For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

# Regulatory Information Transformation Ruleset Expansion to Support Automated Building Code Compliance Checking

3

Xiaorui Xue, S.M.ASCE<sup>1</sup>; Jiansong Zhang, Ph.D., A.M.ASCE<sup>2</sup>

## 4 Abstract

5 The traditional building code compliance checking process mainly relies on design reviewers to review design documents or models manually. The intensive manual effort needed makes this 6 7 process time-consuming, costly, and error-prone. Automated compliance checking (ACC) could be a promising upgrade of the traditional manual code compliance checking. With a reduced workload 8 9 on design reviewers, ACC is cheaper, faster, and immune to human errors. To support ACC, building 10 code requirements need to be represented in a computer-processable format to enable automated 11 reasoning, which will in turn allow an automated assessment of the building design's compliance status with building codes. A major limitation of many existing ACC systems/methods is their 12 limited range of checkable building code requirements. To address that, the state of the art uses 13 14 pattern matching-based rules to transform building code requirements into computable formats 15 automatically, but the ruleset was developed and tested only on a few chapters of building code 16 requirements. An efficient ruleset expansion method is needed to increase its range of checkable building code requirements at a low cost to bring ACC systems closer to full deployment. In this 17 paper, the authors proposed a new regulatory information transformation ruleset expansion method 18 19 for expanding an existing rules et. This method can expand the range of checkable code requirements of ACC systems without significant manual effort. The proposed ruleset expansion method takes an 20 iterative approach to ensure the generality and validity of new pattern matching-based rules and the 21 22 quality of information transformation results. The expanded ruleset was tested on generating logic 23 clauses from Chapter 5 of the International Building Code 2015. Compared to the baseline ruleset, the expanded rules et increased the predicate-level precision, recall, and F1-score of the logic clause 24 generation by 10.44%, 25.72%, and 18.02%, to 95.17%, 96.60%, and 95.88%, respectively. 25

## 26 **1. Introduction**

A broad range of building codes govern the architecture, engineering, and construction (AEC) industry. There are building codes that are applicable to an entire building, such as the International Building Code (IBC) [1] and the International Fire Code (IFC) [2], or building codes that are applicable to a part/component of a building, such as the American Society of Heating, Refrigerating and Air-Conditioning Engineers (ASHRAE) Standard 90.1 and the International Energy

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

Conservation Code (IECC). The number of existing building codes is large. For example, the State 32 33 of Indiana enforces over 50 different construction-related codes, standards and other regulations [3]. 34 The large amount of existing building codes makes traditional manual code compliance checking 35 suffer from three major limitations: (1) a long review process, (2) high review cost, and (3) tendency 36 to have human errors and omissions in review results [4-6]. A traditional code compliance checking 37 process involves multiple review cycles and interactions between design reviewers and building 38 designers. In each design review cycle, building designers submit their designs to design reviewers. 39 If the designs have omissions or mistakes, design reviewers will return them to the building designers and ask for additional information or revisions. The interaction between building 40 41 designers and design reviewers continues until the design reviewers are confident that there is no omission or mistake in the design anymore, at which point a building permit will be issued. 42 43 Traditional code compliance checking is a long and expensive process. For example, the average turnaround time for issuing a building permit is 73 days and the minimum cost is 302 dollars in 44 45 Chicago [7,8]. A potential solution to address the limitations of traditional manual code compliance checking is 46 automated compliance checking (ACC). In ACC, software and algorithms reduce the workload of 47 48 design reviewers by automatically checking building designs against building codes. By making the 49 job of design reviewers easier, ACC is expected to cost less money, consume less time, and generate 50 fewer errors than the traditional code compliance checking approach [9]. The potential benefits of

- ACC also include: (1) facilitating the adoption of building information modeling (BIM) in the AEC industry [10,11], (2) promoting a common information exchange format in the AEC industry [12], and (3) encouraging the use of computing techniques such as natural language processing (NLP) and logic reasoning in the AEC industry [13]. Given all the benefits, the exploration of ACC dates
- back to the 1960s, when Fenves [14] developed a decision table system to check the design of steel structures against the American Institute of Steel Construction (AISC) specifications. The success of Fenves paved the road to more complex ACC systems. For example, Garrett and Fenves [15] proposed the use of information networks to represent design standards in the ACC systems. At the same time, the boom of expert systems gave birth to expert ACC systems, such as the Standard Interface for Computer-Aided Design (SICAD) [16], the Standards Processing Expert (SPEX) [17],
- 61 and the Design Prototype system [18].

Prior attempts at expert ACC systems ceased in the 1990s due to their high maintenance cost, but 62 63 more complex expert ACC systems, such as the BCA ider and the Design Check [19], emerged in the 64 early 2000s [20]. In addition to expert ACC systems that aim to check general building codes, some expert ACC systems focus on building codes in a specific domain, such as (1) the Fire-Code 65 66 Analyzer that checks the compliance of fire protection codes in New Zealand [17], (2) the Life 67 Safety Code Advisor that checks against the National Fire Protection Association (NFPA) safety 68 code in the United States, and (3) the TALLEX that checks compliance of tall buildings in the United 69 Arab Emirates (UAE) [21]. The introduction of BIM as a digital representation of buildings led to 70 the advancement of BIM-based ACC systems. Solibri Model Checker (SMC) started as a BIM validation tool, but automated code compliance checking plugins soon evolved [22]. The 71

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

72 DesignCheck checks BIM of a building against accessibility regulations [19]. The Construction and 73 Real Estate Network (CORENET) project by the Singapore government checks BIM, instead of 74 paper-based design documents, against barrier-free access codes and fire codes. The KBimCode can 75 check BIM against South Korea's building codes [23].

76 Despite all the development and advancement of ACC systems, existing ACC systems still rely on domain experts to extract building code requirements and formalize them into a computer-77 78 processable format [24,25], such as decision tables [9], knowledge models [26], or structured 79 rulesets [27]. The high cost and long duration involved in this manual process prevent a comprehensive testing of the regulatory information coverage of ACC systems. Many existing ACC 80 systems only focus on a specific range of building code provisions for this reason. For example, 81 Nguyen [28] proposed an ACC computing framework by manually converting two chapters of the 82 83 International Building Code (IBC) to a hierarchical data structure. Nawari [29] manually translated a portion of the National Green Building Code Standard (ICC 700-2008) to the extensible markup 84 85 language (XML) format. These efforts in academia successfully proved the concepts of many promising ACC systems, but their limited range of checkable building code requirements prevented 86 their full-fledged industry adoption. The large amount of manual effort required to convert building 87 codes in natural language to a computer-processable structured format hindered the full deployment 88 of ACC systems in the AEC industry. Vendors of commercial ACC systems believed that they could 89 convert enough building code requirements and achieve automated code compliance checking by 90 91 hiring domain experts. However, the high cost of manual conversion led to this idea tested uneconomic, and most related attempts ceased in the 1990s. 92

To enhance the practicality of ACC systems, the extraction and transformation of building code 93 requirements from a wider range of sources need to be automated. Fully automated processing of a 94 95 wide range of building code requirements is a key functionality of envisioned future ACC systems. The authors proposed a ruleset expansion method to expand the range of checkable building code 96 97 requirements of ACC systems that use a pattern matching-based regulatory information transformation ruleset to automatically transform building codes from natural language to 98 99 computable logic clauses that support automated reasoning. This type of system [30], although fullyautomated, has a limited range of checkable building code requirements in its current 100 implementations. The existing automated building code transformation method [31] was established 101 based on the assumption that a set of common patterns exist that could be used to analyze the various 102 103 information elements encoded in building code requirements and transform them to an unambiguous 104 and computable representation. It was not clear, however, how large this set of common patterns should be. As a result, based on this assumption, although theoretically there exists a "superset" of 105 common patterns, in practice, it is likely that we will need to adopt a data-driven approach to 106 107 accumulatively grow an existing set of common patterns to asymptotically approach that theoretical 108 "superset." In line with that, the existing set of common patterns (or ruleset of pattern matching-109 based rules) needs to be continuously expanded to cover otherwise missed regulatory information when they are applied to unfamiliar building codes. More pattern matching-based rules are needed 110 to bring them to practice in the AEC industry. 111

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

In this paper, the authors proposed a rules et expansion method to expand the range of building code 112 113 requirements that ACC systems can extract and transform in an efficient manner. The ruleset 114 expansion method needs to control two aspects: (1) the validity of the ruleset, and (2) the generality 115 of the ruleset. The proposed method therefore takes an iterative approach to pursue a balanced validity and generality of added rules while maintaining the compatibility of added rules with 116 existing rules. The proposed method is designed to support ACC systems that check 3D digital BIM 117 118 models of building designs. The design information extraction module of the ACC systems extracts 119 building design information through a series of information extraction and transformation 120 algorithms [32]. Then, the extracted building design information and automatically generated logic clauses are sent to the compliance checking reasoning module to generate compliance checking 121 results [25,33,34]. Methods that extract building design information from BIM or IFC files have 122 123 been well developed [35,36]. While the proposed method in this paper focuses on producing logic 124 clauses that work with BIMs, it can potentially be used with building plans and drawings as well. 125 The extraction of building design information from building plans and drawings could be a 126 promising direction for future research.

## 127 **2. Background**

# 128 2.1 Natural Language Processing

Chowdhury defines natural language processing (NLP) as "an area of research and application that 129 explores how computers can be used to understand and manipulate natural language text or speech 130 to do useful things [37]." NLP includes a wide range of tasks, such as (1) information retrieval [38], 131 (2) information extraction [39], (3) text classification [40], (4) text generation [41], (5) text 132 summarization [42], (6) question answering [43], (7) machine translation [44], and (8) speech 133 recognition [45]. There are two main approaches to accomplishing NLP tasks: the rule-based 134 approach and the machine learning-based approach [46]. Rule-based NLP systems may require 135 136 manual effort in rule generation, but usually outperform machine learning-based NLP systems in a specific task or in a specific domain [47]. Machine learning-based NLP systems can be further 137 classified into "shallow" learning systems and "deep" learning systems based on the types of 138 139 machine learning models they use. "Shallow" learning systems use traditional machine learning algorithms, such as support vector machines (SVMs) or decision trees, and require manual feature 140 141 engineering. Deep learning systems use neural networks and do not require manual feature 142 engineering [48]. There is no lack of efforts to use NLP in the AEC domain. For example, Tixier et 143 al. [49] used NLP to extract the reasons for accidents from construction injury reports. Lin et al. [50] 144 used NLP technologies to extract information from BIM. ACC research also uses NLP techniques to process building codes, for matching between concepts in building codes and concepts in BIM 145 [51], and for converting building codes to logic clauses that support automated reasoning [52]. 146

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

## 147 2.2 Part-of-Speech

148 Part-of-speech (POS) of a word represents its lexical and syntactic function in a sentence [43]. 149 English words have eight basic POS categories: (1) noun, (2) verb, (3) adjective, (4) adverb, (5) 150 pronoun, (6) preposition, (7) conjunction, and (8) interjection [54]. The same word may have different POS categories in different contexts. For example, the word "run" could be a verb in its 151 152 simple present tense or past perfect tense depending on the context. In NLP systems, words are categorized into more specific POS categories to represent text more informatively. For example, 153 the Penn Treebank Corpus classifies words into 36 POS categories [55] and the Brown corpus has 154 179 POS categories [56]. In the development of the Penn Treebank Corpus and the Brown Corpus 155 156 above, human annotators manually assigned words to different POS categories according to their 157 understanding of the English language and the contexts of the words. POS tagging software, which 158 is commonly called "POS taggers," could replace annotators' manual effort in this task. POS taggers 159 automatically determine the POS category of a word using its contextual information in an algorithmic manner [57]. POS taggers began with rule-based taggers that used a set of rules to 160 161 determine the POS categories of words. These rules can be compiled by experts [58] or extracted 162 from text algorithmically [53]. With the development and integration of machine learning, POS taggers shifted to the use of statistical models. For example, Giménez and Marquez [59] used one 163 SVMs model to determine POS categories of known words and another SVMs model to predict 164 those of unknown words. Brants [60] developed a POS tagger which uses Hidden Markov Modek 165 166 (HMM) to capture dependencies among words and determine the POS categories of words by their inter-dependencies. Plank et al. [61] proposed the use of bi-directional neural networks to 167 168 accomplish multilingual POS tagging. POS tagging is an important early step of many NLP systems 169 [59].

