

Language Analysis in Library OPAC: Designing an Open Source Software based Framework for Bibliographic Records in Mainstream and Tribal Languages

Parthasarathi Mukhopadhyay* and Anirban Dutta

Department of Library and Information Science, University of Kalyani, Nadia - 741 235, India

**E-mail: psmukhopadhyay@gmail.com*

ABSTRACT

It reports the development of an enhanced library OPAC prototype through integration of language analysis tool and book reader in the retrieval interface. Language analysis or text analytics is considered as one of the components of language documentation and when integrated with library OPAC can extend supports to analyse corpus of the retrieved document in terms of word/phrase frequency, term circles, term links, term context etc through visual representation in a single-window along with the other datasets generally expected in a typical library OPAC. The open source software based integration mechanism is tested with English and Bengali as mainstream languages and a Unicode-compliant Indian official tribal language Santali (Ol Chiki script) as minority language.

Keywords: Language analysis; Text analytics; Enhanced OPAC, Book reader; Ol Chiki; Santali.

1. INTRODUCTION

As opined by Boerger, B.H., *et al.*¹, the general workflow of language documentation includes recording, metadata maintenance, transcribing, annotation & analysis, translation, archiving and dissemination. In simple words, language documentation refers to documenting and describing verbal, written, and gesture (for sign language) forms of a language systematically in their suitable sociocultural context with emphasis on long-term archiving and access^{2,3,4}. Library professionals can prove themselves helpful in many of these aspects like annotation, metadata encoding & standardisation, long-term archiving and effective dissemination. One of the major purposes of language documentation or documentary linguistics, apart from the activities as given, is to provide scope for linguistic analysis in any given language from the corpora of textual materials available in that language^{5,6}. The advancement of the NLP tools and techniques, availability of the open source language analysis tools/services and increasing use of the open source software in library management are the factors that provide an opportunity to enhance library OPAC by rendering not only bibliographic metadata in retrieval interface but also to extend supports to obtain results of language analysis such as word/phrase frequency, term circles, term links, term context and many more from the corpus of retrieved document through a managed visual interface.

2. REVIEW OF LITERATURE AND OVERVIEW OF SYSTEMS

Language documentation is a multidisciplinary research

domain and consists of facets like specific stories, samples of speeches, writing codes etc⁷. Furbee⁸ said that the theories of this emerging research area are still premature but Hill⁹ suggested that language ideologies are integral part of the domain. Varela & Lee¹⁰ reported that inclusion of language analysis ensures better utilities in language documentation and Austin¹¹ pointed out four primary constraints related to this: i) ambiguous nature of content (makes it difficult to determine quality of the resources); ii) multidisciplinary approaches (includes applied linguistics, anthropology, information science, etc.); iii) non-availability of comprehensive metadata schema; and iv) sustainability (trainings and capacity building programmes are required for achieving combined skills). Many models of language analysis have been developed by researchers to address the issues raised by domain experts like cognitive folk model for schema analysis^{12,13}, schema analysis model¹⁴ and content analysis model¹⁵. A few researchers came up with frameworks and algorithms to solve problems of syntactic and semantic structures of languages¹⁶, morphological alterations and morphotactics¹⁷, and language analysis techniques compatible with Natural Language Processing ('meaning extraction' and 'named entity extraction')¹⁸. There are some ground-breaking software models in the domain of language documentation specifically meant for language analysis and results visualisation like 'Computer Aided Text Analysis (CATA)'¹⁹, 'Corpus Linguistics (CL)'²⁰ and text analyses of indigenous literature²¹. The primary targets of these software solutions were schema analysis and identifications of narratives, metaphors, proverbs and other linguistic and paralinguistic features of languages. Many text analysis tools and guidelines have also been developed by Illinois Library²². Only a few studies are so far carried out on Indian tribal

languages documentation till date. The Santali language is one of the constitutionally approved Indian tribal languages. It is a sub-group of Austroasiatic-Mundari language and the only tribal language which has extensive web presence (Santali version of Wikipedia was launched in 2018 in Unicode-compliant Ol Chiki script)²³. The ‘Ol Chiki’ script, invented by Pandit Raghunath Murmu in the 1920s is now used as an official script of Santali language. It consists of six vowels, twenty four consonants, and five sound symbols²⁴.

