

1 Knowledge Acquisition and Question Answering 2 via Controlled Natural Language

3 **Tiantian Gao**

4 Department of Computer Science, Stony Brook University
5 [Stony Brook, NY, USA]
6 tiagao@cs.stonybrook.edu

7 — Abstract —

8 Knowledge acquisition from text is process of automatically acquiring, organizing and structur-
9 ing knowledge from text which can be used to perform question answering or complex reasoning.
10 However, current state-of-the-art systems are limited by the fact that they are not able to con-
11 struct the knowledge base with high quality as knowledge representation and reasoning (KRR)
12 has a keen requirement for the accuracy of data. Controlled Natural Languages (CNLs) emerged
13 as a way to improve the quality problem by restricting the flexibility of the language. However,
14 they still fail to do so as sentences that express the same information may be represented by
15 different forms. Current CNL systems have limited power to standardize the sentences into the
16 same logical form. We solved this problem by building the Knowledge Acquisition Logic Machine
17 (KALM), which performs semantic analysis of English sentences and achieves superior accuracy
18 of standardizing sentences that express the same meaning to the same logical representation. Be-
19 sides, we developed the query part of KALM to perform question answering, which also achieves
20 very high accuracy in query understanding.

21 **2012 ACM Subject Classification** I.2.1 Applications and Expert Systems

22 **Keywords and phrases** Knowledge Acquisition, Question Answering, Controlled Natural Lan-
23 guage

24 **Digital Object Identifier** 10.4230/OASISs.ICLP DC.2018.23

25 **1** Introduction

26 Knowledge acquisition is the process of extracting, organizing, and structuring knowledge from
27 data sources such that the constructed knowledge base can be used for question answering
28 or performing complex reasoning. Traditional ways of knowledge acquisition largely reply
29 on domain experts to encode the knowledge base in rule-based systems such as XSB [12]
30 and Clingo [4]. However, this requires too much domain specific knowledge and eligible
31 engineers are in very short supply. Information extraction systems emerged as the tools to
32 extract knowledge frame text (i.e., OpenIE [1], SEMAFOR [2], Stanford CoreNLP/KBP
33 [8], SLING [10]). They achieved admirable results in processing free text, however, their
34 accuracy is far from meeting the requirement of knowledge acquisition. In addition, they are
35 only designed to extract the knowledge from text, but not intended to represent it in a way
36 suitable for reasoning. Controlled Natural Languages (CNLs) [7] emerged as a technology
37 that bridges this gap. Representative systems include Attempto Controlled English (ACE) [3]
38 and Processable English (PENG) [11]. They are designed to process English sentences with
39 restricted grammar but unambiguous interpretations and translate the sentences into logic
40 for reasoning. The main issue with CNLs is that they have limited power of standardizing
41 sentences that express the same information via different syntactic forms into the same logical
42 representation. For instance, the sentences *a customer buys a phone*, *a customer makes*
43 *a purchase of a phone*, *a customer is a buyer of a phone* are mapped to different logical



© Tiantian Gao;
licensed under Creative Commons License CC-BY
14th Doctoral Consortium (DC) on Logic Programming.

Editors: Neda Saeedloei and Paul Fodor; Article No. 23; pp. 23:1–23:8

OpenAccess Series in Informatics

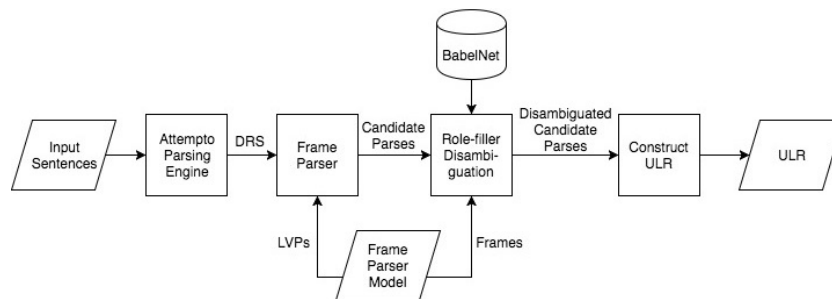


OASIS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

44 representations. Therefore, they are not suffice for question answering or complex logical
45 reasoning.