# 170 2.3 Ontology

Ontology is the explicit and formal description of knowledge through relationships among concepts 171 in a domain [62]. In 1999, the World Wide Web Consortium (W3C) first developed the Resource 172 173 Description Framework (RDF) language for ontology [63]. Then, it collaborated with the Defense 174 Advanced Research Projects Agency (DARPA) to extend the RDF into a more expressive DARPA 175 Agent Markup Language (DAML) [64] [65]. After that, many ontologies emerged, either for a specific domain (e.g., medical) [66] or for general-purpose [67]. Ontology is used to: (1) analyze 176 177 and reuse domain knowledge, (2) share structured domain knowledge among people and software, 178 (3) specify domain assumptions, and (4) distinguish domain knowledge from operational knowledge 179 [68].

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

## 180 2.4 Text Similarity Measurements

181 Text similarity is an important benchmark in NLP. There are many ways to measure the similarity 182 between text strings [69]. Text similarity can be measured by comparing words or characters in text 183 strings. For example, the Levenshtein Distance [70] [71] measures the minimum number of single character transformations needed to convert one string to another. In Levenshtein Distance, 0 means 184 185 two strings are identical, and the larger it is, the less similar the two strings are, with no strict upper bound. The Jaccard Distance, on the other hand, measures the number of items shared by two sets 186 187 [72]. In the Jaccard distance, 0 means two sets are identical, and 1 means two sets share no common items. The Jaro Winkler Similarity [73] is an extension of the Levenshtein Distance. By normalizing 188 189 the Levenshtein Distance with the length of the text string, the Jaro Winkler Similarity ranges from 190 0 to 1. In the Jaro Winkler Similarity, 0 means two text strings are completely different and 1 means two text strings are the same. 191

192 The inability to measure similarity between word meanings is one limitation of measuring text 193 similarity at the word and character levels. One potential solution to the problem is representing 194 words (and their contexts) as vectors in high-dimensional spaces. Popular text vectorization 195 techniques include, for example, Word2Vec [74], FastText [75], and Glove [76]. The distance 196 between meanings of the two words and their contexts can be measured by the cosine distance 197 between the two vectors.

## 198 **3. Methodology**

The proposed method expands the range of processable building code requirements by adding new 199 pattern matching-based rules to an existing ruleset. The pattern matching-based rules capture 200 regulatory information in building codes and convert the captured information to logic clauses. The 201 pattern matching-based rules consider both syntactic information, which is provided by POS tags 202 (e.g., the word "height" in the phrase "building height" is a noun because it has a POS tag "NN"), 203 204 and sematic information, which is provided by an ontology (e.g., the phrase "less than" after an 205 attribute and before a value means that the attribute value must be smaller than the specified value, 206 and the phrase "minimum clearance" means the attribute "clearance" must be greater than or equal 207 to a specified value). For example, the pattern "subject, conjunction, subject" can be used to extract the regulatory information of which two subjects are in equivalent status. The conjunction (i.e., and, 208 or) is a label by the POS tagger. Subjects, which were labeled as noun by the POS tagger, were 209 210 further labeled as subjects after feature enhancement by the ontology. The pattern matching -based 211 rules regulate how this method extracts building code requirements and converts them to logic 212 causes. The logic clauses represent building code requirements in a strict horn clause (HC) format 213 in B-Prolog syntaxto avoid ambiguity in natural language and facilitate automated reasoning by the 214 logic reasoner.

Each logic clause has a left-hand side and a right-hand side, separated by the delimiter ":-". The

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

left-hand side is the head of the logic clause that represents the subject of compliance checking in 216 217 the logic clause, i.e., the building design component that this code requirement governs. The subject 218 of compliance can be an entire building, a component of a building, a certain attribute of a building, 219 or an attribute of a building component. The predicates on the right-hand side of the delimiter ":-" 220 (i.e., in the body of the logic clause) are conditions that the subject of compliance need to meet to comply with the building code requirement. This logic clause indicates that if all predicates on the 221 right-hand side of the delimiter ":-" are evaluated to True, then the predicate on the left-hand side 222 of the delimiter ":-" will also be evaluated to True. In the context of this study, it means that if the 223 224 subject of compliance meets all the conditions of the corresponding building code requirement, it is then considered to be compliant with the building code requirement. In the logic clauses, the 225 conjunction relation (i.e., AND) is represented as a comma "," and the disjunction relation (i.e., OR) 226 227 is represented as a semicolon ";". One manually generated logic clause example as part of the gold standard (see details in Section 4.2 228 229 - Gold Standard Generation) is provided in Figure 1. The "Travel distance" in Figure 1 is the subject of compliance and the predicates on the right-hand side describe the building code requirement that 230 the "Travel distance" need to comply with. Each predicate describes one condition that the subject 231 needs to satisfy to comply with in the building code requirement described in the logic clause. If all 232 of the predicates on the right-hand side are evaluated to true, the ACC system will then determine 233 234 the building design to be in compliance with the corresponding building code requirement. More 235 specifically, the predicates "from(Travel distance, Accessible space), to(Travel distance, Area of refugee)" describe that the travel distance is measured from the accessible space to the 236 refugee area. The predicate "in accordance with (Travel distance 2, section 1017 1)" describes 237 that the travel distance is specified in Section 1017.1 of the IBC 2015. The predicates "not 238 239 greater than(Travel distance, Travel distance 2)" require the travel distance from the accessible space to the area of refugee to be no greater than the travel distance specified in Section 1017.1 of 240 241 the IBC 2015. Other predicates in the logic clause are required by the strict HC format in B-Prolog syntax for this logic rule to execute. Overall, this logic clause describes the building code 242 requirement that the maximum travel distance from an accessible space to an area of refugee should 243 244 not be greater than the travel distance specified in Section 1017.1 of the IBC 2015.



For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

247

253

254 255

In the manual transformation of building code, domain experts complete the transformation based on their understanding of building code requirements. In our automated transformation, a pattern matching-based regulatory information transformation ruleset is used to complete this transformation automatically. To support the matching pattern development in the ruleset, building codes undergo feature enhancement by POS tagging and ontology matching (Figure 2).



The goal of this research is to develop an efficient and effective method to extend an existing pattern 256 matching-based regulatory information transformation ruleset. Although it is possible to develop a 257 new ruleset from scratch, the authors chose to expand an existing ruleset developed by Zhang and 258 El-Gohary [31], based on the assumption that asymptotic full coverage of building codes could be 259 achieved by expanding an existing ruleset. In addition, expanding an existing ruleset, instead of 260 generating new rulesets, has the benefits of (1) reducing rule generation workload, and (2) allowing 261 the expanded ruleset to capture patterns absent in the training dataset, while maintaining the 262 compatibility of the expanded rules et with the automated building code compliance checking system 263 264 The expansion of an existing rules trequires new rules to be added. The added rules must meet 265 certain standards or have certain characteristics to realize the benefit goals above. For example, the amount of effort/time spent on new rules development should be significantly less than (e.g., no 266 greater than 20% of) the development of the original ruleset. To achieve the first goal, the number 267 of added rules should be small. To achieve the second goal, the process of adding new rules needs 268 269 to be selective. The added rules should be valid and general. A rule is valid when it correctly extracts the regulatory information it is designed to extract. A rule is considered general when it has been 270 271 applied at least twice in the training dataset. The combination of multiple valid and simple pattern 272 matching-based rules can be used to represent more complex patterns in building codes. The added 273 rules also need to be general to capture patterns that are not in the training data. Building codes are legal documents composed by a panel of experts following strict guidelines. Therefore, the same 274 275 patterns may be shared by different chapters of the building code. The generality requirement allows 276 pattern matching-based rules to capture common patterns shared by different chapters of building 277 codes, or different building codes. Different building codes, or at least different chapters in one 278 building code, should follow a set of common patterns, according to English grammar and the way 279 building codes were compiled. The ruleset expansion method proposed in this research is designed

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

to ensure the generality and validity of added rules. The transformation rule generation and validation are manually conducted in the proposed method. However, once the transformation rules are generated and validated, they can be used to automatically transform building code requirements into logic clauses.

## 284 3.1 Ruleset Expansion Method

The proposed rules expansion method takes an iterative approach to add new rules into an existing 285 286 regulatory information transformation ruleset. The goal of this research is to develop an efficient 287 method to expand an existing pattern matching-based regulatory information transformation ruleset and ensure the generality and expandability of the added rules. To achieve this goal, the authors 288 should identify missing regulatory information and generate new rules to capture it. There are two 289 approaches to completing this task: (1) identify all missed regulatory information and generate 290 corresponding rules in a single pass, or (2) identify one piece of missed regulatory information at a 291 292 time, generate a rule to capture the missed information, review the performance of the new rule. 293 modify the new rule, and then proceed to the next iteration of identifying missed information. Both 294 approaches can generate a ruleset that captures all regulatory information. However, the first 295 approach does not consider the validity and generality of the added rules and the interaction among 296 them (i.e., multiple rules may match the same regulatory information). The second approach iteratively adds new rules and tests them before they are eventually added to the ruleset. The second 297 298 approach has the potential to generate more valid and general rules than the first approach. What 299 really distinguishes the first and second approaches is the granularity of pattern matching-based 300 rules. The first approach aims to extract regulatory information at the chapter level or even at the whole building code level. Whereas the second approach extracts regulatory information at the 301 sentence level. Because logic clauses are generated at the sentence level, the second approach 302 naturally fits better. In addition, the shorter the pattern lengths are, the more flexible they become 303 304 and the better scalability the whole ruleset will have. Patterns may match the whole sentence of a 305 building code requirement or (most likely) part of a sentence. Data-driven expansion of the existing ruleset is also made possible through the second approach, whereas it would not have been possible 306 with the first approach. Previous rule-based NLP applications [58,77,78] with manually developed 307 rules also supported this point. Testing rules one-by-one before they are added to the rules et is also 308 309 a more rigorous rule development process than generating all the rules and testing them together. 310 Because of the above-mentioned advantages, the rules et expansion method proposed in this research 311 takes the second approach, which is shown in Figure 3. One candidate pattern matching-based rule is generated to capture missing regulatory information one piece at a time, until all the missed 312 regulatory information in the training dataset is captured. A version of logic clauses is first generated 313 314 by the ruleset to identify missing regulatory information. If an instance of missing regulatory information is identified in this version of logic clauses, a candidate pattern matching-based rule is 315 316 generated to extract it. The candidate pattern matching-based rule will be added to the ruleset if it is

317 proven to be general and valid. The generality and validity of the candidate rule are tested by

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

inspecting a new version of logic clauses generated by the ruleset when the candidate rule is included. 318 319 A valid rule must correctly extract the information it is designed to extract and does not introduce 320 errors when it is applied to other parts of the training text. A general rule needs to be able to be 321 applied at least twice in the training dataset, which means it should be applied at least once outside 322 of the sentence it is extracted from. The validity of a rule has higher priority than its generality. If 323 the rule introduces any errors in the training text, it will be modified until it introduces no errors. At 324 the same time, the validity of rules will not be sacrificed to make a new rule general. It is possible 325 that a pattern appears only once in the training text, and it is still necessary to capture an instance of 326 regulatory information. The rules et expansion process continues until the expanded rules et captures all the regulatory information in the training data. 327



329 330

328

Figure 3. Ruleset Expansion Method.