3. OBJECTIVE

Language documentation means creating or collecting, processing, preserving and disseminating linguistic data (such as, texts, audios, videos)²⁵. Hundesberger²⁶ said way back in

2009 that “Libraries and language labs should be two closely cooperating entities, possibly even under the same roof, to make it easier for students to follow up links”. This research work is aimed to support views of the domain experts that a library should support both retrieval and analysis from the same search interface and the main objective of this research is to design and development of an enhanced OPAC in an open source ILS (Integrated Library System) where end-users can access retrieved bibliographic resource in full-text systematically through a book reader and can visualise different results of language analysis in the same retrieval interface. It requires seamless integration of ILS, book reader software and language analysis tool to provide such enhanced information services in OPAC. The mechanism as proposed can support

Table 1. Tasks and steps related to framework design

Task	Steps	Result
Identification and selection of front-end ILS	Major considerations 1. Open source, Web-enabled, Unicode-compliant and REST/API support; 2. Scope for API development without hampering future upgradation; and 3. Ability to handle selective occurrences of a repeatable MARC tag/subfield and conditional display control	Possible candidate software: Koha, NewGenLib, Evergreen Selection Koha (keeping in view its flexibility in API handling and user base)
	Major considerations 1. Open source, Web-enabled, Unicode-compliant and data visualization support; 2. Scope for selective export of result display URL for linking in bibliographic framework selectively; and 3. Ability to generate a series of analysis services, separate URL for each analysis and control at the user end to navigate from one analysis to another.	Possible candidate software: Voyant, Lexos, Textalyser, Word & Phrase and WordSeer Selection Voyant (keeping in view its support for wide array of filetypes, different analysis services, improved visualization facilities and local installation options)
Identification and selection of back-end language analysis tool/service	Major attributes 1. The architecture of the mechanism must support URL based linking between Koha catalogue dataset and Voyant corpus; 2. At Voyant end the corpus (Unicode-compliant full-text document associated with a MARC record in Koha) must produce unique URLs for available analysis results; 3. A repeatable MARC tag/subfield needs to be selected which can record corpus data URL (analysis service URL) along with other links; and 4. Linking script at Koha end needs to create a “Language analysis” tab conditionally in OPAC, if MARC tag/subfield provides analysis URL, and then can fetch data visualization from corresponding corpus available at Voyant end in the given OPAC tab with facility to navigate other analysis URLs.	Features: 1. As it is targeting OPAC, the logical selection of API end in Koha are - OPACUserCSS, and OPACUserJS (these two system preferences may be populated with the script that will link Koha and Voyant); 2. MARC tag 856 sub-field u is a logical choice to include locally available URL for language analysis services; as \$u is repeatable, any number of analysis URLs and link to book reader URL can be included here; and 3. The script has two logical section – a) linking algorithm to connect book reader, if first occurrence of 856 \$u is populated in MARC record; and b) linking algorithm to connect language analysis, if second occurrence of 856 \$u is available in MARC record.
	Linking mechanism	

Hindu gods and heroes : studies in the history of the religion of India, by Lionel D. Barnett.
by Barnett, Lionel D. (Lionel David), 1871-1960.

Material type: Text; Format: print ; Literary form: Not fiction

Publisher: London : J. Murray, 1922

Available from Many Books. Click book title to load the book in reader.

Availability: **Items available for loan:** (1).

Figure 1. Indication of corpus availability in OPAC.

The gods of India, a brief description of their history, character & worship, by E. Osborn Martin. With 68 illustrations and map.

By: [Martin, E. Osborn](#)

Material type: Text

Publisher: London, New York, [Dent](#); [Dutton](#), 1914

Description: xviii, 330 p. front., illus. (map) plates. 20 cm

Subject(s): [Hindu gods](#) | [Hindu mythology](#)

DDC classification: 294.5/211

LOC classification: BL2001 | .M3

Online resources: [Click here to access online](#) | [Click here to access online](#)

Summary: This book has ventured to classify and group the under three separate heads, deities of a giving Hinduism introductory brief chapter to the third section, and concluding with a chapter summing up the whole subject and defining the scope and value of Hindu mythology.

