# Logical English to annotate metadata and augment semantics descriptions

Jacinto Dávila[1], Miguel Calejo[2] and Andrew Noble[3]

## Introduction

We present Logical English (Kowalski et al, 2023), LE, as a suitable alternative to annotate metadata of CSV datasets and augment semantics descriptions. Logical English is a controlled natural language, designed as syntactic sugar for logic programming languages such as Prolog and Datalog. It can used to write documents with a formal structure that can be mapped to computer code, while preserving the readability of a text in English.

For the purposes of the 2025 OMG Semantic Augmentation Challenge[4], we have prepared an example of a LE for the given dataset FDIC Insured Banks, containing descriptions, in clear English, for each label of the CSV files (taken from the ArcGIS platform) and a collection of statements with the connections between those labels and a number of ontologies, including the suggested (see Figures 1 and 2).

```
1   the target language is: prolog.
2
3   the templates are:
4   *a label* is a label for *a description*.
5   the dataset ID is *a dataset*.
6   the dataset *a dataset* is described at *a location*.
7   the dataset *a dataset* can be downloaded from *a url*.
8   the provider of *a dataset* is *a provider*.
9   the platform of *a dataset* is *a platform*.
10  the method of delivery of *a dataset* is *a method*.
11  the product version of *a dataset* included here is *a date*.
12
13  the ontology is:
14
15  the dataset ID is FDIC Insured Banks.
16  %the dataset FDIC Insured Banks is described at https://hub.arcgis.com/datasets/geoplatform::fdic-insured-banks/about .
17  %the dataset FDIC Insured Banks can be downloaded from https://hub.arcgis.com/datasets/geoplatform::fdic-insured-banks/ .
18  the provider of FDIC Insured Banks is GeoPlatform ArcGIS Online.
19  the platform of FDIC Insured Banks is ArcGIS.
20  the method of delivery of FDIC Insured Banks is http.
21  the product version of FDIC Insured Banks included here is 2018-06-29:00:00.
22
23  %Column Label to Meaning Mappings:
24  X is a label for Latitude.
25  Y is a label for Longitude.
26  OBJECTID is a label for ID.
27  ACQDATE is a label for Acquisition Date.
28  ADDRESS is a label for Branch Address.
29  ADDRESS2 is a label for Street Address Line 2.
30  BKCLASS is a label for Institution Class.
31  CBSA is a label for Core Based Statistical Areas (Branch).
32  CBSA_DIV is a label for Metropolitan Divisions Name (Branch).
33  CBSA_DIV_FLG is a label for Metropolitan Divisions Flag (Branch).
34  CBSA_DIV_NO is a label for Metropolitan Divisions Number (Branch).
35  CBSA_METRO is a label for Metropolitan Division Number (Branch).
36  CBSA_METRO_FLG is a label for Metropolitan Division Flag (Branch).
37  CBSA_METRO_NAME is a label for Metropolitan Division Name (Branch).
38  CBSA_MICRO_FLG is a label for Micropolitan Division Flag (Branch).
39  CBSA_NO is a label for Core Based Statistical Area Name (Branch).
40  CERT is a label for Institution FDIC Certificate #. CITY is a label for Branch City.
41  COUNTY is a label for Branch County.
42  CSA is a label for Combined Statistical Area Name (Branch).
43  CSA_FLG is a label for Combined Statistical Area Flag (Branch).
44  CSA_NO is a label for Combined Statistical Area Number (Branch).
45  ESTYMD is a label for Branch Established Date.
46  FI_UNINUM is a label for FDIC UNINUM of the Owner Institution.
47  ID is a label for ID. LATITUDE is a label for Latitude.
48  LONGITUDE is a label for Longitude.
49  MAINOFF is a label for Main Office.
50  MDI_STATUS_CODE is a label for Minority Status Code.
51  MDI_STATUS_DESC is a label for Minority Status Description.
52  NAME is a label for Institution Name.
53  OFFNAME is a label for Office Name.
54  OFFNUM is a label for Branch Number.
55  RUNDATE is a label for Run Date.
56  SERVTYPE is a label for Service Type Code.
57  SERVTYPE_DESC is a label for Service Type Description.
58  STALP is a label for Branch State Abbreviation.
59  STCNTY is a label for State and County Number.
60  STNAME is a label for Branch State.
61  UNINUM is a label for Unique Identification Number for a Branch Office.
62  ZIP is a label for Branch Zip Code.
63
64  an object is of a type
65      if a label is a label for the type
66      and the object is of the label.  % connection to the dataset
67
```