## 331 3.2 Feature Enhancement

The building codes are first enhanced by POS tagging and ontology matching, to generate extra 332 features. The enhanced building codes are more informative and support more complex operations 333 334 than the original building codes without extra features. Building codes are labeled with information 335 tags in this step. Extra features can provide more information about building code expression 336 patterns and, therefore, increase the performance of pattern matching-based rules. To understand building codes, it requires knowledge both of the English language and of the AEC domain. The 337 feature enhancement makes the ACC system better at processing building codes by introducing such 338 knowledge to the ACC system. The authors apply POS tagging to generate syntactic features (i.e., 339 340 for knowledge of the English language) and apply ontology matching to introduce AEC domain knowledge in this step. 341

The ACC systemuses POS tagging, which captures the grammatical roles of words in a sentence, to generate syntactic features from building code. The same word in different POS categories can have distinct meanings. For example, when the word "run" is a verb, it means an action of moving through a space. When the word "run" is a noun, it refers to a physical object. Therefore, syntactic features together with semantic features can disambiguate words in building codes. For example, when the word "runs" is a noun in the sentence "The extensions of handrails shall be in the same direction of the flights of stairs at stairways and the ramp runs at ramps" [1] (Section 1014.6 of IBC

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

2015), it represents a physical object with attributes governed by the building code. However, when 349 350 the word "runs" is a verb in the sentence "Where a partition containing piping runs parallel to the 351 floor joists" [1] (Section 2308.5.8 of IBC 2015), such a possibility can be ruled out. In this research, 352 the authors used the A Nearly-New Information Extraction System (ANNIE) POS tagger in the General Architecture for Text Engineering (GATE) [80] with proven performance in tagging 353 building codes to generate accurate syntactic features. Such external POS taggers, which were 354 355 trained on a larger body of text and fine-tuned by domain experts, can bring additional grammatical 356 knowledge to the ACC system.

- The ACC system also uses an ontology to introduce AEC domain knowledge for logic clause 357 generation. In manual code compliance checking, reviewers already have domain knowledge 358 needed to understand building codes, based on their education, training, and experience. However, 359 360 in ACC systems, such knowledge needs to be explicitly provided. An ontology allows the ACC 361 system to access domain knowledge and consider domain knowledge in rule generation. For 362 example, with an ontology, the method can treat "International Fire Code" and "automatic sprinkler system" as integral phrases instead of multiple individual words. It also makes the method treat 363 "inches" and "feet" as units specifying a numerical constraint, instead of regular nouns. In addition, 364 365 the ontology also supports the disambiguation of vague terms.
- The used ontology has two main types of items: (1) essential information, and (2) secondary 366 367 information. Essential information includes: (1) subject of a particular regulatory requirement (e.g., 368 building), (2) attribute (e.g., building height), (3) comparative relationship (less than, greater than), 369 (4) quantity (e.g., value or range of value), (5) quantity unit (e.g., inch, feet), and (6) reference to other quantity. Secondary information includes restrictions and exceptions. Restrictions mean 370 constraints to subjects and attributes [26] [31]. For example, in the sentence "Exterior exit stairways 371 372 and ramps serving as an element of a required means of egress shall be open on not less than one side, except for required structural columns, beams, handrails and guards [1]," "serving as an 373 374 element of a required means of egress" is a constraint to "Exterior exit stairways and ramps." 375 Exceptions are the conditions where a requirement does not apply. The ontology was created and tested in a previous study [31]. With an expansion for this specific task, its comprehensiveness is 376 377 ensured in the context of this application by enumerating all covered concepts in the corresponding 378 code requirements. The ontology is also scalable. Similar to the ruleset itself, the ontology could also be accumulatively and continuously developed to fulfill the need for processing different 379 380 building codes, until it reaches or asymptotically approaches the saturated state where any 381 potentially related concept to building codes is included. It is editable in GATE [80] or using a plain text editor. Ontology editing tools provide support for the scalability of the ontology as the size and 382 complexity of the ontology increases. 383

## 384 3.3 Pattern Extraction

There are two approaches to extracting regulatory information from building codes: the top-down approach, and the bottom-up approach. In the top-down approach, the information extraction

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

algorithm constructs a global logic clause framework that matches the overall structure of a sentence 387 388 and fills in the building code requirements into the framework. In the bottom-up approach, the 389 information extraction algorithm captures building code requirements at a local level and assembles 390 them into a logic clause. The building codes consist of a lot of long and complex sentences with 391 diverse structures. The versatility and complexity of building code sentences means developing enough sentence-level frameworks to accommodate all sentences in building codes may require a 392 393 similar amount of manual effort as directly converting building codes to logic clauses. It is possible 394 that every sentence, or at least every few sentences, requires a different framework. The authors 395 propose the use of the bottom-up approach. The complex and versatile structures of building code sentences may require many complex sentence-level frameworks, but pattern matching-based rules 396 can assemble these structures from simple local patterns. 397

## 398 **4. Experiment**

## 399 4.1 Ruleset Expansion Experiment

The authors tested the effectiveness of the proposed ruleset expansion method by measuring its 400 401 precision, recall, and F1-score in logic clause generation, which in turn tested the generality of the 402 original ruleset. Chapter 10 of the International Building Code 2015 (IBC 2015) was selected as the 403 training data for the experiment and Chapter 5 of the IBC 2015 was selected as the testing data. Zhang and El-Gohary developed the original rules et based on Chapters 12 and 23 of IBC 2006 [31]. 404 405 The authors used the ruleset expansion method to generate new rules based on Chapter 10 of the IBC 2015 and tested the expanded ruleset on Chapter 5 of the IBC 2015, in comparison with the 406 407 original ruleset.

In the first part of the experiment, the original ruleset generated a baseline version of logic clauses from the training data. The training data was first pre-processed by a POS tagger and an ontology to generate enhanced features. The POS tagger used in this research is the ANNIE tagger from the GATE tool [79]. The authors used the ontology developed in [31] with expansions on Chapters 5 and 10 of the International Building Code 2015. After that, the authors used the ruleset expansion method to expand the original ruleset.

First, the authors identified missing regulatory information and updated the original ruleset with a 414 candidate pattern matching-based rule to capture the missing regulatory information. For example, 415 416 the original ruleset did not have enough patterns to extract all the essential information for requirements that are described using negation together with a past participle verb, so a 417 418 corresponding pattern and candidate rule was added. The expanded rules et also includes all rules in 419 the original ruleset. The authors added 64 new rules to the original ruleset, a much smaller number 420 compared to the 306 rules already in the original ruleset. Two of the 64 new rules were developed 421 to extract missed regulatory information in the example type mentioned above. One rule with the 422 pattern "modal verb, negation, base form verb, [adjective, past participle verb, past tense verb]" was

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

generated to extract the regulatory requirement of a subject. This rule can process building code 423 424 requirement sentences like "A basement (candidate subject) provided with one exit shall (modal 425 verb) not (negation) be (base form verb) located (past participle verb) more than one story below grade plane" (Section 1006.3.2.2 of IBC 2015) [1], and "The area of a Group F-2 or S-2 building 426 427 (candidate subject) no more than one story in height shall (modal verb) not (negation) be (base form verb) limited (past participle verb) where the building is surrounded and adjoined by public ways 428 429 or yards not less than 60 feet in width" (Section 507.3 of IBC 2015) [1]. The first subject is extracted 430 by identifying the first subject candidate to the left of (not necessarily immediately next to) the relationship. This arrangement makes pattern matching flexible. The rule is both general and valid, 431 because it was applied twice in the training dataset and correctly extracted the regulatory 432 information it was designed to extract. Another pattern "candidate subject, preposition, comparative 433 434 relation, value, unit" was generated to extract the quantitative regulatory requirement of a subject. This rule can process building code requirement sentences like, "The ladder or steps shall not 435 436 encroach into the required dimensions of the window well (candidate subject) by (preposition) more 437 than (comparative relation) 6 (value) inches (unit)." (Section 1030.5.2 of IBC 2015) [1]. Only the "of" relationship between "the required dimension" and "the window well" was extracted by the 438 original ruleset. The newly added rule was not general in the current training dataset (i.e., it was 439 applied only once in the training dataset), but it was still valid, because it correctly extracted the 440 regulatory information it was designed to extract. Although the rule was not general, it was still 441 442 needed to capture an instance of regulatory in formation. Therefore, it was still added to the ruleset. The complete set of the 64 rules can be found in Appendix A. In the second part of the experiment, 443 the expanded ruleset was tested on Chapter 5 of the IBC 2015 to automatically convert building 444 445 codes to logic clauses. The automatically generated logical clauses were compared against a gold standard. 446

447 4.2 Gold Standard Generation

Chapters 10 and 5 of the IBC 2015 were transformed into logic clauses by three annotators semi-448 automatically to create a gold standard of information transformation and logic clause generation. 449 All three annotators have background AEC knowledge to understand building codes, and necessary 450 skills to transform building codes into logic clauses. The authors provided the annotators a clear 451 452 annotation protocol, a brief training section before annotation, and machine-generated logic clauses as reference during their annotation. They worked independently without access to the logic clauses 453 454 generated by other annotators. However, they were presented with the machine-generated logic clauses, which could help annotators align with the rule generation mechanism of pattern matching-455 based rules, achieve higher inter-annotator agreement, and reduce rule generation time. It also 456 457 ensures the compatibility of human-generated logic clauses with the automated code compliance checking system. 458

Annotators were required to use the exact words that came from the building code in their generated
 logic clauses. For example, if the building code uses the word "exterior" for exterior walls,

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

461 annotators must also use the word "exterior" in their generated logic clauses to represent exterior 462 wall, rather than using "external wall" or other names. Therefore, the problem that annotators may 463 use different words for the same meaning is prevented. The product of the manual transformation 464 process was three versions of logic clauses, with each version independently and manually 465 generated by one of the annotators. After that, annotators reviewed each other's work and 466 collectively generated the final gold standard. All annotators agreed that the gold standard represents 467 the meaning of building codes accurately and approved it.