46 In this work, we build Knowledge Acquisition Logic Machine (KALM), which conducts
47 semantic analysis of CNL sentences and achieves superior accuracy of standardizing English
48 sentences that express the same information via different forms to the same logical form.
49 The system is built based on utilizing linguistic knowledge bases (BabelNet [9] and FrameNet
50 [5]) and our frame-based parsing and disambiguation algorithms. Besides, we developed
51 the query part of KALM which supports high accuracy query parsing and answer retrieval.
52 The following is organized as follows: Section 2 describes the KALM system for knowledge
53 acquisition, Section 3 describes the query part of KALM, Section 4 shows the evaluation
54 results of KALM, Section 5 discusses the next steps of work, Section 6 concludes the paper.

55 2 Knowledge Acquisition Logic Machine (KALM)



■ Figure 1 Pipeline for translating a sentence into ULR

56 Figure 1 shows the pipeline of KALM that translates a CNL sentence into *unique logical*
57 *representation (ULR)*, the semantic form of CNL sentences. The KALM framework consists
58 of five components:

59 **Syntactic Parsing.** We use Attempto Parsing Engine (APE)¹ to parse CNL sentences
60 and translate them into Discourse Representation Structure (DRS) [6], which represents
61 the syntactic and dependency information of the sentences. DRS relies on 7 predicates:
62 `object/6`, `predicate/4`, `property/3`, `modifier_adv/3`, `modifier_pp/3`, `relation/3`, and
63 `has_part/2`. For example, the `object-predicate` represents an entity which corresponds to
64 a noun word in the sentences. A `predicate-predicate` represents an event and the subject
65 and object of the events. `predicate-predicate` corresponds to a verb word in a sentence. For
66 example, given the sentence *A customer buys a phone*, it is parsed into DRS as

```
67     object(A, customer, countable, na, eq, 1)
68     object(B, phone, countable, na, eq, 1)
69     predicate(C, buy, A, B)
```

70 **Frame-based Parsing.** Based on the DRS, the frame-based parser generates a list of
71 *candidate parses*, which represent the frame-semantic relations the sentences entail. For
72 instance, given the sentence *A customer buys a phone*, the frame-based parser generates
73 the following parse result: `Frame(Commerce_Buy, Roles: Buyer = customer, Goods =`
74 `phone)`. The parse says there are two entities: *customer* and *phone*, which are involved in

¹ <https://github.com/Attempto/APE>

75 the `Commerce_Buy` relation. The *customer* serves as the `Buyer` role of this frame relation and
 76 *phone* serves as the `Goods` role of the frame relation. The parser is constructed based on two
 77 components: *logical frames* and *logical valence patterns (lvps)*. The logical frames represent
 78 the definition of the frame relations via Prolog facts. For instance, the `Commerce_Buy`² frame
 79 is represented as

```
80 fp(Commerce_Buy, [
81     role(Buyer, [bn:00014332n], []),
82     role(Seller, [bn:00053479n], []),
83     role(Goods, [bn:00006126n, bn:00021045n], []),
84     role(Recipient, [bn:00066495n], []),
85     role(Money, [bn:00017803n], [currency])]).
```

86 where for each `role-term`, the first argument represents the name of the *frame role*, the
 87 second argument represents the BabelNet synsets associated which capture the meaning of
 88 the role, and the third argument specifies some data type constraints. The lvps represent
 89 the grammatical context of a sentence that could potentially entail a frame. Consider the
 90 following lvp for extracting an instance of the `Commerce_Buy` frame:

```
91 lvp(buy, v, Commerce_Buy, [
92     pattern(Buyer, verb->subject, required),
93     pattern(Goods, verb->object, required),
94     pattern(Recipient, verb->pp(for)->dep, optnl),
95     pattern(Money, verb->pp(for)->dep, optnl),
96     pattern(Seller, verb->pp(from)->dep, optnl)]).
```

