

Entity Search with NECESSITY

Ekaterini Ioannou[‡]

Anshul Jain[§]

Claudia Niederée[‡]

Saket Sathé[§]

Srikanth Bondalapati[§]

Zoltan Miklos[§]

Nicolas Bonvin[§]

Gleb Skobeltsyn[§]

[§]Ecole Polytechnique Fédérale de Lausanne (EPFL)
{name.surname}@epfl.ch

[‡]L3S Research Center/Leibniz Universität Hannover
{surname}@L3S.de

ABSTRACT

Loosely structured heterogeneous information spaces are typically created by merging data from a variety of different applications and information sources. A common problem these information spaces need to address is that various data describe the same real-world entities (e.g., conferences, organizations). In this demo, we introduce NECESSITY, an efficient and scalable entity store. NECESSITY is able to handle a large number of entities and at the same time provide an efficient and highly accurate entity search functionality for heterogeneous and partially structured queries that follow the vision of dataspace.

1. MOTIVATION AND OUTLINE

We are currently witnessing a rapid increase in the number of loosely structured *heterogeneous information spaces* - collections of data coming from a variety of different applications and information sources. One common problem these information spaces face, is managing their entities (e.g., organizations, events), since there will be given various representations for the same real world *entities*. The NECESSITY entity store¹ is able to address this challenge. Our system can handle a large number of entities and at the same time provides an efficient and highly accurate entity search functionality.

NECESSITY stores entity profiles composed by a set of attribute-value pairs. It allows efficient entity search over these loosely structured heterogeneous information spaces, with queries being conditions on the entity's attributes or values. As such, our approach contributes to the idea of realizing dataspace as envisioned in [3] in both the data model for storing data but also the challenge for developing searching methods for querying a large number of diverse and interrelated data.

NECESSITY was implemented and evaluated in the con-

¹The name of the system was inspired by philosopher William of Ockham's famous principle: "Entities must not be multiplied beyond *necessity*".

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

text of the Entity Name System (ENS), for the OKKAM project². The aim of ENS is to foster the global re-use of entity identifiers and to mediate between existing identifiers for individual entities (details available in [1]). ENS receives queries and checks whether the entity described in each query exists in OKKAM. If the entity exists OKKAM returns the corresponding identifier. The core benefit of ENS is the easing integration of external applications. For example, repositories from personal information management systems can now rely on the service provided by the ENS for creating their URIs for entities. As such, the integrating challenge of knowing which representations in different repositories refer to the same entity, would be resolved by the use of shared IDs as issued by OKKAM.

There are further application types that can profit from our proposed approach. One example is entity search in collaboratively authored information spaces, such as Wikipedia. Each Wikipedia entry is composed by various contributors who are not enforced to follow a specific format or schemata in a consistent way. As such, processing a human query over Wikipedia data could benefit from matching the query with the heterogeneous data of the entities. An other example of targeted applications for NECESSITY is entity search engines. These applications are built upon information extracted from Web pages. Effectively integrating this extracted data imposes a matching challenge of effectively identifying and merging the existing data that refer to the same real world entities. In addition, searching for a specific entity through this plethora of entities requires advanced matching functionalities.

The rest of this paper is organized as follows. Section 2 presents the NECESSITY entity store, with main focus on the entity search process. Section 3 describes the demo scenario, and Section 4 conclusions along with future work.

2. ENTITY SEARCH PROCESS

Entities in NECESSITY are modeled as a set of attribute-value pairs; a representation similar to dataspace proposal [3, 2]. As such, a person entity will be represented by name, affiliation, email address, and whatever else is available. Entity search allows users or applications to retrieve the entities —ideally only one— already in NECESSITY that best match an entity description provided as a query to the system. An entity query is a set of predicates, where each predicate is a keyword or an attribute value pair, e.g., Q_1 : name="John Smith" EPFL, and Q_2 : name=Smith affiliation=EPFL.

²<http://www.okkam.org>

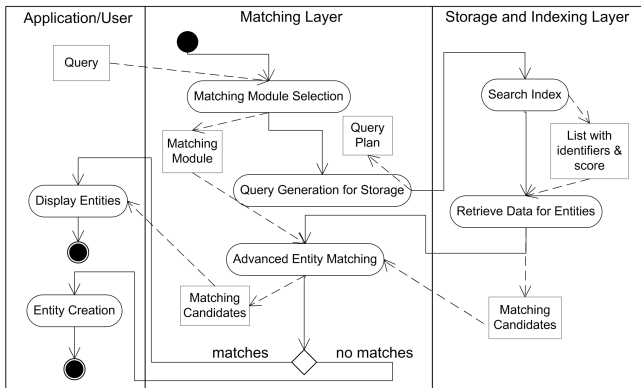


Figure 1: The search process of NECESSITY.

