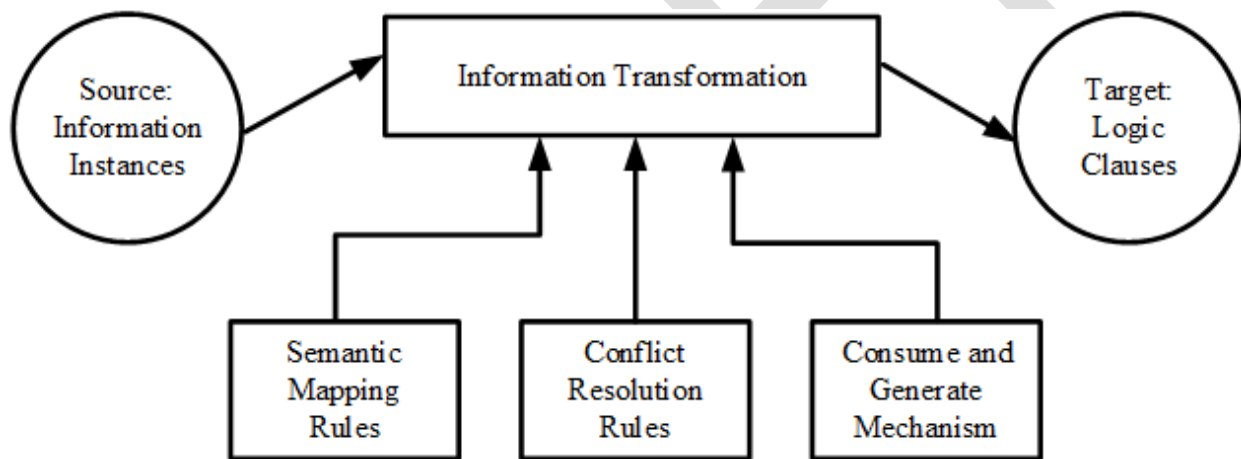
960
961

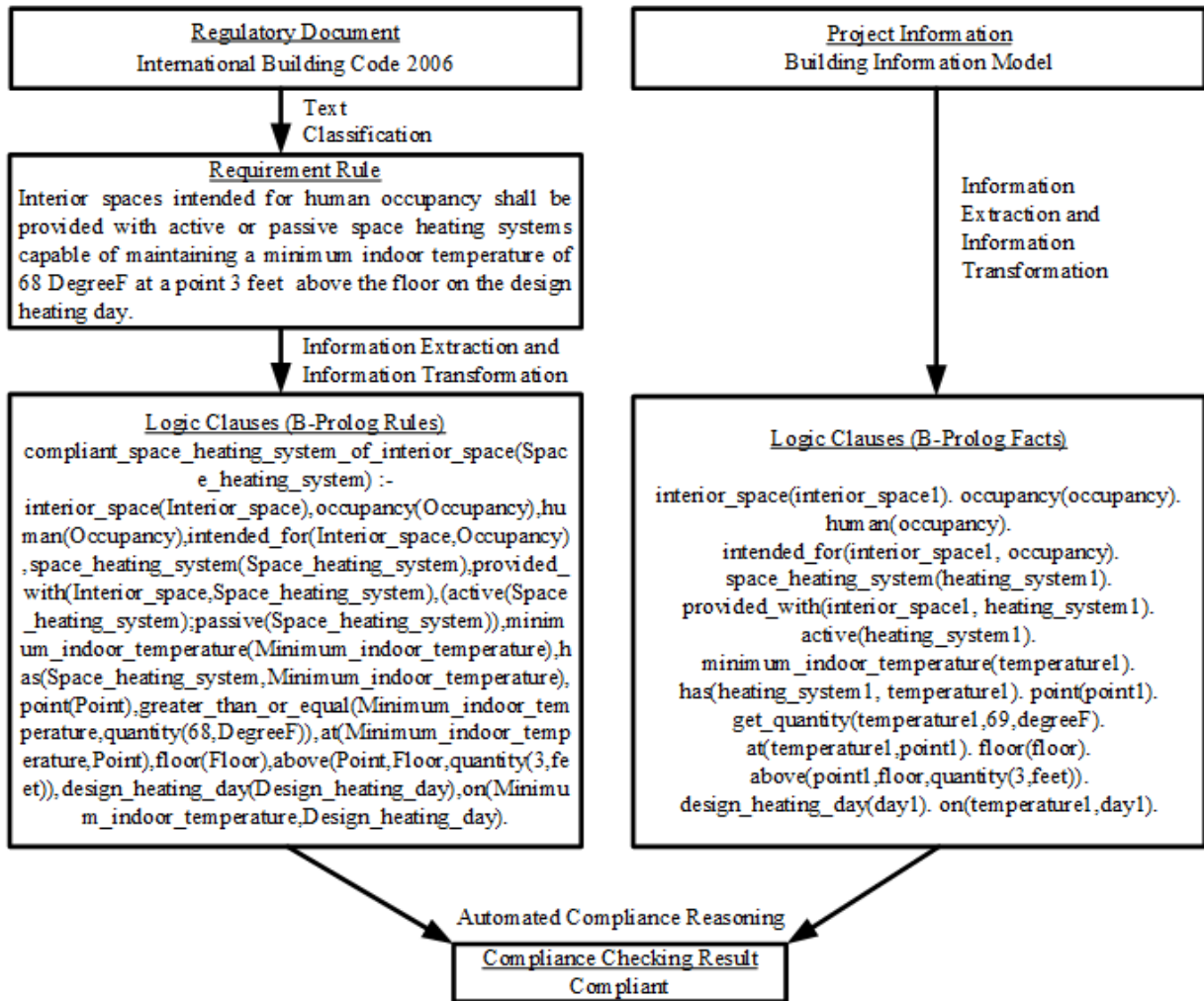
962 **Figures**

963 Figure 1. Proposed approach for automated rule extraction



965 Figure 2. Proposed information transformation methodology

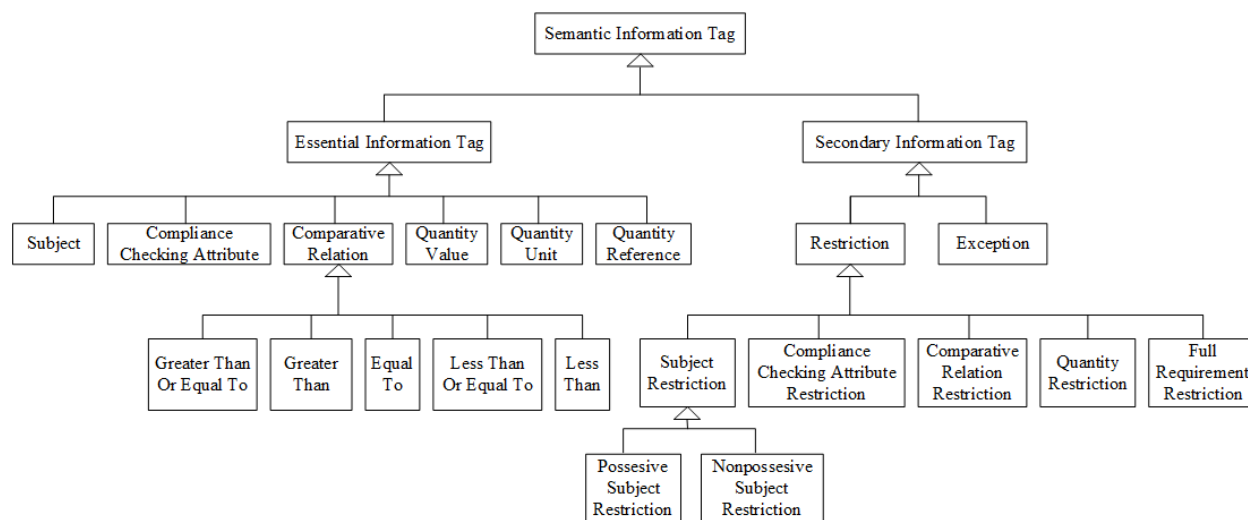




978
 979
 980
 981
 982
 983
 984
 985

986 Figure 6. Semantic information tags

987



988

989

990

991

992

993

994

995

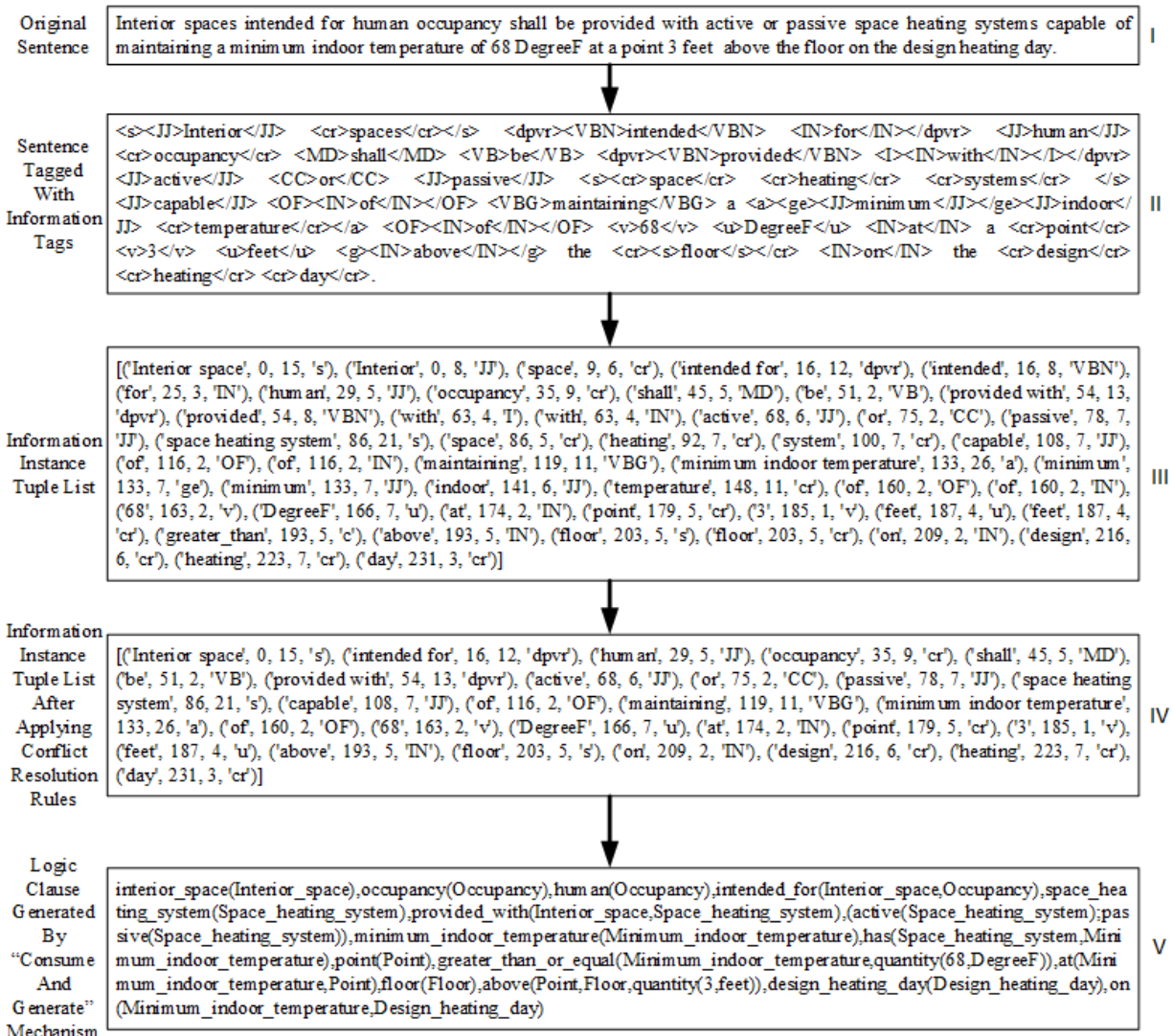
996

997

998

999

1000 Figure 7. Example illustrating the processing of a sample sentence: (a) original sentence; (b) sentence
 1001 tagged with information tags; (c) information instance tuple list; (d) information instance tuple list after
 1002 applying conflict resolution rules; (e) logic clause generated by consume and generate mechanism



1003
 1004
 1005
 1006
 1007

1008 Figure 8. Pseudocode for main algorithm

```
open file_tagged_with_information_tags

for line in file_tagged_with_information_tags:
    initialize list_for_line to []