97 The first and second arguments represents a *lexical unit* (a word + *part-of-speech*) that
 98 could trigger an instance of the `Commerce_Buy` frame. Next, it comes with a list of `pattern-`
 99 `terms`, each represents the syntactical context between the lexical unit, frame role, and
 100 the actual *role-filler* word. The lvps are generated automatically by KALM based on the
 101 annotated training sentences, which contains the frame name, lexical unit, and frame elements
 102 information. When a new sentence comes, we check every word in the sentence and find
 103 whether there exists any lvp whose lexical unit matches the chosen word. If so, we apply the
 104 lvp to the sentence and extract an instance of the frame from the sentence.

105 **Role-Filler Disambiguation.** Doing frame-based parsing is not enough because the
 106 aforementioned frame-based parsing only replies the grammatical information of the sentence.
 107 This way of parsing may generate candidate parses that misidentify the frames, role-filler
 108 words, or assign the wrong roles to the role-filler words. To rule out the wrong parses,
 109 we perform role-filler disambiguation which checks whether the extracted role-filler words
 110 are semantically compatible with the frame roles. For each role-filler and role pair, we
 111 compute a semantic score. Based on the scores for the role-filler and role pairs, we score the
 112 entire candidate parse and removes the ones that falls below a threshold. To compute the
 113 semantic score, we first query the role-filler word against BabelNet and get a list of associated
 114 BabelNet synsets (called *candidate role-filler synsets*). Then, we traverse BabelNet semantic
 115 network and measure the semantic similarity between each candidate role-filler synset and
 116 the corresponding role-synset. Basically, we consider all semantic paths that connect the
 117 synset pair, and then use a heuristic scoring function to score the path. The candidate

² https://framenet2.icsi.berkeley.edu/fnReports/data/frame/Commerce_buy.xml

118 role-filler synset which achieves the highest semantic score is chosen and assigned as the
 119 disambiguated role-filler synset for the respective role-filler word.

120 **Translating into ULR.** Based on the *disambiguated candidate parses* generated from the
 121 role-filler disambiguation step, we translate the parses into ULR. ULR uses `frame/2` and
 122 `role/2` predicates to represent instances of frames and roles. ULR uses `synset/2` and `text/2`
 123 predicates to represent the synset and text information for the role-filler words. For example,
 124 the sentence *a customer buys a phone* is translated into ULR as

```
125     frame(id_1, Commerce_buy).
126     role(id_1, Buyer, id_2).
127     role(id_1, Goods, id_3).
128     synset(id_2, bn:00022095n). % customer synset
129     text(id_2, customer).
130     synset(id_3, bn:00062020n). % phone synset
131     text(id_3, phone).
```

132 3 Question Answering

133 3.1 Issues in CNL-based Queries

134 The ACE query language³ supports two types of queries: *true/false*- and *wh*-queries where
 135 the query words include *who*, *where*, *what*, and so on. A *true/false*-query is translated into
 136 DRS the same way as a definite sentence does. For *wh*-queries, APE uses a special predicate
 137 `query/2` to represent the *wh*-words. For instance, the query *who buys what?* is represented
 138 in DRS as

```
139     query(A,who)-1/1
140     query(B,what)-1/3
141     predicate(C,buy,A,B)-1/2
```

142 where the variables *who* and *what* are captured by the `query`-predicate.

143 However, APE only does shallow syntactic analysis of a query. There are a few issues
 144 to solve before we can precisely capture the meaning of a query and acquire the intended
 145 knowledge. Consider the following query sentences:

- 146 1. *Mary buys which car?*
- 147 2. *Who buys IBM's stocks?*
- 148 3. *Which person buys which car in which place at which price?*
- 149 4. *A \$person buys a \$car in a \$place at a \$price.*

150 First of all, as a *wh*-variable is a placeholder for the entities to be shown in the output
 151 result, the type of entities the variable represents must be disambiguated and also used for
 152 acquiring the related information. As shown in Sentence (1) in the above example, we need
 153 to identify that the type of the entities associated with the *which*-variable is a *car*. Therefore,
 154 if we know *Mary* buys both a *Camry* and a *pen*, only *Camry* should be returned.