468 To evaluate the quality of human-generated logic clauses, the authors measured the similarity between the logic clauses generated by different annotators. Because annotators transformed the 469 same building code, they should generate similar logic clauses. A high similarity among human-470 generated logic clauses of different annotators indicates a high quality of the logic clause generation. 471 472 The authors chose to measure logic clause similarity by comparing characters and words at the string level in this study. While text vectorization and cosine similarity measure the meaning-wise 473 474 similarity of natural language text, because the logic clauses generated in this study are not in natural 475 language and our similarity measurement focuses on the existence of logic clause components rather than the meaning of the logic clauses, vectorization of text and cosine similarity were therefore not 476 477 used. The authors measured the similarity among logic clauses in two different ways: (1) the Levenshtein Distance and the Jaro Winkler Similarity were used to measure the character-level 478 similarity between the human-generated logic clauses; and (2) the Jaccard Distance was used to 479 480 measure the predicate-level similarity between the human-generated logic clauses. For example, the Levenshtein Distance, the Jaccard Distance, and the Jaro Winkler Similarity between the two sample 481 logic clauses in Table 1 were 14, 0.67, and 0.97, respectively. Overall, annotators reached an average 482 Levenshtein Distance of 6.88, an average Jaccard Distance of 0.63, and an average Jaro Winkler 483 Similarity of 0.99. 484

485

486 Table 1. Sample Logic Clauses Generated by Annotators

Building Code Sentence	Logic Clause 1	Logic Clause 2
The maximum width of a	compliance_width_of_swinging_door_leaf257(	compliance_width_of_swinging_door_leaf2
swinging door leaf shall	Swinging_door_leaf):-	57(Swinging_door_leaf):-
be 18 inches (1210 mm)	width(Width), swinging_door_leaf(Swinging_d	width(Width),swinging_door_leaf(Swinging
nominal	oor_leaf), has(Swinging_door_leaf, Width), less_	_door_leaf),has(Swinging_door_leaf,Width),
nonnai.	than_or_equal_to(Width,quantity(48,inches)).	equal_to(Width,quantity(48,inches)).

487

488 Because text similarity is task-specific, there was no universally applicable standard for it. Instead, NLP researchers developed their own metrics according to the needs of tasks [80,81]. The 14 489 490 Levenshtein Distance seems high, but it does not consider the length of the text string. If the length 491 of text string is considered, the 0.97 Jaro Winkler Similarity proves that human-generated logic clauses are similar at the character level. The 0.67 Jaccard Distance is relatively low. It indicates a 492 significant number of predicates are different in human-generated logic clauses. However, the 493 494 difference could be overstated because the Jaccard Distance requires two predicates to be 495 completely identical in order to be considered the same. If two predicates are off by even one character, they are still considered different and accounted for in the Jaccard Distance. Combining 496 the use of all three measures illustrates that human-generated logic clauses are similar in general. If 497

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

- 498 two predicates are different, they usually only differ by a few characters. Measurement results using 499 the three metrics show that the three annotators reached a reasonable alignment and the quality of
- 500 the gold standard was good [82,83].

## 501 **5. Result**

To evaluate the performance of the ruleset expansion method, the machine-generated logic clauses 502 503 were compared against the human-generated gold standard. The closer the machine-generated logic 504 clauses are to the gold standard, the better the pattern matching-based rules are, and therefore, the 505 better performance the rules et expansion method is deemed to have. Predicates in logic clauses can be further broken down into predicate elements (i.e., predicate names or predicate arguments). For 506 example, the predicate "has (Stairway, Clear width)" is formed by three elements, "has," "Stairway," 507 and "Clear width." Therefore, the authors calculated both predicate-level performance and 508 predicate element-level performance of the ruleset expansion method. In the predicate-level 509 510 performance, the minimum unit of measurement is a predicate. In the predicate element-level 511 performance, the minimum unit of measurement is a word or phrase. For example, if the gold 512 standard is "considered by(Code change proposals, International fire code development committee)," and the machine-generated logic clause is 513 "considered by(Designation f, International fire code development committee)." The predicate-514 level performance treats this predicate as one incorrect predicate. The predicate element-level 515 516 performance treats this predicate as three elements, two of which are correct and one is incorrect. 517 Because of the phenomenon that predicates could be partially correct. Predicate element-level accuracy could provide a more accurate evaluation regarding the performance of the ruleset 518 expansion method and pattern matching-based rules. 519

- The performance of the expanded pattern matching-based regulatory information transformation 520 ruleset is summarized in Table 2. The performance was measured at the predicate level and the 521 522 predicate element level. The experiment focuses on logic clauses about quantitative requirements 523 because the original ruleset focused on quantitative requirements. In this research, sentences of building code provisions and generated logic clauses have a one-to-many mapping relationship. 524 Patterns, on the other hand, can match the whole sentence or part of a sentence. Regulatory 525 information that spans over multiple sentences is represented by multiple logic clauses. The original 526 ruleset filters quantitative and non-quantitative requirements automatically. Therefore, there is no 527 528 completely missed logic clause. The logic clauses generated that were not in the gold standard were 529 counted as false positives. The logic clauses generated that functioned in the same way as those in the gold standard were counted as true positives. The logic clause-level performance is reported in 530 Table 3. The predicate element-level performance was higher than the predicate-level performance, 531 532 which indicates some predicates were partially correct. Through error analysis, the authors 533 recognized four main sources of errors:
- The partially correct predicates. After reviewing machine-generated logic clauses and the gold
   standard, the authors found that a significant portion of predicates in machine-generated logic

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

- clauses were partially correct. For example, when the correct predicate in the gold standard is
  "surrounded\_by(Buildings,Public\_ways)," the expanded ruleset generated a partially correct
  predicate "surrounded\_by(Chapter\_9,Public\_ways)." Future development about pattern
  matching-based rules could focus more on capturing the correct elements in predicates.
- 540 The flexibility of B-Prolog logic clauses and the ambiguity of natural language. The flexibility 2. of B-Prolog logic clauses and the inherent ambiguity of natural language left a room of 541 542 interpretation for the annotators. In other words, it was possible to represent the same building 543 code requirement in different predicates. For example, one annotator translated the "minimum fire resistant rating of 1 hour" to "greater than(Minimum fire resistance rating, quantity(1, 544 545 hour))." Another annotator translated the same phrase to "equal to (Minimum fire resistance rating, quantity(1, hour))." One annotator considered that fire 546 547 resistant rating of a subject should be greater than the minimum fire-resistant rating, which is 1 hour. Another annotator considered that this phrase means the minimum fire-resistant rating of 548 549 a subject should be 1 hour. Therefore, the precision of the rule generation was affected. A viable solution to increase inter-annotator agreement may include more detailed and stricter annotation 550 551 guidelines.
- The patterns and terminologies unseen in Chapter 10. Although most regulatory information 552 3. patterns in Chapter 5 were captured in Chapter 10, a small amount of regulatory information 553 patterns were missed. The authors attributed missed building code requirements to unseen 554 555 patterns or unique terminologies in Chapter 5. Chapter 5 and Chapter 10 of the IBC 2015 focus on different topics (i.e., general building height and area, and means of egress, respectively), so 556 some terminologies in Chapter 5 did not occur in Chapter 10. For example, the erroneous 557 predicate "unobstructed to(Be,Room)" was generated instead of "unobstructed to(Room)" due 558 to a pattern that did not occur in Chapter 10. Such error will need to be addressed by 559 accumulatively expanding the training dataset. 560
- 4. The backward compatibility requirement. The generality and validity requirements of the ruleset expansion method ensures the quality of the generated logic clauses and shows a promising future that pattern matching rule-based regulatory information extraction can potentially capture all building code requirements with a sufficiently comprehensive set of rules. However, the compatibility requirement also forbade the removal of any existing rule, which led to some false positives. In the future, the flexibility of modifying existing rules may need to be tested.

## 569 Table 2. Performance of Applying Ruleset Expansion Method

568

	Training				Tes	ting		
	Predicate level		Predicate element level		Predicate level		Predicate element level	
	Before <sup>1</sup>	After <sup>2</sup>	Before <sup>1</sup>	After <sup>2</sup>	Before <sup>1</sup>	After <sup>2</sup>	Before <sup>1</sup>	After <sup>2</sup>
Precision	84.35%	96.31%	87.90%	98.33%	86.17%	95.17%	90.03%	97.48%
Recall	79.00%	98.38%	81.78%	99.39%	76.84%	96.60%	81.77%	98.65%
F1-score	81.59%	97.34%	84.73%	98.86%	81.24%	95.88%	85.70%	98.06%

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

<sup>1</sup> Performance of the original ruleset, which is the ruleset before the application of the ruleset expansion method. <sup>2</sup> Performance of the expanded ruleset, which is the ruleset after the application of the ruleset expansion method.

#### 570 Table 3. Logic Clause-Level Performance

	Tra	ining	Testing		
	Before	After <sup>2</sup>	Before	After <sup>2</sup>	
Precision	89.86%	100.00%	93.81%	100.00%	
Recall	98.27%	99.68%	96.81%	97.98%	
F1-score	93.87%	99.84%	95.29%	98.98%	

<sup>1</sup> Performance of the original ruleset, which is the ruleset before the application of the ruleset expansion method. <sup>2</sup> Performance of the expanded ruleset, which is the ruleset after the application of the ruleset expansion method.

#### 6. Real-World Case Study 571

572 The proposed ruleset expansion method was also evaluated in a real-world test case (Figure 4). The 573 test case was based on a convention store in Texas to be checked against Chapter 5 and Chapter 10 of the IBC 2015. The BIM model and compliance checking report for the convention store were 574 provided by an industry partner of this research. According to the compliance checking report, five 575 applicable logic clauses in those sections were selected in the test case. Building design information 576 was extracted from the BIM model of the test case and further elaborated to fit the needs of the case 577 578 study. To make the test case comprehensive, three non-compliance instances with the applicable 579 logic clauses were added into the design information. In the test run, all three non-compliance instances were successfully detected. Example checked logic clauses and corresponding non-580 581 compliance instances are shown in Table 4.



582 583

Figure 4. A Real-World Test Case for ACC.

584



For the final p	ublished version, p	lease find it at the B	Elsevier Database Here:
https://www.science	direct.com/science/	/article/abs/pii/S092	26580522001030?via%3D ihub

Section in IBC 2015	Building code requirement	Logic clause	Non- compliance instance
1010	"The required capacity of each door opening shall be sufficient for the occupant load thereof and shall provide a minimum clear width of 32 inches (813 mm)." [1]	compliance_door_opening(Door_opening ):- required_capacity(Required_capacity),do or_opening(Door_opening),has(Door_ope ning,Required_capacity),sufficient(Door_ opening),for(Door_opening,Occupant_loa d),occupant_load(Occupant_load),provide (Required_capacity,Minimum_clear_widt h),minimum_clear_width(Minimum_clear_ width),equal_to(Minimum_clear_width, quantity(32,inches)).).	A door opening provided a clear width of 27 inches.
1010	"The height of door openings shall be not less than 80 inches (2032mm)." [1]	compliance_height_of_door_openings(Do or_openings):- height(Height),door_openings(Door_open ings),has(Door_openings,Height),not less_than(Height,quantity(80,inches)).	The height of a door opening is 72 inches.
506	"To qualify for an area factor increase based on frontage, the public way or open space adjacent to the building perimeter shall have a minimum distance (W) of 20 feet (6096 mm) measured at right angles from the building face to any of the following:" [1]	compliance_Frontage(Frontage):- qualify_for(Qualify,Area_factor_increase ),qualify(Qualify),area_factor_increase(A rea_factor_increase),based_on(Area_facto r_increase,Frontage),(frontage(Frontage); public_way(Frontage);open_space(Fronta ge)),minimum_distance(Minimum_distan ce),has(Building_perimeter,Minimum_dis tance),measured_at(Frontage,Right_angle s),right_angles(Right_angles),from(Right _angles,Building_face),building_face(Bui lding_face),equal_to(Minimum_distance, quantity(20,feet)),associated(Frontage,Mi _nimum_distance).	The distance from the building face is 18 feet.