*Figura 1: A Logical English Document as metadata*

1   jd@logicalcontracts.com, jacinto@ula.ve
2   mc@logicalcontracts.com
3   andrew@nobleaccounting.com.au
4   https://www.omg.org/events/omg-semantic-augmentation-challenge/

```
68   %Ontological Mappings:
69   %Spatial/Geographic Mappings (GeoSPARQL/GeoNames):
70
71       X is a geo:lat (latitude coordinate). % as in http://www.opengis.net/ont/geosparql#lat
72       Y is a geo:long (longitude coordinate). % as in http://www.opengis.net/ont/geosparql#long
73       LATITUDE is a geo:lat (latitude coordinate). % as in http://www.opengis.net/ont/geosparql#lat
74       LONGITUDE is a geo:long (longitude coordinate). % as in http://www.opengis.net/ont/geosparql#long
75       OBJECTID is a geo:SpatialObject. % as in http://www.opengis.net/ont/geosparql#SpatialObject
76       ADDRESS is a locn:Address. % as in http://www.w3.org/ns/locn#Address
77       ADDRESS2 is a locn:addressArea. % as in http://www.w3.org/ns/locn#addressArea
78       CITY is a gn:populatedPlace. % as in http://www.geonames.org/ontology#populatedPlace
79       COUNTY is a gn:administrativeDivision. % as in http://www.geonames.org/ontology#administrativeDivision
80       STNAME is a gn:administrativeDivision. % as in http://www.geonames.org/ontology#administrativeDivision
81       STALP is a gn:countryCode. % as in http://www.geonames.org/ontology#countryCode
82       ZIP is a locn:postCode. % as in http://www.w3.org/ns/locn#postCode
83
84   %Financial/Business Mappings (FIBO):
85       BKCLASS is a fibo-be-le-lp:BusinessEntity. % as in https://spec.edmcouncil.org/fibo/ontology/BE/LegalEntities/LegalPersons/BusinessEntity
86       NAME is a fibo-fnd-org-fm:FormalOrganization. % as in https://spec.edmcouncil.org/fibo/ontology/FND/Organizations/FormalOrganizations/FormalOrganization
87       CERT is a fibo-fnd-arr-id:Identifier. % as in https://spec.edmcouncil.org/fibo/ontology/FND/Arrangements/IdentifiersAndIndices/Identifier
88       FI_UNINUM is a fibo-fnd-arr-id:Identifier. % as in https://spec.edmcouncil.org/fibo/ontology/FND/Arrangements/IdentifiersAndIndices/Identifier
89       UNINUM is a fibo-fnd-arr-id:Identifier. % as in https://spec.edmcouncil.org/fibo/ontology/FND/Arrangements/IdentifiersAndIndices/Identifier
90       OFFNAME is a fibo-be-le-fbo:Branch. % as in https://spec.edmcouncil.org/fibo/ontology/BE/LegalEntities/FormalBusinessOrganizations/Branch
91       OFFNUM is a fibo-fnd-arr-id:Identifier. % as in https://spec.edmcouncil.org/fibo/ontology/FND/Arrangements/IdentifiersAndIndices/Identifier
92       MAINOFF is a fibo-be-le-fbo:HeadOffice. % as in https://spec.edmcouncil.org/fibo/ontology/BE/LegalEntities/FormalBusinessOrganizations/HeadOffice
93       SERVTYPE is a fibo-fbc-fct-fse:FinancialService. % as in https://spec.edmcouncil.org/fibo/ontology/FBC/FunctionalEntities/FinancialServicesEntities/FinancialService
94       SERVTYPE_DESC is a fibo-fbc-fct-fse:FinancialServiceDescription. % as in https://spec.edmcouncil.org/fibo/ontology/FBC/FunctionalEntities/FinancialServicesEntities/FinancialServiceDescription
95
96   %Temporal Mappings (W3C Time Ontology):
97       ACQDATE is a time:Instant. % as in http://www.w3.org/2006/time#Instant
98       ESTYMD is a time:Instant. % as in http://www.w3.org/2006/time#Instant
99       RUNDATE is a time:Instant. % as in http://www.w3.org/2006/time#Instant
100
101  %Statistical/Administrative Mappings (SDMX, Dublin Core):
102      CBSA is a sdmx-concept:statisticalClassification. % as in http://purl.org/linked-data/sdmx/2009/concept#statisticalClassification
103      CBSA_DIV is a sdmx-concept:statisticalClassification. % as in http://purl.org/linked-data/sdmx/2009/concept#statisticalClassification
104      CSA is a sdmx-concept:statisticalClassification. % as in http://purl.org/linked-data/sdmx/2009/concept#statisticalClassification
105      STCNTY is a sdmx-concept:administrativeClassification. % as in http://purl.org/linked-data/sdmx/2009/concept#administrativeClassification
106
107  %General Data Mappings (DCAT, Dublin Core):
108      ID is a dct:identifier. % as in http://purl.org/dc/terms/identifier
109      CBSA_DIV_FLG is a dcat:Dataset. % as in http://www.w3.org/ns/dcat#Dataset
110      CBSA_METRO_FLG is a dcat:Dataset. % as in http://www.w3.org/ns/dcat#Dataset
111      CBSA_MICRO_FLG is a dcat:Dataset. % as in http://www.w3.org/ns/dcat#Dataset
112      CSA_FLG is a dcat:Dataset. % as in http://www.w3.org/ns/dcat#Dataset
113      MDI_STATUS_CODE is a skos:Concept. % as in http://www.w3.org/2004/02/skos/core#Concept
114      MDI_STATUS_DESC is a skos:prefLabel. % as in http://www.w3.org/2004/02/skos/core#prefLabel
115
116  %Demographic/Social Mappings (FOAF, Schema.org):
117      MDI_STATUS_CODE is a schema:demographicGroup. % as in http://schema.org/demographicGroup
118      MDI_STATUS_DESC is a schema:description. % as in http://schema.org/description
119
```