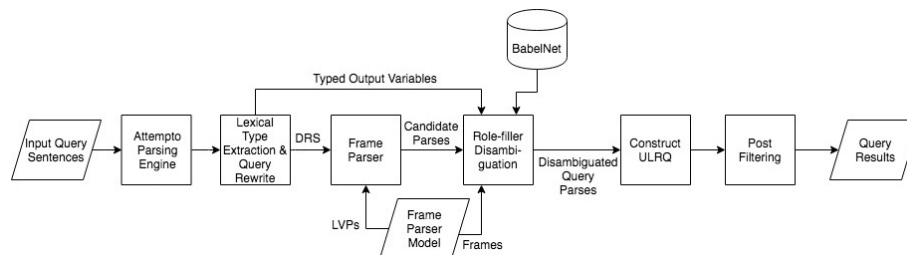
155 Second, ACE's query language has constrained power of denoting types in the query. As
 156 shown in Sentence (2) in the above example, *emphWho* could refer to either a *company* or a
 157 *person*. However, it is ambiguous whether the user intends to acquire *company* or *person*
 158 entities or both. One solution to that is to use the query word *which* and rewrite the sentence

³ http://attempto.ifi.uzh.ch/site/docs/ace/6.7/ace_constructionrules.html

159 as *Which person buys IBM's stocks* if the user intends to acquire *person* entities. However,
 160 the sentence may become cumbersome when there are many such typed variables as shown
 161 in Sentence (3). To solve this problem, we introduce *typed output variables* in the query
 162 language as the form $\$type$. Hence, Sentence (3) will be rewritten to Sentence (4) which is
 163 expressed in a more precise and concise way.

164 Third, to acquire the associated instances of frames from the knowledge base which is
 165 constructed in the knowledge acquisition phase, we also need to perform frame-semantic
 166 parsing based on queries. However, as shown in the previous example, the DRS for query is
 167 not exactly the same as the DRS used to represent definite sentences. Therefore, the existing
 168 lyps are not applicable for parsing queries. One way to solve this problem is to construct
 169 an additional set of training sentences for queries. However, this will requires a lot of work.
 170 Besides, since FrameNet doesn't contain any sentences related to queries, it would require a
 171 lot of manual work to construct CNL queries. To solve this issue, we perform a DRS rewrite
 172 to queries such that we can reuse the existing lyps for definite sentences to parse queries.

173 3.2 Question Answering



■ **Figure 2** Pipeline for translating a query into ULRQ and answer retrieval and filtering

174 Figure 2 shows the pipeline that translates a CNL query into the logical form, *Unique*
 175 *Logical Representation for Queries (ULRQ)*, which is used to query the knowledge base to
 176 retrieve the answers. The question answering part consists of the following components:

177 **Syntactic Parsing.** This is the same as the knowledge acquisition part.

178 **Query Parsing.** We also perform frame-based parsing to generate several candidate parses
 179 which represent the frame relations the query belongs to. However, as mentioned in the
 180 previous subsection, the DRS for queries are different from the DRS for definite sentences.
 181 Therefore, we perform a DRS adaptation of the DRS corresponding to the query such that
 182 the existing lyps for definite sentences can be reused to do frame-based parsing for queries.
 183 Besides, we perform a syntactic analysis of the queries and identify the lexical types of the
 184 query words (e.g., *which*).

185 **Role-Filler Disambiguation.** This is the same as the knowledge acquisition part.