Summary: This corpus has 1 document with 106,490 total words and 16,191 unique word forms. Vocabulary Density: 0.152. Average Words Per Sentence: 26.6. Most frequent words in the corpus: **Gods** (437); Worship (338); India (325); Hindu (313); **God** (293)

Tags from this library: No tags from this library for this title. [Log in to add tags.](#)

★★★★★ Average rating: 0.0 (0 votes)

Holdings (1)
Title notes (3)
Book reader
Text analysis
Comments (0)
Similar Items

Item type	Current location	Call number	Copy number	Status	Notes	Date due
 Books	Library2 <i>General Stacks</i> (<i>Library 2</i>)	294.5/211 MAR/G (Browse shelf)	c1	Available	Text analysis available	

Figure 2. Integration of book reader and text analysis tabs in OPAC.

library materials in mainstream languages as well as for Unicode-compliant minority languages (like tribal languages of India).

4. STATEMENT OF PROBLEM

The above-stated objective related to design and development of an enhanced OPAC to support language analysis as well as visualisation of the analysis results is based on the following three interlinked research questions:

RQ 1: How and to what extent is it possible to customise

OPAC of a selected web-enabled open source ILS so that it can display results of language analysis (provided through a back-end text analytics tool) for text corpus of a retrieved bibliographic record in the same interface along with other OPAC data (like holdings, availability, notes etc)?

RQ 2: What should be the parameters for selection of open source ILS and open source language analysis tool? How can these two software be linked to generate language analysis services in OPAC for documents having corpus and availability of the corpus is indicated at the metadata

1913) if available in the library catalogue can lead to Indic-script based language analysis services in OPAC (Fig. 5).

7. EXTENSION OF MECHANISM TO A TRIBAL LANGUAGE

The definition of minority language varies but in most of the cases a minority language is a language which is spoken by less than 50 per cent of a population in a given region, state or country³⁰. Obviously the size of the speaker population in a conferred geographic area is the key characteristic of a minority language and interestingly a minority language in one area can be a majority language in another region³⁰⁻³¹. The report on language (Table C-6, published in 2018) by Census of India includes two parts – a) 22 languages that are included in the 18th Schedule of the Indian Constitution; and b) those not included in it (total 99). The group ‘a’ includes two major tribal languages (listed in the Eighth Schedule) namely Bodo and Santali. Many of the tribal languages in India have shown negative growth in compare with 2001 census including two major tribal languages Bodo and Santali (included in the Eighth Schedule).

This paper has selected Santali to extend the language analysis mechanism as discussed in the foregoing sections (from two tribal languages included in the Eighth Schedule of the Constitution namely Bodo and Santali) as a minority tribal language of India. The reasons are :

- Santali is spoken by 69,73,345 people whereas Bodo is spoken by 14,54,547 people^{32,33};
- Santali has now its own script Ol Chiki (included in Unicode standard 5.1.0, released in April, 2008 and continues - [http://unicode.org/versions/Unicode 5.1.0](http://unicode.org/versions/Unicode%205.1.0) – Fig. 6) whereas Bodo is written in Devnagri script;
- Ol Chiki script has now many Open Type Fonts (OTF) supports like Sakal Bharati (from Technology Development for Indian Languages, Ministry of Electronics & Information Technology, Govt. of India), Noto Sans Ol Chiki (developed by Google) etc. All the major web browsers now support display in Ol Chiki and thereby full-text display in OPAC is an easy task now with the help of a local flip-book reader tool (Fig. 7);
- Santali has many literary and cultural heritage objects and recently made significant stride towards mainstream Indian languages;
- Santali language Wikipedia in Ol Chiki script (sat.wikipedia.org) was launched in August 2018 after obtaining necessary approval from Wikimedia Language Committee.

The technical factors as mentioned in point 2,3 and 4 are key ingredients in extending the language analysis mechanisms as applied in mainstream languages (Fig. 1 through Fig. 5) to Santali language resources in Unicode-complaint Ol Chiki script without much trouble (Fig. 8).

8. CONCLUSIONS

Language documentation helps to stop the endangerment of languages and holds its cultural diversity. UNESCO atlas of the world's languages in danger³⁴ reported that India has

most endangered languages in the world (197 endangered languages) and the Endangered Languages Project³⁵ identified 201 endangered Indian languages. Krauss³⁶ predicted that up to 90 per cent of all languages will be extinct by the end of this century. Library can play a vital role for the activities related to integration of information retrieval with language analysis and other facets of language documentation. All the tools and techniques applied in this research work are available as open source and can be implemented in libraries of any type or size in developing enhanced OPAC.