*Figura 2: The Logical English Metadata includes ontological, (is a), information*

The LE document also contain a set of queries which represent typical questions that can be posted to the document, considered as a knowledge base. In this case, questions about the structure and content of the described dataset. These are only referential questions and can be modified or extended at will by any competent writer (Figure 3):

```
120   query thing is:
121     which thing is a which type.
122
123   query label is:
124     which label is a label for Institution Class.
125
126   query description is:
127     CBSA is a label for which description.
128
129   query dataset is:
130       the dataset ID is which one.
131
132   query about is:
133       the dataset FDIC Insured Banks is described at which text.
134
135   query download is:
136       the dataset FDIC Insured Banks can be downloaded from which url.
137
138   query provider is:
139       the provider of FDIC Insured Banks is which provider.
140
141   query platform is:
142       the platform of FDIC Insured Banks is which platform.
143
144   query method is:
145       the method of delivery of FDIC Insured Banks is which method.
146
147   query version is:
148       the product version of FDIC Insured Banks included here is which version.
```

*Figura 3: Prototypical Questions for the LE Metadata*

# The output generated by the format when used with the source dataset.

Whenever is queried, an LE document is translated into a target language. In this case, we selected Prolog in order to combine it with a transformation of the dataset into Datalog. This could, of course, be done in many other ways, including consulting the dataset on the flight with just in time translations.

```prolog
1  :-module('augmentedsem-prolog', []).
2  source_lang(en).
3  local_dict([the_product_version_of_included_here_is, A, B], [dataset-dataset, date-date], [the, product, version, of, A, included, here, is, B]).
4  local_dict([the_method_of_delivery_of_is, A, B], [dataset-dataset, method-method], [the, method, of, delivery, of, A, is, B]).
5  local_dict([the_dataset_can_be_downloaded_from, A, B], [dataset-dataset, url-url], [the, dataset, A, can, be, downloaded, from, B]).
6  local_dict([the_dataset_is_described_at, A, B], [dataset-dataset, location-location], [the, dataset, A, is, described, at, B]).
7  local_dict([the_provider_of_is, A, B], [dataset-dataset, provider-provider], [the, provider, of, A, is, B]).
8  local_dict([the_platform_of_is, A, B], [dataset-dataset, platform-platform], [the, platform, of, A, is, B]).
9  local_dict([is_a_label_for, A, B], [label-label, description-description], [A, is, a, label, for, B]).
10 local_dict([the_dataset_ID_is, A], [dataset-dataset], [the, dataset, 'ID', is, A]).
11 local_meta_dict([],[],[]).
12 prolog_le(verified).