186 **Translating Queries into ULRQ.** Queries are represented in a similar way as definite
 187 sentences except that we use logical variables to denote instances of frames and roles. For
 188 instance, the query *Who buys a phone?* is translated into ULRQ as

```
189 ?- frame(FrameV, 'Commerce_Buy'),
190     role(FrameV, 'Buyer', BuyV), synset(BuyV, BuyerRoleFillerOutV),
191     role(FrameV, 'Goods', GoodV), synset(GoodV, GoodsRoleFillerOutV),
192     check_type(BuyerRoleFillerOutV, bn:00046516n), % person synset
193     check_type(GoodsRoleFillerOutV, bn:00062020n). % phone synset
```

194 **Type Filtering of Query Results.** As is shown from the ULRQ above, the clauses from
 195 lines 1-3 retrieves all instances of frames and the associated roles from the knowledge base.
 196 However, not all role-fillers for **Buyer** and **Goods** may be related to *person* and *phone*. To
 197 rule out the unrelated ones, we perform type filtering of the query results, which calls the
 198 `check_type` predicate in the above ULRQ.

199 **4 Evaluation**

200 At present, KALM contains 50 logical frames with 213 logical valance patterns. We use the
 201 following metrics to measure the performance of the system:

FrSynC	all frames, roles & output variables are identified correctly; all role-filler words & variable types are disambiguated correctly
FrC	all frames, roles and output variables are identified correctly, but some disambiguation mistakes
Wrong	some frames, roles or output variables are misidentified

203 For knowledge acquisition, we achieve an accuracy of 95.6% (FrSynC). This accuracy is
 204 far from the state-of-the-art information extraction systems including SEMAFOR, SLING,
 205 and Stanford CoreNLP. For understanding of the queries, we achieve an accuracy of 94.49%
 206 (FrSynC).

207 **5 Next Steps**

208 The current work focuses on acquisition of definite knowledge from CNL sentences and
 209 question answering. The next step is to acquire rules from CNL sentences and perform more
 210 complex reasoning. This not only requires parsing individual sentences correctly, but also
 211 requires multi-sentence parsing and information in different sentences must be related to
 212 each other properly. This goes well beyond anaphora resolution, which ACE is already able
 213 to handle.

- 214 1. *Every bird is an animal.*
- 215 2. *Every bird flies.*
- 216 3. *Stella is a sea eagle.*
- 217 4. *Penguins do not fly.*
- 218 5. *A violet is not an animal.*
- 219 6. *Sparrow Daffy doesn't fly.*

220 Consider the above example: Sentences (1) and (2) denote rules which say that if we know
 221 there is a *bird*, we can infer the bird is an animal and flies. Therefore, based on Sentences (1),
 222 (2) and (3), we can infer *Stella* is an animal and flies. However, this doesn't hold for *Tweety*
 223 because Sentence (4) is an exception to Sentence (2) and therefore refutes any conclusion
 224 derived from Sentence (2). Moreover, based on Sentence (1) and (5), since a violet is not an
 225 animal, we conclude that a violet is not a bird. But, this way of reasoning is not desired for
 226 Sentence (2) and (6) because Daffy may be injured and therefore not being able to fly.

227 To precisely capture the meaning of rules in CNL and perform reasoning for the above
 228 cases, the research issues are three-fold: the first issue is the development of CNL extensions
 229 that are suitable for representing rules and inter-sentence dependencies/references. For
 230 instance, in addition to the form "*every . . .*" as shown in Sentence (1) and (2), we can also
 231 use an "*if . . . then*" statement to represent a rule. Besides, we need a mechanism to indicate
 232 the inter-sentence dependencies as shown in Sentence (1) and (4) where Sentence (4) is an

233 exception case to Sentence (2). This could be done either by specifying the inter-sentence
 234 dependencies explicitly or by an automatic mechanism to recognize these dependencies
 235 without explicit mentioning.

236 The second issue is the actual nature of the logic to be used for capturing rules. As
 237 shown in the above example, when there is a fact that a *violet* is not an animal, it is natural
 238 to infer that it is not an animal. But, it is not reasonable to infer the Daffy is not a bird
 239 if Daffy doesn't fly. To distinguish the differences, we can use a *first-order logic* rule to
 240 represent Sentence (1) where *contrapositive inference* is desired and use a *Prolog* rule to
 241 represent Sentence (2) where contrapositive inference is not required. As to the inter-sentence
 242 dependency between Sentence (2) and (4), we believe *defeasible logic* is a good fit. Basically,
 243 rules have priorities in defeasible logic where the rule with a higher priority can defeat a default
 244 rule which has a lower priority. For the above example, we label Sentence (4) as a rule with
 245 a higher priority than the rule corresponding Sentence (2) and also defeat the low priority
 246 rule when incompatible conclusions are derived.