## 586 **7. Robustness**

In the experiment illustrated above, the proposed method was tested for expanding the range of 587 checkable building code requirements to new chapters of building codes, which are in different 588 domains/topics of the building code from which the original ruleset was initially developed. To 589 590 further evaluate the robustness of the expanded ruleset, the authors also tested it on processing 591 construction contracts, a fundamentally different type of construction documents compared to building codes. Nine free and openly available construction contracts or construction contract 592 templates were collected. In total, 185 sentences were extracted from these contracts. The expanded 593 ruleset was then executed to convert the extracted sentences in construction contracts to logic 594 595 clauses. The performance of the expanded ruleset is illustrated in Table 5. Examples of contract 596 contents and corresponding logic clauses generated are shown in Table 6. The results show the

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

597 robustness of the expanded ruleset is promising (although not perfect) for processing construction

598	documents	beyond	the original	l intent of	fbuilding	codes.
-----	-----------	--------	--------------	-------------	-----------	--------

599

#### 600 Table 5. Performance on Processing Construction Contract

	Predicate level	Predicate element level
Precision	90.52%	97.20%
Recall	92.92%	98.42%
F1-score	91.70%	97.81%

601

602 Table 6. Examples of Contract Sentences and Corresponding Logic Clauses Generated

Contract sentence	Logic clause
Two copies of the Contract Documents shall be signed by the Owner and the Contractor.[84]	compliance_Owner3(Owner):- number_prep(Number),copies(Copies),has(Copies,Number),contract_do cuments(Contract_Documents),has(Contract_Documents,Copies),(owne r(Owner);contractor(Owner)),signed_by(Contract_Documents,Owner),e qual_to(Number,quantity(2,one)).
Contractor shall maintain in a safe place at the Property one record copy of all drawings, specifications, addenda, written amendments, and the like in good order and annotated to show all changes made during construction, which will be delivered to Owner upon completion of the Work. [85]	compliance_Number1(Number):- maintain_in(Contractor,Safe_place),contractor(Contractor),safe_place(S afe_place),at(Safe_place,Property),property(Property),number_prep(Nu mber),record_copy(Record_copy),has(Record_copy,Number),drawings( Drawings),has(Drawings,Record_copy),like_in(Like,Good_order),like( Like),good_order(Good_order),to_show(Good_order,Changes),changes( Changes),made_during(Changes,Construction),construction(Construction n),delivered_to(Good_order,Owner),owner(Owner),upon(Owner,Compl etion),completion(Completion),work(Work),has(Work,Completion),equ al_to(Number,quantity(1,one)),associated(Safe_place,Good_order).
If the final amount of the ALLOWANCE work is less than the ALLOWANCE line item amount listed in the Agreement, a credit will be issued to Owner after all billings related to this particular line item ALLOWANCE work have been received by Contractor.	final_amount/Final_amount/if(Final_amount). <sup>2</sup> final_amount(Final_amount),if(Final_amount),allowance_work(ALLO WANCE_work),has(ALLOWANCE_work,Final_amount),allowance_li ne_item_amount(ALLOWANCE_line_item_amount),listed_in(ALLOW ANCE_line_item_amount,Agreement),agreement(Agreement),credit(Cr edit),owner(Owner),issued_to(Credit,Owner),after(Owner,Billings),billi ngs(Billings),related_to(Billings,This_particular_line_item_ALLOWAN CE_work),this_particular_line_item_allowance_work(This_particular_line_item_ant_LOWANCE_work),received_by(This_particular_line_item_ant_compared)
86	ALLOWANCE work, Contractor), contractor(Contractor), less than (Fin

al amount,quantity(1,ALLOWANCE line item amount)).

603

## **8. Contributions to the Body of Knowledge**

This research contributes to the body of knowledge in four main ways. First, this research proves the feasibility of expanding the range of checkable building code requirements by expanding an existing regulatory information transformation ruleset. The authors expanded the range of checkable building code requirements of an automated code compliance checking system to cover Chapter 5 and Chapter 10 of the IBC 2015. This expansion was achieved by 64 new rules. It shows that different chapters of the IBC share similar patterns, and the number of new pattern matching-based

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

rules needed to expand the range of checkable code requirements is small. Second, this research 611 612 proposed a new ruleset expansion method. This method ensures the quality of added pattern-613 matching-based rules and, therefore, the quality of logic clauses generated by the pattern matchingbased rules. In a previous study, three hundred and six rules were developed to cover two chapters 614 of building code. In comparison, only sixty-four new rules were developed to cover two new 615 chapters of building code. This research shows that the marginal cost of expanding the range of 616 617 checkable building code requirements is low. It provides a workable and low-cost method to expand 618 the range of checkable code requirements of ACC systems. The cost of expanding the range of checkable building codes by expanding an existing regulatory information transformation ruleset 619 could further decrease in the future as the number of existing rules increases, because building codes 620 share similar patterns and the number of unseen patterns in new building codes could decrease as 621 622 existing pattern matching-based rules cover more patterns in building codes. Future researchers and developers can adopt this method to expand the range of checkable code requirements of the ACC 623 624 system and bring the ACC system to full deployment in the AEC industry. While the research 625 focuses on processing building codes, the testing results of transforming construction contracts show that the proposed ruleset expansion method is potentially robust in processing different types of 626 construction documents. Third, this research also generated a dataset of building codes in logic 627 clauses. This dataset can facilitate other regulatory information transformation research, such as 628 629 machine learning-based logic clause generation. Last but not least, this research facilities the 630 adoption of ACC in the AEC industry. With an expanded range of checkable code requirements, the utility of ACC is enhanced. ACC can reduce the time, cost, and human-errors in code compliance 631 checking and encourages the AEC industry to shift towards a digital paradigm. 632

## 633 **9. Conclusion**

The paper presents a ruleset expansion method that can extend the range of checkable code 634 635 requirements of ACC systems to different chapters of the International Building Code (IBC), which 636 can potentially be applied to other codes beyond the IBC and other construction documents such as contracts. It takes an iterative approach to ensure the generality and validity of the added pattern 637 matching-based rules and the generated logic clauses. Experimental results on Chapters 5 and 10 of 638 IBC 2015 showed the expanded ruleset generated logic clauses with 95.17% predicate-level 639 640 precision, 96.60% predicate-level recall, and 95.88% predicate-level F1-score. This performance 641 proved the effectiveness of the ruleset expansion method and the expanded ruleset. Through error 642 analysis, the authors attributed the remaining errors to the flexibility of B-Prolog language, the ambiguity in natural language, missed building code requirement patterns, and the compatibility 643 requirement. The authors also suggested solutions to further increase the performance of the ruleset 644 645 expansion method, such as expanding the training dataset and providing stricter annotation guidelines. This research also presents a dataset of logic clauses. This dataset has the potential to 646 647 facilitate research on different regulatory information transformation approaches, such as machine learning-based logic clause generation. The research findings in this study can be used to build fully 648

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D ihub

automated building code compliance checking systems with wider code requirements coverage than 649 650 the state of the art. The demonstrated decreasing marginal cost of transformation rule development 651 and high predicate-level performance renders the rule-based processing of building code requirements promising to bring fully automated building code compliance checking to real-world 652 applications. Future research is needed to discover the boundary of the theoretical "superset" of 653 common patterns used in building code transformation rules, for guiding the practical 654 655 implementation of the demonstrated rule-based processing of building code requirements in real 656 ACC systems. Furthermore, the successful demonstration of such processing in construction contracts in this study helps open the door to rule-based processing of a variety of textual documents 657 in the AEC industry, to support future automation and AI applications in the AEC industry. 658

## 659 **10.Acknowledgement**

660 The authors would like to thank the National Science Foundation (NSF). This material is based on 661 work supported by the NSF under Grant No. 1827733. Any opinions, findings, and conclusions or 662 recommendations expressed in this material are those of the authors and do not necessarily reflect 663 the views of the NSF.

## 664 **11.Reference**

- 665 [1] International Code Council, International Building Code, International Code Council,
  666 2015.1609837398, <u>https://codes.iccsafe.org/content/IBC2015</u>, Last Access: Feb 12<sup>th</sup>, 2022.
- 667 [2] International Code Council, International Fire Code, International Code Council,
  668 2015.1609837398, <u>https://codes.iccsafe.org/content/IFC2015</u>, Last Access: Feb 12<sup>th</sup>, 2022..
- Indiana Department of Homeland Security, Codes, Standards and Other Rules, in: Secondary
  Indiana Department of Homeland Security (Ed.), Secondary Codes, Standards and Other Rules,
  https://www.in.gov/dhs/2490.htm#top, Last Access: Jan 21<sup>st</sup>,2020.
- [4] S. Moon, G. Lee, S. Chi, H. Oh, Automatic Review of Construction Specifications Using
  Natural Language Processing, in: Proceedings of Computing in Civil Engineering 2019: Data,
  Sensing, and Analytics, American Society of Civil Engineers, Atlanta, Georgia, 2019, pp. 401407. https://doi.org/10.1061/9780784482438.051.
- 676 [5] C. Preidel, A. Borrmann, Refinement of the visual code checking language for an automated
  677 checking of building information models regarding applicable regulations, in: Proceedings of
  678 Computing in Civil Engineering 2017, American Society of Civil Engineers, Seattle,
  679 Washington, 2017, pp. 157-165. https://doi.org/10.1061/9780784480823.020.
- A. Alghamdi, M. Sulaiman, A. Alghamdi, M. Alhosan, M. Mastali, J. Zhang, Building 680 [6] 681 accessibility code compliance verification using game simulations in virtual reality, Computing 682 in Civil Engineering 2017, Seattle, Washington, 2017, pp. 262-270. https://doi.org/10.1061/9780784480830.033. 683

