

# Introducing HypoGator: generating coherent hypotheses from biased geo-political data

?

**Abstract**—[ABSTRACT]

**Index Terms**—[Items]

## I. INTRODUCTION

Modern KNOWLEDGE BASE MANAGEMENT SYSTEMS allow the ingestion of various data representations, such as relational, graph and full-text data [1]. Current KBMS usually ingest reliable data sources, such as encyclopaedic data [1], business data [2] and medical journals and clinical data [3]. Despite such KBMS are able to reconcile different representation towards an uniform one [4], no technique is currently exploited for detecting *contradicting* facts for hypothesis generation: in particular both data collected from on-line social network [5] or even medical diagnoses [6] may contain contradictions.

*Example 1: One of the simplest ways to find contradictions is to check whether there are facts that are both affirmed and denied at the same time. E.g., fact “Casu Marzu is not a good cheese” is a rebuttal of the fact “Casu Marzu is a exquisite cheese”. Alternatively, we can focus on factual representations that do not allow the contemporaneity of two alternative hypotheses, thus violating a functional dependency. E.g., fact “Yesterday Alice flew to Berlin” contradicts “Yesterday Alice took a trip to Kearney” because Alice cannot be found in two different places at the same time.*

As previously stated, the inherent inconsistency of spurious data sources leads to the generation of conflicting hypotheses in response to a query. The inability of detecting such inconsistencies prohibits to weight how many KBMS facts either support or discard a given hypothesis, thus preventing from correctly ranking the generated hypotheses. This whole scenario demands for a query language which interpretation is able to detect such inconsistencies from the data, and an hypothesis generation task that is able to generate hypotheses which do not contain contradictions.

In order to solve this practical problem, we focus on querying domain-specific knowledge bases containing the aforementioned form of inconsistencies. Such setting allows to use ontologies to represent knowledge, thus allowing to represent both structured (tables), semistructured (wikis, web pages) and unstructured (full texts, videos, images) documents with a uniform representation [7], [8]. MeSH<sup>1</sup> is an example of an ontology for medical settings. Ontologies prompt the intended semantics and properties of such final representations: they differentiate *entities* from *fillers* (e.g., properties), where the

formers are individual units that can be described by different qualities represented as the latter. They also distinguish the (*binary*) *relationships* involving such entities, and the *facts* providing a description of the observed world, which may involve multiple entities and fillers. Events are a specific case of facts. As previously mentioned, this ideal data description does not prevent the data to be biased: as an example, a different point of view may provide biased factual descriptions, that need to be detected. In order to do so, taxonomies and semantic networks focus more on describing entities and fillers through *relationships*, thus allowing to detect data inconsistencies [9]. ICD-11<sup>2</sup> is a well-known medical taxonomy used to identify and categorize specific diseases and medical procedures uniquely. Ontologies and taxonomies provide the pillars sustaining the HypoGator framework, which consists of a pipeline composed of the following blocks, which are described in more detail in the homonymous paper sections (§):

- **Query Interpretation** (§III): any generic query (SPARQL, CYPHER, SQL) is represented as a graph  $q$  [10] and then used to perform an approximated graph matching with the knowledge base [11].
- **Hypothesis Generation** (§IV): the previous approximated graph matching generates the preliminar hypotheses  $h$ : similar or coherent hypotheses are merged (“synthesis”), while conflicting hypotheses are kept separated.
- **Hypothesis Evidence and Scoring** (§V): before returning such hypotheses, we must know which knowledge base element support the provided hypotheses and which does not. Such evidence extraction is used to provide a first scoring function counting how many supporting facts validate the hypothesis against the remaining ones (either supporting or discarding the hypothesis).
- **Hypothesis Query Answering** (§VI): after the previous filter phase, each hypothesis  $h$  is possibly rewritten into an hypothesis  $h'$  maximizing the similarity with  $q$ . Finally, we rank the hypotheses by combining the previous section’s ranking and the edit distance between  $h$  and  $h'$ .

Before providing details of such phases, we analyse the state of the art of query interpretation (§II-A), information extraction (§II-B-II-C) and knowledge base construction (§II-D).