13 the_dataset_ID_is('FDIC Insured Banks').
14 the_provider_of_is('FDIC Insured Banks', 'GeoPlatform ArcGIS Online').
15 the_platform_of_is('FDIC Insured Banks', 'ArcGIS').
16 the_method_of_delivery_of_is('FDIC Insured Banks', http).
17 the_product_version_of_included_here_is('FDIC Insured Banks', '2018 - 6 - 29 : 0 : 0').
18 is_a_label_for('X', 'Latitude').
19 is_a_label_for('Y', 'Longitude').
20 is_a_label_for('OBJECTID', 'ID').
21 is_a_label_for('ACQDATE', 'Acquisition Date').
22 is_a_label_for('ADDRESS', 'Branch Address').
23 is_a_label_for('ADDRESS2', 'Street Address Line 2').
24 is_a_label_for('BKCLASS', 'Institution Class').
25 is_a_label_for('CBSA', 'Core Based Statistical Areas ( Branch )').
26 is_a_label_for('CBSA_DIV', 'Metropolitan Divisions Name ( Branch )').
27 is_a_label_for('CBSA_DIV_FLG', 'Metropolitan Divisions Flag ( Branch )').
28 is_a_label_for('CBSA_DIV_NO', 'Metropolitan Divisions Number ( Branch )').
29 is_a_label_for('CBSA_METRO', 'Metropolitan Division Number ( Branch )').
30 is_a_label_for('CBSA_METRO_FLG', 'Metropolitan Division Flag ( Branch )').
31 is_a_label_for('CBSA_METRO_NAME', 'Metropolitan Division Name ( Branch )').
32 is_a_label_for('CBSA_MICRO_FLG', 'Micropolitan Division Flag ( Branch )').
33 is_a_label_for('CBSA_NO', 'Core Based Statistical Area Name ( Branch )').
34 is_a_label_for('CERT', 'Institution FDIC Certificate #').
35 is_a_label_for('CITY', 'Branch City').
36 is_a_label_for('COUNTY', 'Branch County').
37 is_a_label_for('CSA', 'Combined Statistical Area Name ( Branch )').
38 is_a_label_for('CSA_FLG', 'Combined Statistical Area Flag ( Branch )').
39 is_a_label_for('CSA_NO', 'Combined Statistical Area Number ( Branch )').
40 is_a_label_for('ESTYMD', 'Branch Established Date').
41 is_a_label_for('FI_UNINUM', 'FDIC UNINUM of the Owner Institution').
42 is_a_label_for('ID', 'ID').
43 is_a_label_for('LATITUDE', 'Latitude').
44 is_a_label_for('LONGITUDE', 'Longitude').
45 is_a_label_for('MAINOFF', 'Main Office').
46 is_a_label_for('MDI_STATUS_CODE', 'Minority Status Code').
47 is_a_label_for('MDI_STATUS_DESC', 'Minority Status Description').
48 is_a_label_for('NAME', 'Institution Name').
49 is_a_label_for('OFFNAME', 'Office Name').
50 is_a_label_for('OFFNUM', 'Branch Number').
51 is_a_label_for('RUNDATE', 'Run Date').
52 is_a_label_for('SERVTYPE', 'Service Type Code').
53 is_a_label_for('SERVTYPE_DESC', 'Service Type Description').
54 is_a_label_for('STALP', 'Branch State Abbreviation').
55 is_a_label_for('STCNTY', 'State and County Number').
56 is_a_label_for('STNAME', 'Branch State').
57 is_a_label_for('UNINUM', 'Unique Identification Number for a Branch Office').
58 is_a_label_for('ZIP', 'Branch Zip Code').
```

*Figura 4: Logical English translated into Prolog/Datalog*

Figure 4 shows the first part of the file with the translation from the original LE document.