Figure 1 is an illustration of the search process incorporated in NECESSITY for answering entity queries. When NECESSITY receives an entity query, it selects the appropriate matching module to process it. We consider the possibility of different methods for selecting the module, for example explicit selection by user/application (e.g., based on previous experience), or selection using query information such as restrictions in execution time. Once the matching module is selected, the query is reformulated to realize the semantics of the entity storage and to cope with incomplete or imprecise information. More specifically, this action can generate a disjunction of predicates over the entity profiles in the storage layer, or include missing schema information. The generated query is then sent to the storage and indexing layer for evaluation. The storage will use the distributed index to retrieve a *constant* number of the most relevant entities, named “matching candidates”, and then their corresponding entity profiles. The matching module will receive these candidates and perform advanced entity matching. This will select the most relevant entities for the given query by computing the matching probability between each candidate and the entity described in the given query.

There are two possible outcomes of the entity search process. The first is an empty list, which indicates that the entity described by the query is not in the NECESSITY, and thus the application/user should create it. The second possible result is a ranked list of entities which were found to match the query.

We evaluated NECESSITY with 1,000,000 entities, and 500 queries extracted from various real data and web pages. Both entities and queries, were generated using Cogito extractor by Expert System³, and each query was manually processed to identify the corresponding entity. NECESSITY returned the correct entity in the first five results for 89% of the queries, and in the first ten results for 93% of the queries.

3. DEMONSTRATION

In the demo, users will be able to perform their own entity search queries on the NECESSITY entity store. NECESSITY will contain at least 1M entities, including people and organizations from Wikipedia⁴, locations from GeoNames⁵, and

³<http://www.expertsystem.net/page.asp?id=1515/>

⁴<http://www.wikipedia.com/>

⁵<http://www.geonames.org/>

proteins from UniProt⁶. In addition, we will present the details and discuss the various aspects of NECESSITY, particularly in relation to the topics shortly described in the following paragraphs.

Entity Queries. NECESSITY can handle various types of queries made from different sources. The various types include; queries which are over-specified or under-specified, contain only keywords, or incomplete information.

Queries Generated for the Storage. User queries are internally reformulated in several ways. This includes query extension, and incorporation of higher weights to attributes which are considered more important than others (e.g., name for person entities).

Identification of the Matching Candidates. NECESSITY uses a key-value store for storing the entity profiles, along with an inverted index for processing queries. These are used to efficiently drill down top-k matching candidates from the storage that best match the given query.

Advance Matching. Our current implementation of NECESSITY contains two matching modules. The first is group linkage, which returns entities those attribute value pairs have high similarity with query predicates [4]. The second module is generic matching, an extension of group linkage that in addition considers domain specific similarity functions and the selectivity of predicates included in the query.

4. CONCLUSIONS & FUTURE WORK

This paper introduced NECESSITY, an entity storage for managing entities. NECESSITY is able to handle a large number of entities while providing an efficient and highly accurate search functionality, with a longer-term goal to be able to handle entities of web-scale. Our future plans include performing a new set of experiments that aim at testing the system in a distributed setup. In addition, we plan better translation of the requests into queries for the storage with special focus on reducing the time needed for answering these requests and increasing the result quality.

5. ACKNOWLEDGMENTS

The authors thank Oleksandr Druzhynin for helping in implementation; Dr. Peter Fankhauser, and Juri Luca De Coi for allowing the use of their generic matcher; and the participants of the OKKAM project for valuable discussions. This work is partially supported by FP7 EU Project OKKAM (contract no. ICT-215032).

6. REFERENCES

- [1] P. Bouquet, H. Stoermer, C. Niederee, and A. Mana. Entity name system: The backbone of an open and scalable web of data. In *ICSC*, 2008.
- [2] G. H. L. Fletcher, J. V. den Bussche, D. V. Gucht, and S. Vansummeren. Towards a theory of search queries. In *ICDT*, 2009.
- [3] A. Y. Halevy, M. J. Franklin, and D. Maier. Principles of dataspace systems. In *PODS*, 2006.
- [4] B.-W. On, N. Koudas, D. Lee, and D. Srivastava. Group linkage. In *ICDE*, 2007.

⁶<http://www.uniprot.org/>