

A Natural Language Interface over the MUSICBRAINZ Database

Johan Granberg and Michael Minock

Department of Computing Science: Umeå University

Abstract. This paper demonstrates a way to build a natural language interface (NLI) over semantically rich data. Specifically we show this over the MUSICBRAINZ domain, inspired by the second shared task of the QALD-1 workshop. Our approach uses the tool C-PHRASE [4] to build an NLI over a set of views defined over the original MUSICBRAINZ relational database. C-PHRASE uses a limited variant of X-Bar theory [3] for syntax and tuple calculus for semantics. The C-PHRASE authoring tool works over any domain and only the end configuration has to be redone for each new database covered – a task that does not require deep knowledge about linguistics and system internals. Working over the MUSICBRAINZ domain was a challenge due to the size of the database – quite a lot of effort went into optimizing computation times and memory usage to manageable levels. This paper reports on this work and anticipates a live demonstration¹ for querying by the public.

Keywords: Natural Language Interfaces, Relational Databases, MUSICBRAINZ, C-PHRASE

1 Introduction

It has often been noted that natural language interfaces to databases suffer when the manner in which data is stored does not correspond to the user’s conceptual view of such data [1, 2, 5]. This mismatch between the way data is structured and the user’s conceptual model can be for a variety of reasons, but here we speculate that the three most common reasons are:

1. The database is highly normalized (e.g. to BCNF) for the sake of eliminating update anomalies.
2. The database is highly abstracted, including many attributes per relation so as to avoid cost associated with joins.
3. The database is stored in a semi-structured form corresponding to RDF triples.

It is the thesis of this paper that the user would prefer to query the data in a conceptual form similar to what is represented in the entity-relationship diagram

¹ <http://www.cs.umu.se/~johang/research/brainz/public/>

of the database domain (see figure 1). Naturally if we are given a database in one of the three forms above, we assume that it is possible, via standard view definitions, to transform the data so that it may be accessed (and perhaps even updated) via the conceptual model.

This paper explores these ideas and shows some preliminary results over the MUSICBRAINZ database paired with the 50 natural language queries in the shared task for MUSICBRAINZ in the QALD-1 workshop. In section 2 we present our approach to building a natural language interface supporting queries over MUSICBRAINZ. Section 3 discusses our initial results and some anecdotes from our development efforts. Section 4 summarizes our findings and points toward near term and longer term plans, including the fielding of a live interface to MUSICBRAINZ for querying by the public.

2 Approach

2.1 The conceptual model

After looking at the set of 50 natural language queries to be supported over MUSICBRAINZ, we defined the conceptual model appearing in figure 1. This diagram is traditional except that it shows a form of conceptual aggregation. For example a group or a person can release an album, soundtrack or single. A traditional ER diagram would require six relationship diamonds instead of one to depict this.

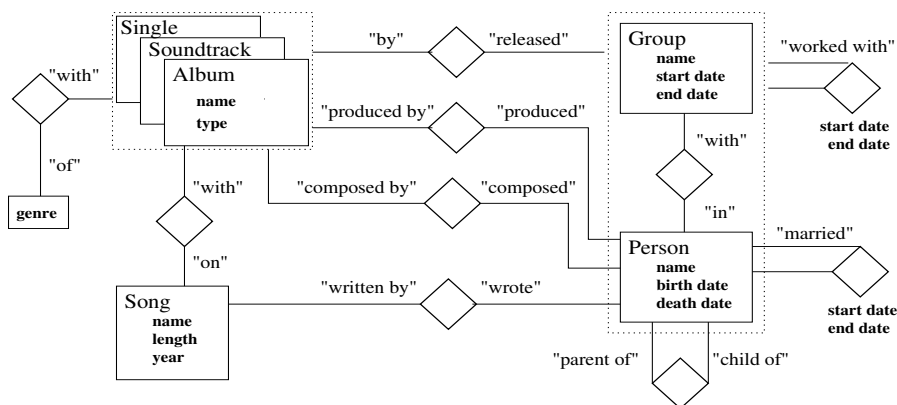


Fig. 1. Conceptual model of MUSICBRAINZ.

2.2 The data source

We had two choices for what we used as the data source for the actual MUSICBRAINZ data. The first was the original data stored in a PostgreSQL

database and the second was the RDF dump of the data built for the QALD-1 shared task. Because of our rooting in relational databases, we chose to simply draw on the data in the original MUSICBRAINZ POSTGRES SQL database.

The schema of the underlying MUSICBRAINZ database is primarily designed according to option 2 from section 1 above. That is, several highly abstracted relations with many attributes represent abstract entities such as `L_ARTIST_ARTIST` and `LT_ARTIST_ARTIST` which has tuples that represent both individual artists (e.g. Bob Dylan) as well as bands (e.g. REM). The database itself is rather large. There are approximately 10.6 million songs, 0.8 million albums and 0.6 million individual and bands. In our experiment we remove the tuples containing non-ASCII characters and arrive at 9.7 million songs, 0.7 million albums and 0.4 million artists.

2.3 The view definitions

Views were defined in the standard way over the base relations of the MUSICBRAINZ database. One consideration was whether to materialize these views for quicker access to the data. This essentially doubles the database size. Our findings are that this leads to a speed up factor of approximately 3. For example using regular views the test query, “which singles did the Dead Kennedys release?” took 2.0 ms on average. Using materialized views it took an average of 0.7 ms. Considering other performance issues in the system (e.g. natural language parsing times), we concluded that querying through the views is not currently a bottleneck. Thus we did not elect to materialize views.