247 The third issue is standardization. Same as knowledge acquisition for definite sentences
 248 and queries, we will also standardize rules that express the same meaning via different
 249 syntactic forms.

250 6 Conclusion

251 In this work, we described the KALM system, which achieves superior accuracy in knowledge
 252 acquisition and question answering. The system is built on our frame-based parsing and
 253 disambiguation algorithms and the use of external linguistic knowledge bases including
 254 BabelNet and FrameNet. As the next step, we plan to work on extracting rules from
 255 sentences and perform common sense reasoning.

256 ——— References ———

- 257 1 Gabor Angeli, Melvin Jose Johnson Premkumar, and Christopher D. Manning. Leveraging
 258 linguistic structure for open domain information extraction. In *53rd Annual Meeting of the*
 259 *Association for Computational Linguistics*, pages 344–354, Beijing, China, 2015.
- 260 2 Dipanjan Das, Desai Chen, André F. T. Martins, Nathan Schneider, and Noah A. Smith.
 261 Frame-semantic parsing. *Comp, Linguistics*, 40(1):9–56, 2014.
- 262 3 Norbert E. Fuchs, Kaarel Kaljurand, and Tobias Kuhn. Attempto controlled english for
 263 knowledge representation. In *Reasoning Web, 4th Intl. Summer School, Sept. 7-11*, pages
 264 104–124, Venice, Italy, 2008.
- 265 4 Martin Gebser, Benjamin Kaufmann, Roland Kaminski, Max Ostrowski, Torsten Schaub,
 266 and Marius Thomas Schneider. Potassco: The potsdam answer set solving collection. *AI*
 267 *Commun.*, 24(2):107–124, 2011.
- 268 5 Chrstopher R. Johnson, Charles J. Fillmore, Miriam R.L. Petruck, Collin F. Baker, Mi-
 269 chael J. Ellsworth, Josef Ruppenhofer, and Esther J. Wood. FrameNet: Theory and Prac-
 270 tice, 2002.
- 271 6 Hans Kamp and Uwe Reyle. *From discourse to logic: Introduction to modeltheoretic se-*
 272 *mantics of natural language, formal logic and discourse representation theory*, volume 42.
 273 Springer Science & Business Media, 2013.
- 274 7 Tobias Kuhn. A survey and classification of controlled natural languages. *Comp. Linguistics*,
 275 40(1):121–170, 2014.
- 276 8 Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard,
 277 and David McClosky. The Stanford CoreNLP natural language processing toolkit. In

- 278 *Association for Computational Linguistics (ACL) System Demonstrations*, pages 55–60,
279 2014. URL: <http://www.aclweb.org/anthology/P/P14/P14-5010>.
- 280 **9** Roberto Navigli and Simone Paolo Ponzetto. BabelNet: The automatic construction, eval-
281 uation and application of a wide-coverage multilingual semantic network. *Artificial Intelli-*
282 *gence*, 193:217–250, 2012.
- 283 **10** Michael Ringgaard, Rahul Gupta, and Fernando C. N. Pereira. SLING: A framework for
284 frame semantic parsing. *CoRR*, 1710.07032:1–9, 2017. URL: [http://arxiv.org/abs/1710.](http://arxiv.org/abs/1710.07032)
285 [07032](http://arxiv.org/abs/1710.07032), [arXiv:1710.07032](http://arxiv.org/abs/1710.07032).
- 286 **11** Rolf Schwitter. English as a formal specification language. In *13th Intl. Workshop on*
287 *Database and Expert Systems Appl. (DEXA 2002)*, pages 228–232, Aix-en-Provence, France,
288 2002.
- 289 **12** T. Swift and D.S. Warren. XSB: Extending the power of prolog using tabling. *Theory and*
290 *Practice of Logic Programming*, 2011.