    initialize the_logic_clause_elements_list to []
    initialize generated_logic_clause_elements_list to []
    initialize the_result_two_tuple
    initialize pointer to 0
    initialize step_length to 0
    initialize the_logic_clause to ''
    initialize list_of_touched_pointers to []
    list_of_items_in_line = collect_items_into_list(line)
    for information_instance in list_of_items_in_line:
        start_index = extract_start(information_instance, line)
        length = len(information_instance)
        tag = extract_tag(information_instance)
        append to list_for_line make_tuple(information_instance, start_index, length, tag)
    list_for_line = resolve_conflicts(list_for_line)
    the_result_two_tuple = consume_and_generate(pointer, list_for_line)
    generated_logic_clause_elements_list = first element in the_result_two_tuple
    append to list_of_touched_pointers pointer
    step_length = second element in the_result_two_tuple
    while pointer < len(list_for_line):
        the_logic_clause_elements_list.extend(generated_logic_clause_elements_list)
        the_result_two_tuple = consume_and_generate(pointer, list_for_line)
        step_length = second element in the_result_two_tuple
        generated_logic_clause_elements_list = first element in the_result_two_tuple
        append to list_of_touched_pointers pointer
        if step_length != -1:
            pointer = pointer + step_length
        else if pointer == 0:
            pointer = pointer + 1
        else if (pointer + step_length) not in list_of_touched_pointers:
            pointer = pointer + step_length
        else:
            pointer = pointer + 1
        if pointer < -len(list_for_line):
            break
    the_logic_clause_elements_list.extend(generated_logic_clause_elements_list)
    the_logic_clause_elements_list =
    remove_duplicated_elements(the_logic_clause_elements_list)
    the_logic_clause = build_logic_clause_from_elements(the_logic_clause_elements_list)
    print the_logic_clause

close file_tagged_with_information_tags
```

1009

1010 Figure 9. Pseudocode for consume and generate mechanism

```
define function consume_and_generate(pointer, tuple_list):  
    initialize result_two_tuple to [ [ ], -1 ]  
    initialize count_for_variable to 1  
    if pointer >= len(tuple_list):  
        return result_two_tuple  
    else if PATTERN1  
        look-back search and look-ahead search as needed  
        result[0].extend(Right_hand_side_of_semantic_mapping_rule1)  
        result[1]=length of PATTERN1  
    else if PATTERN2  
        look-back search and look-ahead search as needed  
        result[0].extend(Right_hand_side_of_semantic_mapping_rule2)  
        result[1]=length of PATTERN2  
    else if PATTERN3  
        look-back search and look-ahead search as needed  
        result[0].extend(Right_hand_side_of_semantic_mapping_rule3)  
        result[1]=length of PATTERN3  
  
    ...  
  
    return result_two_tuple
```

1011

1012

1013