2.4 Authoring with the C-PHRASE administration interface

Once several errors were dealt with (see section 3.1) the authoring process proceeded well. It followed the name, tailor and define method laid out in [4]. Well over half the training set can be authored for within 90 minutes. With our new enhancements to the parser to better handle WH-movement, this may be further reduced. We intend to produce series of YOUTUBE videos that demonstrate this process.

3 Preliminary Results

3.1 Difficulties

There were several difficulties we encountered that, while perhaps anecdotal, are still worth mentioning. To date the C-PHRASE system has been applied only over small databases. MUSICBRAINZ is a sizable database, so this brought up some scalability issues that we had not earlier experienced. We assume other NLI systems attempting to scale to this size of a database might face similar problems.

Large main-memory hash tables The first issue related to memory management issues in CLISP, the version of LISP that C-PHRASE is implemented over. There are some unfortunate memory bugs that corrupt CLISP memory when utilization climbs over a certain threshold. It is difficult to track, because the corruption generally causes a **Segmentation Fault** at later steps when memory is accessed. Under normal circumstances this problem does not surface. However to allow for named entity recognition, C-PHRASE materializes string values from the database into hash tables to scan for matching values in the user's typed request. In the case of MUSICBRAINZ this means building main memory hash tables containing millions of constants. This was too much for CLISP to handle.

After several false starts, the solution to this problem was to implement a remote hash facility that maintained the main memory hash tables remotely outside of CLISP. The overhead access time for these hash tables is negligible and the implementation is stable and scalable, bounded ultimately by the size of virtual memory.

Limitations of the PostgreSQL query optimizer C-PHRASE maps English to logical expressions in Codd's tuple calculus. From such logical expressions, SQL is in turn generated. Before our experiment we would generate SQL such as the following to answer the query, "Which singles did the Dead Kennedys release?"

```
SELECT DISTINCT NAME
FROM SINGLE AS x
WHERE
  EXISTS(
    SELECT *
    FROM BAND as y1
    WHERE
      x.artist = y1.id AND
      y1.name = 'Dead Kennedys').
```

Unfortunately such queries are not taken up by PostgreSQL's optimizer. This query in fact takes 23 minutes to answer on an older Solaris server where we run our database. In contrast the equivalent query

```
SELECT DISTINCT x.NAME
FROM SINGLE AS x,BAND as y1
WHERE x.artist = y1.id AND y1.name = 'Dead Kennedys'"
```

takes 0.7 seconds on the same server. Here the query time is entirely dominated by the time to establish the connection, the actual query execution is reported as 2 ms. We assume that the speed-up is due to the fact that the optimizer has access to both relations on the second line of the query and can plan accordingly. Instead of pestering PostgreSQL about 'improving their optimizer', we altered our translator to produce SQL of the later variety.

3.2 Performance

We are still in the process of collecting performance data. The running example query of "Which singles did the Dead Kennedys release?" shows that the performance is adequate in the case of single join queries. We have run additional tests on queries that exercise more joins and are convinced that the current approach is feasible. The time it takes to parse natural language queries is typically in the range of one to three seconds.

3.3 Current coverage

Unfortunately recent extensions to C-PHRASE have introduced system instability and have blocked a systemic precision, recall and f-measure study on the 50 unseen queries released as part of the QALD-1 shared task. The problems are mostly due to unanticipated parser bugs and performance problems when extending our parser to handle gap threading. Work continues to resolve these problems. In the meantime we have elected not to read the 50 new unseen queries in anticipation of doing a clean coverage test in the future.

As for the original 50 queries for the MUSICBRAINZ example, in our prior working version of C-PHRASE, we covered all but 5 of these queries. The 5 problematic queries (in order of estimated level of difficulty) are those involving implied time intervals ("How many bands broke up in 2010"), types (e.g. "Is Liz Story a person or a group?"), complex time calculations ("Which artists have their 50th birthday on May 30, 2011?"), computed comparison values with ellipsis (e.g. "Which artists died on the same day as Michael Jackson"), and queries involving non-quoted, non-domain string values (e.g. "Are the members of the Ramones that are not called *Ramone*"). Time permitting, we will extend C-PHRASE to handle these types of queries. It will be interesting to learn about how other groups at QALD-1 approached these queries.

4 Conclusions

We were very pleased when we heard of the QALD-1 shared task. We decided to focus our efforts on the more closed-domain task of queries over MUSICBRAINZ. We based our data on the original MUSICBRAINZ relational database and, after confronting several technical difficulties in scaling C-PHRASE, we managed to build a natural language interface that covered 45 of the 50 training examples in the QALD-1 shared task for MUSICBRAINZ.

Unfortunately technical problems have delayed a systematic evaluation. We still intend to perform and report an evaluation of how well we cover the 50 unseen queries, but perhaps more tellingly we intend to field a natural language interface² for real-time querying of MUSICBRAINZ by the public.

² <http://www.cs.umu.se/~johang/research/brainz/public/>

References

1. I. Androutsopoulos and G.D. Ritchie. Database interfaces. In R. Dale, H. Moisl, and H. Somers, editors, *Handbook of Natural Language Processing*, pages 209–240. Marcel Dekker Inc., 2000.
2. A. Copestake and K. Sparck Jones. Natural language interfaces to databases. *The Natural Language Review*, 5(4):225–249, 1990.
3. R. Jackendoff. *X-bar-Syntax: A Study of Phrase Structure, Linguistic Inquiry Monograph 2*. MIT Press, 1977.
4. M. Minock. C-phrase: A system for building robust natural language interfaces to databases. *Journal of Data and Knowledge Engineering (DKE)*, 3(69):290–302, 2010.
5. W. Ogden and P. Bernick. Using natural language interfaces. Technical Report MCCS-96-299, Computing Research Laboratory, New Mexican State University, Las Cruces, May 1996.