<sup>1</sup><https://www.nlm.nih.gov/mesh/>

<sup>2</sup><https://icd.who.int/>

TABLE I  
ASSOCIATION BETWEEN THE VARIABLES' INSTANTIATION (**Candidate**) TO  
THE DATA ASSOCIATING IT (**Knowledge Base Facts**).

Candidate	Knowledge Base Facts
Rome	Abigail came back to Rome.
Rome	On September 1988, Abigail moved from Bologna's to La Sapienza University, Rome.
Latium	Yesterday (2018/06/22) Abigail was seen in Latium.
Turin	Abigail was rushed to the "Umberto I" hospital in Turin.
Minneapolis	Abigail did a trip to Minneapolis on 1987.
Duluth	Then, Abigail travelled from Minneapolis to Duluth.

## II. RELATED WORK

### A. Natural language query answering

- Natural language querying in relational databases: [2], [12]
- Natural language querying in graph databases: [10]

### B. Rule Mining

Miguel, here we should quote the paper on temporal rule mining.

### C. Approximated graph pattern matching

- Graph matching over dimensions: [13]
- Approximated graph matching: [11], [14]

### D. Knowledge Bases

- KB with precise information: [1]
- KB with uncertainty: [4], [8]

## III. QUERY INTERPRETATION

Given that (natural language) queries can be expressed either as a declarative language [2], [12] or using the same format as the data [10], [15] (e.g., graphs), this paper uses the latter approach because its representation is language independent. Moreover, the latter approach is also able to express most of the natural language queries of interest within our scenario [TODO].

With reference to Figure 1, each entity and filler can be expressed as a vertex, while each binary relationship can be expressed as an edge, and facts can be expressed as hyperedges connecting different vertices [16]. Negations can be expressed by juxtaposing the  $\neg$  symbol. Natural language uses adverbs or pronouns to identify which are the elements to be returned by the query, namely *variables*, which are represented as nodes (or edges) which name starts with a question mark. Such variable instantiation can be performed by approximate graph matching the query graph to the data represented in the knowledge base [11], thus creating a set of *morphisms* linking each variable to the corresponding matched KB elements (*candidate*). The KB's subgraph containing all the candidates for a given morphism is an *hypothesis h*. E.g., Table I provides in its left column a set of possible hypotheses' candidates generated from the knowledge base data and answering the query "Where did Abigail go?".