Figure 5 shows the next part which contains a partial rendition of the is_a/2 relation. Those correspond mainly to the ontological information provided by users. Some other essential components of the same relation are depicted in Prolog and are not shown to the final user, like the transitivity rule and the rule that connects the actual data in the dataset with the their types.

```
59 is_a(A, B) :-
60     is_a_label_for(C, B),
61     is_a(A, C).
62 is_a('X', 'geo : lat ( latitude coordinate )').
63 is_a('Y', 'geo : long ( longitude coordinate )').
64 is_a('LATITUDE', 'geo : lat ( latitude coordinate )').
65 is_a('LONGITUDE', 'geo : long ( longitude coordinate )').
66 is_a('OBJECTID', 'geo : SpatialObject').
67 is_a('ADDRESS', 'locn : Address').
68 is_a('ADDRESS2', 'locn : addressArea').
69 is_a('CITY', 'gn : populatedPlace').
70 is_a('COUNTY', 'gn : administrativeDivision').
71 is_a('STNAME', 'gn : administrativeDivision').
72 is_a('STALP', 'gn : countryCode').
73 is_a('ZIP', 'locn : postCode').
74 is_a('BKCLASS', 'fibo - be - le - lp : BusinessEntity').
75 is_a('NAME', 'fibo - fnd - org - fm : FormalOrganization').
76 is_a('CERT', 'fibo - fnd - arr - id : Identifier').
77 is_a('FI_UNINUM', 'fibo - fnd - arr - id : Identifier').
78 is_a('UNINUM', 'fibo - fnd - arr - id : Identifier').
79 is_a('OFFNAME', 'fibo - be - le - fbo : Branch').
80 is_a('OFFNUM', 'fibo - fnd - arr - id : Identifier').
81 is_a('MAINOFF', 'fibo - be - le - fbo : HeadOffice').
82 is_a('SERVTYPE', 'fibo - fbc - fct - fse : FinancialService').
83 is_a('SERVTYPE_DESC', 'fibo - fbc - fct - fse : FinancialServiceDescription').
84 is_a('ACQDATE', 'time : Instant').
85 is_a('ESTYMD', 'time : Instant').
86 is_a('RUNDATE', 'time : Instant').
87 is_a('CBSA', 'sdmx - concept : statisticalClassification').
88 is_a('CBSA_DIV', 'sdmx - concept : statisticalClassification').
89 is_a('CSA', 'sdmx - concept : statisticalClassification').
90 is_a('STCNTY', 'sdmx - concept : administrativeClassification').
91 is_a('ID', 'dct : identifier').
92 is_a('CBSA_DIV_FLG', 'dcat : Dataset').
93 is_a('CBSA_METRO_FLG', 'dcat : Dataset').
94 is_a('CBSA_MICRO_FLG', 'dcat : Dataset').
95 is_a('CSA_FLG', 'dcat : Dataset').
96 is_a('MDI_STATUS_CODE', 'skos : Concept').
97 is_a('MDI_STATUS_DESC', 'skos : prefLabel').
98 is_a('MDI_STATUS_CODE', 'schema : demographicGroup').
99 is_a('MDI_STATUS_DESC', 'schema : description').
100
```

*Figura 5: The is a relation (partially) in Prolog*

But, in our opinion, a more important contribution of LE to the metadata management is the possibility of quering the document (in English) and obtaining information (also in English) about the related dataset:

```
answer("thing").                                                    ⊕ ━ ⊗

Query thing with noscenario: unknown is an LLM

Answer: -122.117956002298 is a geo : long ( longitude coordinate )

true                                                                      1
                                                            0.498 seconds cpu time

Answer: -122.117956002298 is a LONGITUDE

Answer: -122.117956002298 is a Longitude

Answer: 0 is a CBSA_DIV_FLG

Answer: 0 is a Metropolitan Division Number ( Branch )

Answer: 0 is a dcat : Dataset

Answer: 0 is a CBSA_NO

Answer: 0 is a MAINOFF

Answer: 0 is a CSA_NO

Answer: 0 is a CSA_FLG

Answer: 0 is a CBSA_DIV_NO
```