684	[7]	City of Chicago, Average Time for Permit Issuance, in: Secondary City of Chicago (Ed.),
685		Secondary Average Time for Permit Issuance,
686		https://www.chicago.gov/city/en/depts/bldgs.html, Last Access: Jan, 21st, 2020.
687	[8]	City of Chicago, Permit Fee Calculator, in: Secondary City of Chicago (Ed.), Secondary Permit
688		Fee Calculator, https://ipiweb.cityofchicago.org/DynamicPortal/Forms/FeeCalculator.aspx,
689		Last access: Last Access: Jan, 21 <sup>st</sup> , 2020.
690	[9]	X. Tan, A. Hammad, P. Fazio, Automated code compliance checking for building envelope
691		design, Journal of Computing in Civil Engineering 24 (2) (2010) pp. 203-211,
692		https://doi.org/10.1061/(ASCE)0887-3801(2010)24:2(203).
693	[10]	J. Martins, V. Abrantes, Automated code-checking as a driver of BIM adoption, International
694		Journal for Housing Science, 34 (2010) pp. 287-295,
695		http://www.housingscience.org/html/publications/pdf/34-4-6.pdf, Last Access: Feb 12th, 2022.
696	[11]	D. Greenwood, S. Lockley, S. Malsane, J. Matthews, Automated compliance checking using
697		building information models, in: Proceedings of the The Construction, Building and Real Estate
698		Research Conference of the Royal Institution of Chartered Surveyors, Royal Institution of
699		Chartered Surveyors (RICS), Paris, France, 2010, pp. 266-274.
700		https://www.irbnet.de/daten/iconda/CIB20057.pdf.
701	[12]	J. Choi, I. Kim, An approach to share architectural drawing information and document
702		information for automated code checking system, Tsinghua Science and Technology 13 (S1)
703		(2008) pp. 171-178, https://doi.org/10.1016/S1007-0214(08)70145-7.
704	[13]	J. Zhang, N.M. El-Gohary, Integrating semantic NLP and logic reasoning into a unified system
705		for fully-automated code checking, Automation in Construction 73 (2017) pp. 45-57,
706		https://doi.org/10.1016/j.autcon.2016.08.027.
707	[14]	S.J. Fenves, Tabular decision logic for structural design, Journal of the Structural Division 92
708		(6) (1966) pp. 473-490, https://doi.org/10.1061/JSDEAG.0001567.
709	[15]	J.H. Garrett, S.J. Fenves, A knowledge-based standards processor for structural component
710		design, Engineering with Computers 2 (4) (1987) pp. 219-238,
711		https://doi.org/10.1007/BF01276414.
712	[16]	L. Lopez, S. Elam, K. Reed, Software concept for checking engineering designs for
713		conformance with codes and standards, Engineering with Computers 5 (2) (1989) pp. 63-78,
714		https://doi.org/10.1007/BF01199070.
715	[17]	E.A. Delis, A. Delis, Automatic fire-code checking using expert-system technology, Journal of
716		Computing in Civil Engineering 9 (2) (1995) pp. 141-156,
717		https://doi.org/10.1061/(ASCE)0887-3801(1995)9:2(141).
718	[18]	J.S. Gero, Design prototypes: a knowledge representation schema for design, AI Magazine 11
719		(4) (1990) pp. 26-36, https://doi.org/10.1609/aimag.v11i4.854.
720	[19]	L. Ding, R. Drogemuller, M. Rosenman, D. Marchant, J. Gero, Automating code checking for
721		building designs-DesignCheck, in: Proceedings of the Construction Research Congress (CRC)
722		for Construction Innovation. American Society of Civil Engineerings, Grand Bahama Island,
723		Bahamas, 2006 pp. 1-16, https://doi.org/10.1108/ecam.2006.28613faa.002, Last Access: Feb
724		12 <sup>th</sup> , 2022.

725	[20]	J. Dimyadi, R. Amor, Automated building code compliance checking-where is it at?, in:
726		Proceedings of the 19th International World Building Congress, (6), International Council for
727		Research and Innovation in Building and Construction, Brisbane, Australia, 2013.
728		https://doi.org/10.13140/2.1.4920.4161.
729	[21]	A.R. Sabouni, O.M. Al-Mourad, Quantitative knowledge based approach for preliminary
730		design of tall buildings, Artificial Intelligence in Engineering 11 (2) (1997) pp. 143-154,
731		https://doi.org/10.1007/s00216-016-0157-x.
732	[22]	C. Eastman, J.M. Lee, Y.S. Jeong, J.K. Lee, Automatic rule-based checking of building designs,
733		Automation in Construction 18 (8) (2009) pp. 1011-1033,
734		https://doi.org/10.1016/j.autcon.2009.07.002 .
735	[23]	Y. Shao, C. Hardmeier, J. Tiedemann, J. Nivre, Character-based joint segmentation and POS
736		tagging for Chinese using bidirectional RNN-CRF, arXiv preprint arXiv:1704.01314 (2017),
737		https://aclanthology.org/I17-1018, Last Access: Feb 12th, 2022.
738	[24]	B.T. Zhong, L.Y. Ding, H.B. Luo, Y. Zhou, Y.Z. Hu, H.M. Hu, Ontology-based semantic
739		modeling of regulation constraint for automated construction quality compliance checking,
740		Automation in Construction 28 (2012) pp. 58-70, https://doi.org/10.1016/j.autcon.2012.06.006.
741	[25]	J. Zhang, N.M. El-Gohary, Semantic NLP-based information extraction from construction
742		regulatory documents for automated compliance checking, Journal of Computing in Civil
743		Engineering 30 (2) (2016) pp. 04015014, https://doi.org/10.1061/(ASCE)CP.1943-
744		5487.0000346.
745	[26]	J. Dimyadi, G. Clifton, M. Spearpoint, R. Amor, Computerizing Regulatory Knowledge for
746		Building Engineering Design, Journal of Computing in Civil Engineering (2016) pp. C4016001,
747		https://doi.org/10.1061/(ASCE)CP.1943-5487.0000572.
748	[27]	S.M. İlal, H.M. Günaydın, Computer representation of building codes for automated
749		compliance checking, Automation in Construction 82 (2017) pp. 43-58,
750		https://doi.org/10.1016/j.autcon.2017.06.018.
751	[28]	T.H. Nguyen, Integrating building code compliance checking into a 3D CAD system, in:
752		Proceedings of the Computing in Civil Engineering (2005), American Society of Civil
753		Engineers, Cancun, Mexico, 2005, pp. 1-12. https://doi.org/10.13140/2.1.4920.4161.
754	[29]	N.O. Nawari, Automating codes conformance, Journal of Architectural Engineering 18 (4)
755		(2012) pp. 315-323, https://doi.org/10.1061/41182(416)70.
756	[30]	J. Zhang, N.M. El-Gohary, Semantic-based logic representation and reasoning for automated
757		regulatory compliance checking, Journal of Computing in Civil Engineering 31 (1) (2017) pp.
758		04016037, https://doi.org/10.1061/(ASCE)CP.1943-5487.0000583.
759	[31]	J. Zhang, N.M. El-Gohary, Automated information transformation for automated regulatory
760		$compliance\ checking\ in\ construction,\ Journal\ of\ Computing\ in\ Civil\ Engineering\ 29\ (4)\ (2015)$
761		pp.B4015001, https://doi.org/:10.1061/(ASCE)CP.1943-5487.0000427.
762	[32]	J. Zhang, N.M. El-Gohary, Automated extraction of information from building information
763		models into a semantic logic-based representation, Congress on Computing in Civil
764		Engineering, Proceedings 2015 (2015) pp. 173-180,

765		https://doi.org/10.1061/9780784479247.022.
766	[33]	Zhang, J., and El-Gohary, N. Integrating semantic NLP and logic reasoning into a unified
767		system for fully-automated code checking. Automation in Construction, 73, pp. 45-57,
768		https://doi.org/10.1016/j.autcon.2016.08.027.
769	[34]	Zhang, J. and El-Gohary, N. Semantic-based logic representation and reasoning for
770		automated regulatory compliance checking. Journal of Computing in Civil Engineering, 31
771		(1), pp.4016037, https://doi.org/10.1061/(ASCE)CP.1943-5487.0000583.
772	[35]	S.A. Morshed, X. Lv, R.B. Tanvir, Network-based information extraction from IFC files to
773		support intelligent BIM companion (iBcom) technology, in: Proceedings of Construction
774		Research Congress 2020, American Society of Civil Engineers, Tempe, Arizona, 2020, pp. 427-
775		435. https://doi.org/10.1061/9780784482865.046.
776	[36]	P. Zhou, N. El-Gohary, Automated matching of design information in BIM to regulatory
777		information in energy codes, in: Proceedings of the Construction Research Congress 2018,
778		American Society of Civil Engineers, New Orleans, Louisiana, 2018, pp. 75-85.
779		https://doi.org/10.1061/9780784481264.008.
780	[37]	G.G. Chowdhury, Natural language processing, Annual review of information science and
781		technology 37(1) (2003) pp. 51-89, https://doi.org/10.1002/aris.1440370103.
782	[38]	C.D. Manning, P. Raghavan, H. Schütze, Introduction to information retrieval, Cambridge
783		University Press, 2008.0521865719, https://wiltrud.hwro.de/teaching/ir12w/pdf/19web.pdf,
784		Last Access: Feb 12 <sup>th</sup> , 2022
785	[39]	J. Cowie, W. Lehnert, Information extraction, Communications of the Association for
786		Computing Machinery (ACM), 39(1)(1996) pp. 80-91, https://doi.org/10.1145/234173.234209.
787	[40]	X. Zhang, J. Zhao, Y. LeCun, Character-level convolutional networks for text classification, in:
788		Proceedings of the 28th International Conference on Neural Information Processing Systems,
789		Association for Computing Machinery (ACM), Montreal, Canada, 2015, pp. 649-657.
790		https://dl.acm.org/doi/10.5555/2969239.2969312.
791	[41]	K. McKeown, Text generation, Cambridge University Press, 1992.0521438020,
792		https://www.cambridge.org/core/books/text-
793		generation/EC9075AB5F64A7AFDC5A9513237505CA, Last Access: Feb 12th, 2022.
794	[42]	A. Nenkova, K. McKeown, A Survey of Text Summarization Techniques, Mining Text Data
795		(2012) pp. 43-76, https://doi.org/10.1007/978-1-4614-3223-4_3.
796	[43]	R. Barzilay, M. Elhadad, Using lexical chains for text summarization, Advances in Automatic
797		Text Summarization (1999) pp. 111-121, https://aclanthology.org/W97-0703, ,Last Access: Feb
798		12 <sup>th</sup> , 2022.
799	[44]	P. Koehn, Statistical machine translation, Cambridge University Press, 2009.1139483307,
800		https://www.cambridge.org/core/books/statistical-machine-
801		translation/94EADF9F680558E13BE759997553CDE5, Last Access: Feb 12th, 2022.
802	[45]	D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P.
803		Motlicek, Y. Qian, P. Schwarz, The Kaldi speech recognition toolkit, in: Proceedings of the
804		Institute of Electrical and Electronics Engineers (IEEE)2011 Workshop on Automatic Speech
805		Recognition and Understanding, Institute of Electrical and Electronics Engineers (IEEE)Signal

	Processing Society, New York, New York, 2011.
	https://www.danielpovey.com/files/2011_asru_kaldi.pdf, Last Access: Feb 12th, 2022.
[46]	K. Gali, H. Surana, A. Vaidya, P.M. Shishtla, D.M. Sharma, Aggregating machine learning and
	rule based heuristics for named entity recognition, in: Proceedings of the International Joint
	Conference on Natural Language Processing-08 (IJCNLP) Workshop on Named Entity
	Recognition for South and South East Asian Languages, Association for Computational
	$Linguistics, Hyderabad, India, 2008, \underline{https://www.aclweb.org/anthology/I08-5005}, Last Access:$
	Feb 12 <sup>th</sup> , 2022.
[47]	K. Crowston, X. Liu, E.E. Allen, Machine learning and rule-based automated coding of
	qualitative data, in: Proceedings of the American Society for Information Science and
	Technology, (47), Association for Information Science and Technology, Pittsburgh,
	Pennsylvania, 2010, pp. 1-2. https://doi.org/10.1002/meet.14504701328.
[48]	F. Chollet, Deep Learning with Python, 2017.1617294438,
	https://www.manning.com/books/deep-learning-with-python, ,Last Access: Feb 12th, 2022.
[49]	A.J.P. Tixier, M.R. Hallowell, B. Rajagopalan, D. Bowman, Automated content analysis for
	construction safety: A natural language processing system to extract precursors and outcomes
	from unstructured injury reports, Automation in Construction 62 (2016) pp. 45-56,
	https://doi.org/10.1016/j.autcon.2015.11.001.
[50]	J.R. Lin, Z. Hu, J. Zhang, BIM oriented intelligent data mining and representation, in:
	Proceedings of 30th International Council for Research and Innovation in Building and
	Construction W78 International Conference on Applications of IT in the AEC Industry,
	International Council for Research and Innovation in Building and Construction, Beijing, China,
	2013, pp. 280-289. https://itc.scix.net/pdfs/w78-2013-paper-112.pdf, Last Access: Feb 12 <sup>nd</sup> ,
	2022 .
[51]	R. Zhang, N. El-Gohary, A machine-learning approach for emantic matching of building codes
	and Building Information Models (BIMs) for supporting automated code checking, in: the
	Proceedings of the International Congress and Exhibition "Sustainable Civil Infrastructures",
	Springer, Egypt, Egypt, 2019, pp. 64-73. https://doi.org/10.1007/978-3-030-34216-6_5.
[52]	J. Zhang, N. El-Gohary, Extending building information models semiautomatically using
•	semantic natural language processing techniques, Journal of Computing in Civil Engineering
	(2016) pp. C4016004, https://doi.org/10.1061/(ASCE)CP.1943-5487.0000536.
[53]	E. Brill, A simple rule-based part of speech tagger, in: the Proceedings of the Third Conference
	on Applied Natural Language Processing, Association for Computational Linguistics, Trento,
	Italy 1992, pp. 152-155. https://doi.org/10.3115/974499.974526.
[54]	Butte College, The eight parts of speech, in: Secondary Butte College (Ed.), Secondary, The
	Eight Parts of Speech,
	http://www.butte.edu/departments/cas/tipsheets/grammar/parts_of_speech.html, Last Access:
	Sep 11 <sup>st</sup> ,2019.
[55]	M. Marcus, B. Santorini, M.A. Marcinkiewicz, Building a large annotated corpus of English:
-	The Penn Treebank, (1993), <u>https://doi.org/10.5555/972470.972475</u> , Last Access: Feb 12 <sup>nd</sup> ,