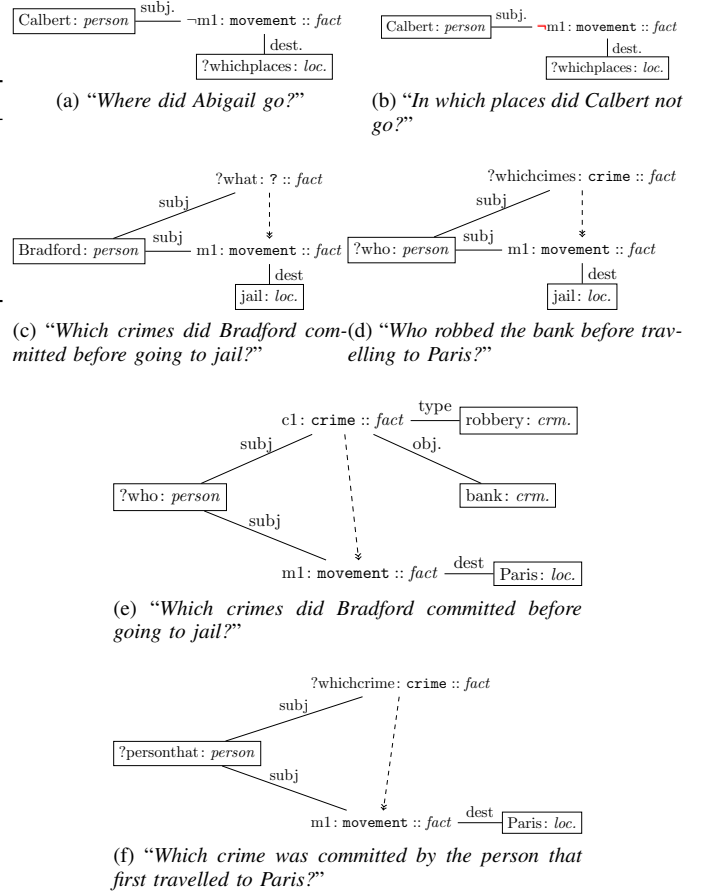


Fig. 1. Each entity/filler is represented by a rectangle, facts have no shape and fact's fields are represented by straight edges. Variables are represented by names preceded by a question mark, ?. Temporal associations between facts are represented by dashed edges.

We now analyse four possible queries of interest within our scenario [TODO]:

- 1) **Return an entity or filler associated to a given fact or relationship.** E.g., for "Where did Abigail go?" (Figure 1a), *go* is the keyword identifying the fact for the specific movement type, while *where* identifies the element to be returned and its type, that is a *location*. Therefore, the query asks for all the facts expressing a movement (*go*) of Abigail towards a specific location (*where*).
- 2) **Return a fact associated to a given fact or relationship.** E.g., in "What did Bradford do before going to jail?" (Figure 1c), the keywords *What* and *do* specify that we need to return a fact which is linked via a temporal association (*before*) to another fact or relationship (*Bradford ... going to jail*). In this case, the to-be-returned fact may have any possible type. The same question can be also refined to return a specific fact type as in the former example: e.g., "Which crimes did Bradford committed before going to jail?" (Figure 1d) may specify to return all the facts pertaining to Bradford's criminal record.
- 3) **Variables appearing in multiple facts/relationships.** E.g., "Who robbed the bank before travelling to Paris?"

(Figure 1e) asks to return an entity (*Who*) which first appears in a movement fact (*Who ...travelling to Paris*) and, later on, performed a crime (*Who robbed the bank*). We want to return an entity appearing in two distinct events having a temporal association. This corresponds to a join between two facts where the criminal corresponds to a moving person.

We can rephrase the same query to return a fact having a variable binding as follows: *Which crime was committed by the person that first travelled to Paris?* (Figure 1f). In this case the variable binding is made explicit by the keywords *person* and *that*, while *Which* remarks that we now want to return only the criminal facts.

- 4) **Negated queries.** E.g., for “*In which places did Calbert not go?*” (Figure 1b), we can provide either all the *places* appearing in negated place relationships or movement facts, or return a set of all the possible *places* not appearing in non-negated place relationships and movement facts. Please note that such queries can be directly expressed in graph query languages which, on the other hand, do not usually query data that may contain negations. On the other hand, current graph pattern matching frameworks may only express positive facts.

#### IV. HYPOTHESIS GENERATION

The hypotheses generator acts as the first part for the query evaluation, which identifies and aggregates the data in the Knowledge Base (approximately) matching the query.

Let us now focus on hypotheses containing only one candidate as in Table I; by using some geographical taxonomy as the one in Figure 2, we may also coarsen the former candidates in order to get a better hypothesis’ support: we have that the hypothesis with candidate (h.w.c.) *Italy* is supported by three entities over five, while h.w.c. *USA* or *Minnesota* is supported by two entities over five. Please note that trivial hypotheses should be discarded: an hypothesis is trivial when it contains all the other generated hypotheses (e.g., “*World*”). Please note that single candidates cannot be inconsistent with themselves, because conflicting candidates are separated in different hypotheses by definition (e.g., all the cities are different to each other, except from Rome that appears twice). Nevertheless, not all the hypotheses are inconsistent to each other (e.g., “*Rome*” is compatible with “*Latium*” and “*Italy*”, but not with “*Turin*” and “*Piedmont*”).

Last, this approach can be generalized over multiple candidates by combining multiple hierarchies together as described in [13], where candidates may be extracted at different possible representation levels [Is this enough or do I have to provide more details on how I would do that?].