*Figura 6: Answering the question "which thing is of which type"*

Figure 6 shows one of such type of interaction in which a user ask the document to if it knows about the type of things in it.  But, of course, more specific questions are also possible, like in Figure 7:
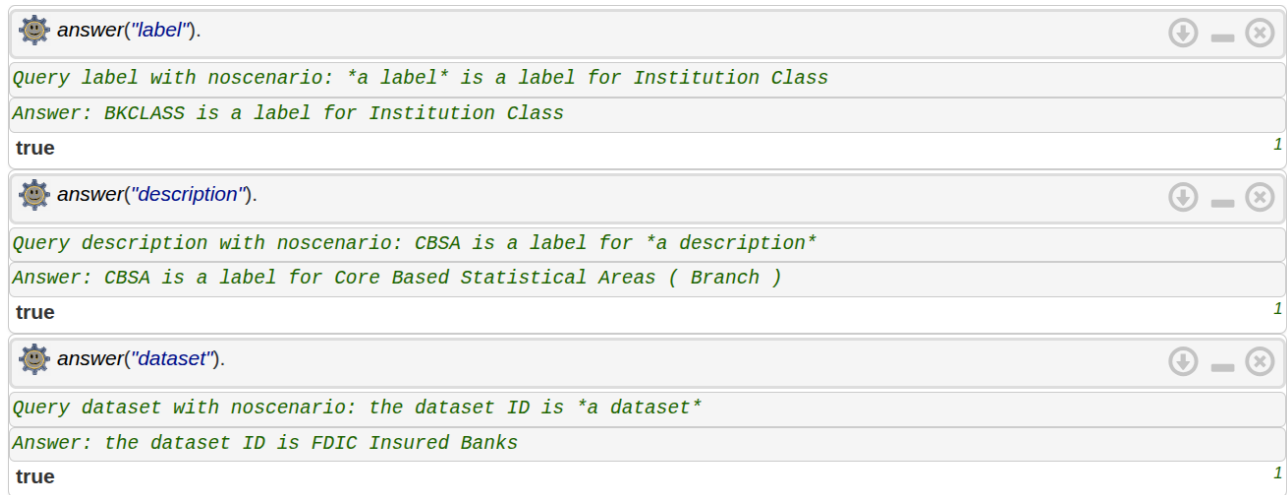


```
answer("label").
Query label with noscenario: *a label* is a label for Institution Class
Answer: BKCLASS is a label for Institution Class
true                                                                    1

answer("description").
Query description with noscenario: CBSA is a label for *a description*
Answer: CBSA is a label for Core Based Statistical Areas ( Branch )
true                                                                    1

answer("dataset").
Query dataset with noscenario: the dataset ID is *a dataset*
Answer: the dataset ID is FDIC Insured Banks
true                                                                    1
```

*Figura 7: Extracting information from the metadata (in English)*

# Describe any features or limitations with respect to the mapping file.

LE is work in progress and there a few details of operation that must be addressed before it can be offer as companion for any dataset. The software is, however, already open source and can be obtained from: https://github.com/LogicalContracts/LogicalEnglish

We have not tested the translation into RDF yet.

# The processing environment(s) in which it was run, including versions of software

Logical English has been developed on SWI-Prolog (threaded, 64 bits, version 9.3.3-200-g7fae34c05). It requires tabling to operate, particularly with the is_a/2 relation. The tests for this submission where done on the swish environment of SWI-Prolog adapted for LE (https://le.logicalcontracts.com/). A serverless version is being tested as well.

# Comments on how the format scales with respect to larger datasets.

We could load the whole dataset requested (one CSV file) into Prolog in about 8 seconds:

*?- time(load_file('FDIC_Insured_Banks.csv')).*
*% 65,506,655 inferences, 8.164 CPU in 8.216 seconds (99% CPU, 8023544 Lips)*
*true.*

And then answer (Prolog) queries from the actual data without any overhead (On a regular Dell Intel® Core™ i7-8650U × 8 PC: with 32,0 GiB RAM, running Ubuntu 24.04.2 LTS)

*?- is_a('St. Louis-St. Charles-Farmington, MO-IL', Type).*
*Type = 'CSA' .*

*?- is_a('18001 Saint Rose Rd', Type).*
*Type = 'ADDRESS' .*

# References

Kowalski, R., Dávila, J., Sartor, G., Calejo, M. (2023). Logical English for Law and Education. In: Warren, D.S., Dahl, V., Eiter, T., Hermenegildo, M.V., Kowalski, R., Rossi, F. (eds) Prolog: The Next 50 Years. Lecture Notes in Computer Science(), vol 13900. Springer, Cham. https://doi.org/10.1007/978-3-031-35254-6_24