For the final published version, please find it at the Elsevier Database Here: <a href="https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D">https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3D</a> ihub

847	[6/]	2022.						
848	[56]	W.N. Francis, H. Kucera, Brown Corpus Manual, Brown University, 1979,						
849	[57]	<u>http://korpus.uid.no/icame/manuals/BROWN/INDEX.HTM</u> , Last Access: Feb 12 <sup>-10</sup> , 2022.						
850	[37]	H. Schmid, Part-or-speech tagging with neural networks, in: the Proceedings of the 15th						
851		Longen 1004 np. 172-176, https://doi.org/10.2115/001886.001015						
052	F <b>£</b> 01	Japan, 1994, pp. 1/2-1/0. https://doi.org/10.5113/991880.991915.						
803 0E4	[38]	S. Bird, E. Klein, E. Löper, Natural language processing with Fython: analyzing text with the network longuage toolly it. O'Beilly Media Inc. 2000 0506555717. https://www.nltls.org/healt/						
004 055		Last A pages Esh 12 <sup>nd</sup> 2022						
000	[50]	Last Access. Feb 12, 2022.						
000	[39]	J. Onnenez, L. Marquez, Fast and accurate part-of-speech tagging. The SVM approach revisited,						
857 050		in: the Proceedings of the Recent Advances in Natural Language Processing III, Association for						
858		Computational Linguistics, Borovets, Bulgaria, 2003, pp. 153-162.						
859	[(0]	nttps://dbip.org/rec/conf/ranip/GimenezM03.						
860	[60]	1. Brants, 1n 1: a statistical part-of-speech tagger, in: the Proceedings of the sixth conference on						
861		Applied natural language processing, Association for Computational Linguistics, Seattle,						
862	5611	washington, 2000, pp. 224-231. https://doi.org/10.3113/974147.974178.						
863	[61]	B. Plank, A. Søgaard, Y. Goldberg, Multilingual part-of-speech tagging with bidirectional long						
864		short-term memory models and auxiliary loss, in: Proceedings of the 54th Annual Meeting of						
865		the Association for Computational Linguistics, (2), Association for Computational Linguistics,						
800	[(0]	Berlin, Germany, 2010, pp. 412–418. https://doi.org/10.18053/v1/P10-2067.						
867	[62]	1.R. Gruber, A translation approach to portable ontology specifications, Knowledge Acquisition $5(2)(1002) = 100,221,144 = 1/(11) = 1002,1002$						
868	[(2]	5 (2) (1993) pp. 199-221, https://doi.org/10.1006/knac.1993.1008.						
869	[63]	D. Brickley, R.V. Guha, A. Layman, Resource description framework (RDF) schema $i_{1}^{c}$ $i_{2}^{c}$ $i_{3}^{c}$ $i_{4}^{c}$ $i_{5}^{c}$						
870		specification, Encyclopedia of Database Systems, (1999), https://doi.org/10.100//9/8-0-38/-						
871	FC 41	39940-9_1319.						
872	[64]	J. Hendler, D.L. McGumness, The DARPA agent markup language, in: Proceedings of Institute						
8/3		of Electrical and Electronics Engineers (IEEE) Intelligent Systems, (15), Institute of Electrical						
874		and Electronics Engineers (IEEE), New York, New York, 2000, pp. 6/-/3. http://www-						
875	[(5]	<u>KSI.stanford.edu/pub/KSL_Reports/KSL-00-10.html</u> , Last Access: Feb 12 <sup>nd</sup> , 2022.						
876	[65]	D.L. Mcgumness, R. Fikes, J. Hendler, L.A. Stein, DAML+OIL: an ontology language for the						
8//		Semantic Web, in: Proceedings of Institute of Electrical and Electronics Engineers						
878		(IEEE)Intelligent Systems 17(5), Institute of Electrical and Electronics Engineers (IEEE), New						
879	5663	York, New York, 2002 pp. 72-80, https://doi.org/10.1109/MIS.2002.1039835.						
880	[66]	L. Amos, D. Anderson, S. Brody, A. Ripple, B.L. Humphreys, UMLS users and uses: a current						
881		overview, Journal of the American Medical Informatics Association 27 (10) (2020) pp.						
882		https://doi.org/1606-1611, 10.1093/jamia/ocaa084.						
883	[67]	M. Hepp, Goodrelations: An ontology for describing products and services offers on the web,						
884		in: Proceeding of the International Conference on Knowledge Engineering and Knowledge						
885		Management, Springer, Berlin, Heidelberg, 2008, pp. 329-346. https://doi.org/10.1007/978-3-						
886	F ( 03	540-87696-0_29.						
887	[68]	N.F. Noy, D.L. McGunness, Ontology development 101: A guide to creating your first ontology,						

888		in: Secondary N.F. Noy, D.L. McGuinness (Eds.), Secondary Ontology development 101: A							
889		guide to creating your first ontology, <u>http://www.ksl.stanford.edu/people/dlm/papers/ontology-</u>							
890		tutorial-noy-mcguinness-abstract.html, Last Access: May 22 <sup>nd</sup> ,2021							
891	[69]	W.H. Gomaa, A.A. Fahmy, A survey of text similarity approaches, International Journal of							
892		Computer Applications 68 (13) (2013) pp. 13-18, https://doi.org/10.5120/11638-7118.							
893	[70]	Z. Su, B. Ahn, K. Eom, M. Kang, J. Kim, M. Kim, Plagiarism Ddetection using the Levenshtein							
894		Distance and Smith-Waterman Algorithm, in: Proceeding of the 2008 3rd International							
895		Conference on Innovative Computing Information and Control, Institute of Electrical and							
896		Electronics Engineers (IEEE), Dalian, China, 2008, pp. 569-569							
897		https://doi.org/10.1109/ICICIC.2008.422.							
898	[71]	VI. Levenshtein, Binary codes capable of correcting deletions, insertions, and reversals, Soviet							
899		Physics Doklady, (10), 1966, pp. 707-710.							
900		https://ui.adsabs.harvard.edu/abs/1966SPhD10707L/abstract, Last Access: Feb 12 <sup>nd</sup> , 2022.							
901	[72]	S. Kosub, Anote on the triangle inequality for the Jaccard distance, Pattern Recognition Letters							
902		120 (2019) pp. 36-38, https://doi.org/10.1016/j.patrec.2018.12.007.							
903	[73]	W.E. Winkler, String comparator metrics and enhanced decision rules in the Fellegi-Sunter sodel							
904		of record linkage,, https://eric.ed.gov/?id=ED325505, Last Access: May 22 <sup>nd</sup> , 2021.							
905	[74]	T. Mikolov, I. Sutskever, K. Chen, G. Corrado, J. Dean, Distributed representations of words							
906		and phrases and their compositionality, in: Proceedings of the 26th International Conference on							
907		Neural Information Processing Systems, (2), Curran Associates Inc., Lake Tahoe, Nevada, 2013,							
908		pp. 3111-3119. https://doi.org/10.5555/2999792.2999959.							
909	[75]	A. Joulin, E. Grave, P. Bojanowski, M. Douze, H. Jégou, T. Mikolov, Fasttext. zip: Compressing							
910		text classification models, arXiv preprint arXiv:1612.03651 (2016),							
911		https://arxiv.org/abs/1612.03651.							
912	[76]	J. Pennington, R. Socher, C.D. Manning, Glove: Global vectors for word representation, in:							
913		Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing							
914		(EMNLP), Association for Computational Linguistics, Doha, Qatar, 2014, pp. 1532-1543.							
915		https://doi.org/10.3115/v1/D14-1162.							
916	[77]	A. Ben Abacha, P. Zweigenbaum, Automatic extraction of semantic relations between medical							
917	•	entities: a rule based approach, Journal of Biomedical Semantics 2 (5) (2011) pp. S4,							
918		https://doi.org/10.1186/2041-1480-2-S5-S4.							
919	[78]	C. Zhang, X. Zhang, W. Jiang, Q. Shen, S. Zhang, Rule-Based Extraction of Spatial Relations							
920		in Natural Language Text, in: Proceedings of the 2009 International Conference on							
921		Computational Intelligence and Software Engineering, Institute of Electrical and Electronics							
922		Engineers (IEEE), Washington, DC, 2009, pp. 1-4. https://doi.org/10.1109/CISE.2009.5363900.							
923	[79]	H. Cunningham, GATE, a general architecture for text engineering, Computers and the							
924		Humanities 36(2) (2002) pp. 223-254, https://doi.org/10.1023/A:1014348124664.							
925	[80]	P. Penumatsa, M. Ventura, A.C. Graesser, M. Louwerse, X. Hu, Z. Cai, D.R. Franceschetti, The							
926		right threshold value: what is the right threshold of cosine measure when using latent semantic							
927		analysis for evaluating student answers?, International Journal on Artificial Intelligence Tools							
928		15 (05) (2006) pp. 767-777, https://doi.org/10.1142/S021821300600293X.							