#### V. HYPOTHESIS EVIDENCE AND SCORING

In contrast to the evidence in the former section, having coherent candidates does not prevent us from having hypotheses that either contain contradictions or are contradicted by other facts or relationships. It is then required to associate each

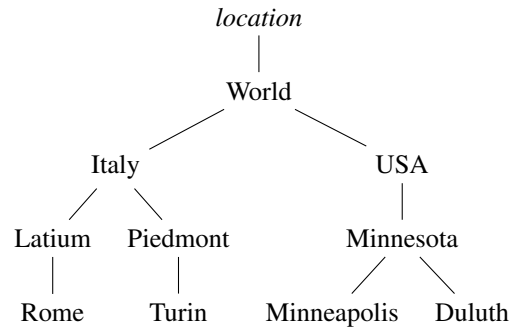


Fig. 2. An example of a geographical taxonomy for geographical dimensions (*location*). Taxonomies may be used to aggregate compatible hypotheses.

generated hypothesis with the set of facts and relationships providing either evidence or counterevidence.

Let us analyse some possible inconsistencies: if we suppose that the relationship “*Abigail has never been in the USA*” is stored within our knowledge base, then we now that this evidence contradicts the facts in Table I generating the candidates “*Minneapolis*” and “*Duluth*”. On the other hand, the relationship “*Abigail was born in Detroit on 1989*” contradicts with the second “*Rome*” candidate, because the entity *Abigail* did not existed prior to 1989. Last, the fact “*On the 22<sup>nd</sup> of June 2018, Abigail died in Turin*” contradicts with the hypothesis “*Latium*” because the same person cannot be located in two different places at the same time.

The former inconsistencies may be observed by comparing different KB facts and relationships with the hypothesis, which combination describes a subset of the whole knowledge graph,  $\kappa$ . We can now represent facts and tuples from  $\kappa$  as tuples, thus allowing to use logical inconsistent detection frameworks which are independent from the graph representation of the data [17]. Rules may be directly provided by the ontologies’ TBox-es [cite:TODO], which describe the correlations between different facts and relationships and provide an homogeneous representation of different facts and relationships. E.g., the TBox may state that each movement fact towards a given destination at a given time implies a *location* relationship of the same subject at the same time. After converting each possible fact and relationships into “finer-grained” facts and relationships, it is now possible to recognize all the spatio-temporal inconsistencies as in [18]. This approach can be therefore extended to any other type of inconsistencies foreseen by the TBox rules in ontologies.

Last, we can now measure its discrepancies via support scores, thus providing a first ranking of the generated hypotheses.

#### VI. HYPOTHESIS QUERY ANSWERING

As previously stated, each “minimal” connected (hyper)graph containing the candidates describes an hypothesis  $h$ . Now, we want to align  $h$  towards the query representation: such alignment can be represented by instantiating the query variables in  $q$  with the candidates, thus obtaining the to-be-

returned hypothesis  $h'$ . Given that  $h$  was obtained through approximate graph matching of  $q$ , we know that the greater the edit distance between the two  $h$  and  $h'$ , the less the reliability on  $h'$  as a possible candidate. By doing so we assign higher scores to the hypotheses that are directly represented within the data, and inferior ones to all the remaining hypotheses that were generated over imprecise(?) inference rules. The combination of the support measure with the edit distance provides the score assigned to  $h'$ .

## VII. CONCLUSION

## ACKNOWLEDGMENT

This work is supported by DARPA under FA8750-12-2-0348-2 (DEFT / CUBISM).

