

LexEQUAL: Supporting Multilexical Queries in SQL

A. KUMARAN *

JAYANT R. HARITSA *

1. Introduction

Current database systems offer support for storing multilingual data [2], but are not capable of querying *across languages*, an important consideration in today's global economy. We therefore propose a new multilexical operator called LexEQUAL, that extends the standard lexicographic matching in database systems to matching of text data across languages, specifically for names, which form close to twenty percent of text corpora.

Our implementation of the LexEQUAL operator is based on transforming matches in language space into parametrized approximate matches in the equivalent *phoneme space*. A detailed evaluation of our approach on a real data set shows that there exist settings of the algorithm parameters with which it is possible to achieve both good *recall* and *precision*.

2. The LexEQUAL Operator

Author	Author_FN	Title	Price	Language
Descartes	René	Les Méditations Métaphysiques	€ 49.00	French
தேகர்டீசு	ரெனே டிகார்டீ	ஆன்மீக மேதாபிசு	INR 250	Tamil
د	بهنسي ، د	العمارة عبر التاريخ	SAR 75	Arabic
Nehru	Jawaharlal	Discovery of India	\$ 9.95	English
Durant	Will	History of Civilization	\$ 149.75	English
नेहरू	जवाहरलाल	भारत एक खोज	INR 175	Hindi

Catalog of Multilingual Books.com

Consider a hypothetical Books.com that sells books in different languages, with a sample product catalog as shown above, and a user query to retrieve all the works of an author, in a set of specified languages. This can be easily achieved using LexEQUAL as shown below:

```
SELECT * FROM Books
WHERE Author LexEQUAL 'Nehru' Threshold 0.3
IN { English, Hindi, Tamil, Arabic }
```

The output for this query will contain all records of the above Books.com table, for whom the Author attribute contains names that are *phonemically close* to "Nehru", with the Threshold parameter in the query determining the match quality tradeoff between precision and recall.

* Department of Computer Science and Automation, Indian Institute of Science, Bangalore, INDIA. {kumaran, haritsa}@csa.iisc.ernet.in

```
LexEQUAL (Sl, Sr, e)
Input: Strings Sl, Sr, Error Threshold, e
       Languages with TTP transformations, SL
1. Ll ← Language of Sl; Lr ← Language of Sr;
2. if Ll ∈ SL and Lr ∈ SL then
3.   Tl ← transform(Sl, Ll); Tr ← transform(Sr, Lr);
4.   Smaller ← (|Tl| ≤ |Tr| ? |Tl| : |Tr|);
5.   if editdistance(Tl, Tr) ≤ (e * Smaller)
6.     then return TRUE else return FALSE;
7. else return NORESOURCE;
```

The pseudocode of our implementation of LexEQUAL is given above. The algorithm transforms the input multilingual strings to their equivalent phoneme strings and flags a match if the edit distance between them is less than the user-specified error limit, Threshold. The transform function uses standard *Text-to-Phoneme (TTP)* converters that convert to a canonical *International Phonetic Association's* phonemic alphabet. Further, the LexEQUAL implementation is parameterized for different cost functions of editdistance function and for domain-specific clustering of similar phonemes.

We have currently implemented LexEQUAL as a *user-defined function*, incorporating *Q-Grams* [1] to filter out *non-matches* inexpensively and *Phonemic Indexes* [4] to narrow the search to *potential matches* using standard database index structures. Our initial results on a commercial database system hold out the promise that an *inside-the-server* implementation of multilingual matching will have runtimes comparable to traditional monolingual matching.

In summary, the LexEQUAL operator employing phonetic matching can complement the standard lexicographic operators, representing a first step towards achieving complete multilingual functionality in database systems. The full version of this paper is available in [3].

Acknowledgements This work was supported in part by a Swarnajayanti Fellowship from the Dept. of Science and Technology, Govt. of India.

References

- [1] L. Gravano et al. Approximate String Joins in a Database (almost) for Free. *Proc. of 27th VLDB Conf.*, 2001.
- [2] A. Kumaran and J. Haritsa. On the Costs of Multilingualism in Database Systems. *Proc. of 29th VLDB Conf.*, 2003.
- [3] A. Kumaran and J. Haritsa. Supporting Multiscript Matching in Database Systems. *Proc. of 9th EDBT Conf.*, 2004.
- [4] J. Zobel et al. Phonetic String Matching: Lessons from Information Retrieval. *Proc. of 19th ACM SIGIR Conf.*, 1996.