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

[81]	N. Rekabsaz, M. Lupu, A. Hanbury, Exploration of a threshold for similarity based on							
	uncertainty in word embedding, in: Proceedings of the European Conference on Information							
	Retrieval, Springer, Scotland, UK, 2017, pp. 396-409. https://doi.org/10.1007/978-3-319							
	56608-5_31.							
[82]	S. Cahyono, Comparison of document similarity measurements in scientific writing using Jaro-							
	Winkler Distance method and Paragraph Vector method, in: Proceedings of the Institute of							
	Physics (IOP)Conference Series: Materials Science and Engineering, (662), Institute of Physics							
	(IOP)Publishing, London, UK, 2019, p. 052016. https://doi.org/10.1088/1757-							
	899X/662/5/052016.							
[83]	I. Kloo, M.F. Dabkowski, S.H. Huddleston, Improving record linkage for counter-threat finance							
	intelligence with dynamic Jaro-Winkler thresholds, in: Proceedings of the2019 Winter							
	Simulation Conference (WSC), Institute of Electrical and Electronics Engineers (IEEE),							
	National Harbor, Maryland, 2019, pp. 2467-2478.							
	https://doi.org/10.1109/WSC40007.2019.9004945.							
[84]	Montrose County, , Construction Contract,							
	https://www.montrosecounty.net/DocumentCenter/View/823/Sample-Construction-Contract,							
	Last Access: December 21 <sup>st</sup> , 2021.							
[85]	Legal Templates, Create a free construction contract							
	agreement, https://legaltemplates.net/form/lt/construction-contract-							
	agreement/?utm_source=google&utm_medium=cpc&utm_campaign=81-							
	Construction+Contractor+Agreement&utm_term=free+construction+contracts&utm_campaig							
	n_id=806343792&utm_adgroup=Construction+Contractor+Agreement+-							
	+Simplified&utm_adgroup_id=44606415274&utm_content=191412796236&gclid=Cj0KCQi							
	Ak4aOBhCTARIsAFWFP9FpR07y2Zo4j5TNjPUCraxOFmc68kb3LNpChyqYuDfo6jMpCY							
	UJIMaAqslEALw_wcB, Last Access: December 21 <sup>st</sup> , 2021.							
[86]	K. Potts, Model construction contract for homeowners https://buildingadvisor.com/project-							
	management/contracts/model-construction-contract_1/, Last Access: December 21st, 2021.							

# Appendix A: Patterns Used in Expanded Pattern Matching-

#### based Rules 959

958

[complementary subject, candidate subject, candidate compliance checking 960 1. attribute], inter clause boundary relation, [complementary subject, candidate 961 subject, candidate compliance checking attribute], indicating "part of" or 962 "belongs to" relation by the term "of", [complementary subject, candidate subject, 963 candidate compliance checking attribute], Conjunctive Term, [complementary 964 subject, candidate subject, candidate compliance checking attribute]. 965

- 2. candidate subject, preposition, complementary subject, inter clause boundary
   relation, candidate subject, adjective, preposition, candidate subject.
- 3. [candidate subject, complementary subject, comparative relation], inter clause
  boundary relation, gerund or present participle verb, [candidate subject,
  complementary subject, comparative relation].
- 971 4. [complementary subject, candidate subject, candidate compliance checking
  972 attribute], past participle verb, comparative relation, value, unit, conjunctive term,
  973 comparative relation, value, unit.
- 5. complementary subject, modal verb, base form verb, value, unit, comparative
  relation, comparative relation, conjunctive term, value, unit, adjective, preposition,
  candidate subject.
- 977 6. [complementary subject, candidate subject, candidate compliance checking
  978 attribute], modal verb, negation, base form verb, comparative relation, value, unit,
  979 preposition, complementary subject.
- 980 7. [candidate subject, complementary subject, comparative relation], gerund or
  981 present participle verb, inter clause boundary relation, [candidate subject,
  982 complementary subject, comparative relation].
- 8. [candidate subject, complementary subject], inter clause boundary relation,
  [candidate subject, complementary subject], inter clause boundary relation
- 985
  98. [candidate subject, complementary subject], slash "/", [candidate subject, complementary subject].
- 10. candidate compliance checking attribute, indicating "part\_of" or "belongs\_to"
  relation by the term "of", for each, candidate subject.
- 989 11. candidate subject, relation verb, inter clause boundary relation, candidate subject
- 12. candidate compliance checking attribute, modal verb, base form verb, negation,
   comparative relation, candidate compliance checking attribute.
- 13. value, complementary subject, preposition, for each, value, unit, indicating
  "part\_of" or "belongs\_to" relation by the term "of", candidate compliance
  checking attribute.
- 995 14. [ candidate subject, complementary subject, candidate compliance checking
  996 attribute], preposition, for each..
- 997 15. [complementary subject, candidate subject, candidate compliance checking
  998 attribute], [non-3rd person singular present verb, modal verb, base form
  999 verb],possessive subject restriction, value,[complementary subject, candidate
  1000 subject, candidate compliance checking attribute].
- 1001 16. [complementary subject, candidate subject, candidate compliance checking
  1002 attribute], [non-3rd person singular present verb, base form verb, 3rd person
  1003 singular present verb], comparative relation, value, [complementary subject,
  1004 candidate subject, candidate compliance checking attribute].
- 1005 17. [candidate subject, complementary subject, comparative relation], inter clause

- boundary relation, conjunctive term, [candidate subject, complementary subject,
  comparative relation].
- 1008 18. complementary subject, preposition, complementary subject, modal verb, negation,
  1009 base form verb, candidate compliance checking attribute.
- 1010 19. [complementary subject, candidate subject, candidate compliance checking
  1011 attribute], indicating "part\_of" or "belongs\_to" relation by the term "of",
  1012 complementary subject, non-3rd person singular present verb,3rd person singular
  1013 present verb, negation, comparative relation, value, unit.
- 1014 20. [candidate subject, complementary subject, candidate compliance checking 1015 attribute], for "with" or "with in" relation, [candidate subject, complementary 1016 subject, candidate compliance checking attribute].
- 1017 21. base form verb, preposition, [candidate subject, complementary subject, candidate
  1018 compliance checking attribute.
- 1019 22. candidate subject, modal verb, negation, base form verb, candidate subject.
- 1020 23. [candidate subject, candidate compliance checking attribute], relation verb,
  1021 [complementary subject, candidate compliance checking attribute], inter clause
  1022 boundary relation, complementary subject.
- 1023 24. [candidate subject, complementary subject, comparative relation], indicating
  1024 "part\_of" or "belongs\_to" relation by the term "of", for each, [candidate subject,
  1025 complementary subject, comparative relation].
- 1026 25. complementary subject, character, cardinal number.
- 1027 26. [candidate subject, candidate compliance checking attribute], indicating "part\_of"
  1028 or "belongs\_to" relation by the term "of", value, unit, adjective.
- 27. candidate compliance checking attribute, gerund or present participle verb, inter
   clause boundary relation, [candidate subject, complementary subject, comparative
   relation, candidate compliance checking attribute].
- 28. [complementary subject, candidate subject, candidate compliance checking
  attribute], preposition, [complementary subject, candidate subject, candidate
  compliance checking attribute], [preposition, the word "to"], [complementary
  subject, candidate subject, candidate compliance checking attribute].
- 1036 29. complementary subject, modal verb, base form verb, preposition, candidate subject,
  1037 value.
- 30. candidate compliance checking attribute, modal verb, negation, base form verb, the
  word "to", candidate subject.
- 1040 31. [complementary subject, candidate subject, candidate compliance checking 1041 attribute], relation verb, value, unit, conjunctive term, value, unit.
- 1042 32. complementary subject, modal verb, base form verb, negation, comparative
   1043 relation, value, unit, preposition, candidate compliance checking attribute.
- 1044 33. candidate compliance checking attribute, conjunctive term, past participle verb,
   1045 candidate compliance checking attribute.

1046	34.	[candidate subject, complementary subject, candidate compliance checking
1047		attribute, comparative relation], 3rd person singular present verb, past participle
1048		verb.
1049	35.	[complementary subject, candidate compliance checking attribute], modal verb,
1050		base form verb, relation verb, [complementary subject, candidate subject,
1051		candidate compliance checking attribute.
1052	36.	candidate subject, modal verb, negation, base form verb, candidate compliance
1053		checking attribute.
1054	37.	preposition, value, unit, candidate subject.
1055	38.	modal verb, negation, base form verb, [adjective, past participle verb, past tense
1056		verb]
1057	39.	value, unit, preposition, candidate subject.
1058	40.	preposition, value, unit, indicating "part_of" or "belongs_to" relation by the term
1059		"of", candidate subject.
1060	41.	relation verb, candidate compliance checking attribute, indicating "part_of" or
1061		"belongs_to" relation by the term "of", cardinal number, conjunctive term,
1062		comparative adjective.
1063	42.	preposition, past participle verb, candidate compliance checking attribute.
1064	43.	complementary subject, the word "to", value, complementary subject.
1065	44.	negation, comparative relation, value, slash "/", unit, candidate compliance
1066		checking attribute.
1067	45.	candidate subject, possessive subject restriction.
1068	46.	complementary subject, candidate compliance checking attribute.
1069	47.	preposition, value, unit, candidate subject, indicating "part_of" or "belongs_to"
1070		relation by the term "of", candidate compliance checking attribute.
1071	48.	complementary subject, modal verb, possessive subject restriction, complementary
1072		subject.
1073	49.	complementary subject, candidate subject, candidate compliance checking
1074		attribute], value, conjunctive term, comparative adjective, [complementary subject,
1075		candidate subject, candidate compliance checking attribute].
1076	50.	adjective, indicating "part_of" or "belongs_to" relation by the term "of",
1077		[candidate subject, complementary subject, candidate compliance checking
1078		attribute].
1079	51.	preposition, value, unit, preposition, candidate compliance checking attribute.
1080	52.	candidate compliance checking attribute, indicating "part_of" or "belongs_to"
1081		relation by the term "of", value, unit, the word "to", value, unit.
1082	53.	[candidate subject, complementary subject, candidate compliance checking
1083		attribute, comparative relation], indicating "part_of" or "belongs_to" relation by
1084		the term "of", gerund or present participle verb, [candidate subject, complementary
1085		subject, candidate compliance checking attribute, comparative relation].

For the final published version, please find it at the Elsevier Database Here: https://www.sciencedirect.com/science/article/abs/pii/S0926580522001030?via%3Dihub

	1086	54.	comparative	relation,	value,	[candidate	subject,	complementary	subject].
--	------	-----	-------------	-----------	--------	------------	----------	---------------	-----------

- 1087 55. [candidate subject, complementary subject, candidate compliance checking
  1088 attribute, comparative relation], indicating "part\_of" or "belongs\_to" relation by
  1089 the term "of", comparative adjective, [candidate subject, complementary subject,
  1090 candidate compliance checking attribute, comparative relation].
- 1091 56. complementary subject, candidate subject, candidate compliance checking
  1092 attribute], negation, comparative relation, value, [complementary subject,
  1093 candidate subject, candidate compliance checking attribute].
- 1094 57. candidate subject, gerund or present participle verb, candidate compliance
  1095 checking attribute, indicating "part\_of" or "belongs\_to" relation by the term "of",
  1096 comparative relation, value.
- 1097 58. negation, comparative relation, value, indicating "part\_of" or "belongs\_to" relation
  1098 by the term "of", candidate compliance checking attribute, indicating "part\_of" or
  1099 "belongs\_to" relation by the term "of", complementary subject.
- 1100 59. candidate compliance checking attribute, modal verb, base form verb, past1101 participle verb.
- 60. negation, comparative relation, value, indicating "part\_of" or "belongs\_to" relation
  by the term "of", candidate compliance checking attribute.
- 1104 61. candidate compliance checking attribute, 3rd person singular present verb,
   1105 negation, comparative relation, value, unit.

1106 62. candidate subject, comparative relation, value, preposition, complementary subject.

- 63. negation, possessive subject restriction, comparative relation, value, unit,
  indicating "part\_of" or "belongs\_to" relation by the term "of", candidate
  compliance checking attribute.
- 1110 64. candidate subject, preposition, comparative relation, value, unit.

1111

1112 Note: A pair of square brackets encloses different options that can be used in that specific slot of the 1113 pattern.