## REFERENCES

- [1] I. J. of Research and Development, *This is Watson*, J. W. Murdock, Ed. IBM Co., May/July 2012, vol. 56(3/4).
- [2] D. Saha, A. Floratou, K. Sankaranarayanan, U. F. Minhas, A. R. Mittal, and F. Özcan, "Athena: An ontology-driven system for natural language querying over relational data stores," *Proc. VLDB Endow.*, vol. 9, no. 12, pp. 1209–1220, Aug. 2016.
- [3] Y. Wang, L. Wang, M. Rastegar-Mojarad, S. Moon, F. Shen, N. Afzal, S. Liu, Y. Zeng, S. Mehrabi, S. Sohn, and H. Liu, "Clinical information extraction applications: A literature review," *Journal of Biomedical Informatics*, vol. 77, pp. 34 – 49, 2018.
- [4] F. Niu, C. Zhang, C. Ré, and J. Shavlik, "Elementary: Large-scale knowledge-base construction via machine learning and statistical inference," *Int. J. Semant. Web Inf. Syst.*, vol. 8, no. 3, pp. 42–73, Jul. 2012.
- [5] D. M. J. Lazer, M. A. Baum, Y. Benkler, A. J. Berinsky, K. M. Greenhill, F. Menczer, M. J. Metzger, B. Nyhan, G. Pennycook, D. Rothschild, M. Schudson, S. A. Sloman, C. R. Sunstein, E. A. Thorson, D. J. Watts, and J. L. Zittrain, "The science of fake news," *Science*, vol. 359, no. 6380, pp. 1094–1096, 2018.
- [6] D. Costa and M. A. Martins, *Measuring Inconsistency in Information*. College Publications, 2018, ch. Inconsistency Measures in Hybrid Logics, pp. 169–194.
- [7] C. Zhang, "Deepdive: A data management system for automatic knowledge base construction." Ph.D. dissertation, University of Wisconsin-Madison, 2015.
- [8] Y. Chen and D. Z. Wang, "Knowledge expansion over probabilistic knowledge bases," in *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD '14. New York, NY, USA: ACM, 2014, pp. 649–660.
- [9] P. Lependu and D. Dou, "Using ontology databases for scalable query answering, inconsistency detection, and data integration," *J. Intell. Inf. Syst.*, vol. 37, no. 2, pp. 217–244, Oct. 2011.
- [10] S. Hu, L. Zou, J. X. Yu, H. Wang, and D. Zhao, "Answering natural language questions by subgraph matching over knowledge graphs," *IEEE Trans. Knowl. Data Eng.*, vol. 30, no. 5, pp. 824–837, 2018.
- [11] R. De Virgilio, A. Maccioni, and R. Torlone, "Approximate querying of rdf graphs via path alignment," *Distributed and Parallel Databases*, vol. 33, no. 4, pp. 555–581, Dec 2015.
- [12] F. Li and H. V. Jagadish, "Understanding natural language queries over relational databases," *SIGMOD Rec.*, vol. 45, no. 1, pp. 6–13, Jun. 2016.
- [13] A. Petermann, G. Micale, G. Bergami, A. Pulvirenti, and E. Rahm, "Mining and ranking of generalized multi-dimensional frequent subgraphs," in *Twelfth International Conference on Digital Information Management, ICDIM 2017, Fukuoka, Japan, September 12-14, 2017*, pp. 236–245.
- [14] J. Aligon, E. Gallinucci, M. Golfarelli, P. Marcel, and S. Rizzi, "A collaborative filtering approach for recommending olap sessions," *Decision Support Systems*, vol. 69, pp. 20 – 30, 2015.
- [15] M. P. Consens and A. O. Mendelzon, "Graphlog: A visual formalism for real life recursion," in *Proceedings of the Ninth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, ser. PODS '90. New York, NY, USA: ACM, 1990, pp. 404–416. [Online]. Available: <http://doi.acm.org/10.1145/298514.298591>
- [16] R. Fagin, "Degrees of acyclicity for hypergraphs and relational database schemes," *J. ACM*, vol. 30, no. 3, pp. 514–550, Jul. 1983.
- [17] M. Minoux, "Ltur: a simplified linear-time unit resolution algorithm for horn formulae and computer implementation," *Information Processing Letters*, vol. 29, no. 1, pp. 1 – 12, 1988. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/002001908890124X>
- [18] J. Grant, M. V. Martinez, and F. Molinaro, Cristian Parisi, "On measuring inconsistency in spatio-temporal databases." in *Measuring Inconsistency in Information*, J. Grant and M. V. Martinez, Eds. College Publications, 2018, ch. 11, pp. 313–342.