REFERENCES

1. Boerger, B.H.; Self, S.N.; Moeller, S.R. & Reiman, D.W. Language and culture documentation manual, 2020. <https://leanpub.com/languageandculturedocumentationmanual> (accessed on 13 June 2020).
2. Gippert, J.; Himmelmann, N.P. & Mosel, U. Essentials of language documentation. Mouton de Gruyter, Berlin, 2006. doi: 10.1515/9783110197730.
3. Austin, P.K. Data and language documentation. In Essentials of language documentation, edited by Jost Gippert, Nikolaus P. Himmelmann and Ulrike Mosel. Mouton de Gruyter, Berlin, 2006, 87–112. doi: 10.1515/9783110197730.
4. Austin, P.K. Training for language documentation: The SOAS experience. In The International Conference on Endangered Austronesian Languages, 5–7 June, 2007, Providence University, Taiwan. 2007. pp. 223-34. https://www.researchgate.net/publication/238506014_Training_for_Language_Documentation_the_SOAS_experience (accessed on 03 June 2020).
5. Himmelmann, N.P. Documentary and descriptive linguistics. *Linguistics*, 1998, **36**(1), 161–95. doi: 10.1515/ling.1998.36.1.161.
6. Woodbury, A.C. Defining documentary linguistics. In Language documentation and description (volume 1), edited by Peter K. Austin. Hans Rausing Endangered Languages Project, SOAS, London, 2003, 35–51. http://www.dl.cnr.fr/Colloques/3L_2008/3LCourseMaterial/Woodbury_2003_Documentation.pdf (accessed on 03 June 2020)
7. Woodbury, A.C. Language documentation. In the Cambridge handbook of endangered languages, edited by Peter K. Austin & Julia Sallabank. Cambridge University Press, Cambridge, 2011, 159-86. doi: 10.1017/CBO9780511975981.009.
8. Furbee, N.L. Language documentation: Theory and practice. In Language documentation: Practice and values, edited by Lenore A. Grenoble & N. Louanna Furbee. John Benjamins Publishing Company, Amsterdam, 2010, 3-24. doi: 10.1075/z.158.04fur.
9. Hill, J.H. The ethnography of language and language documentation. In Essentials of language documentation, edited by Jost Gippert, Nikolaus P. Himmelmann & Ulrike Mosel. Mouton de Gruyter, Berlin, 2006, 113–28. doi: 10.1515/9783110197730.113.

10. Varela, M.E. & Lee, N.H. Language documentation: A reference point for theatre and performance archives? *Int. J. Perform. Arts Digital Media*, 2018, **14**(2), 17-33. doi: 10.1080/14794713.2018.1453242.
11. Austin, P.K. Current issues in language documentation. In *Language documentation and description*, vol. 7, edited by Peter K. Austin. SOAS, London, 2010, 12-33. <http://www.e-publishing.org/PID/080> (accessed on 04 June 2020).
12. D'Andrade, R. A folk model of the mind. In *Cultural models in language and thought*, edited by D. Holland & N. Quinn. Cambridge University Press, 1987, 112-48. doi: 10.1017/CBO9780511607660.006.
13. D'Andrade, R. The identification of schemas in naturalistic data. In *Person schemas and maladaptive interpersonal patterns*, edited by Mardi J. Horowitz. University of Chicago Press, Chicago, 1991, 279-301
14. Mckee, A. A beginner's guide to textual analysis. *Metro magazine: Media & education Magazine*, 127/128, 2001, 138-49. <https://eprints.qut.edu.au/41993/> (accessed on 11 June 2020).
15. Frey, L.; Botan, C. & Kreps, G. *Investigating communication: An introduction to research methods* (2nd ed.). Allyn & Bacon, Boston, 1999. http://mason.gmu.edu/~afinn/html/teaching/courses/f03_comm250/fbk_chapters/09.pdf (accessed on 11 June 2020)
16. Zaenen, A. & Uszkoreit, H. Overview of language analysis and understanding. In *Survey of the state of the art in human language technology*, edited by Ron Cole *et al.* Cambridge University Press, 1997, 112-13. <http://www.dfki.de/~hansu/HLT-Survey.pdf> (accessed on 04 June 2020).
17. Karlsson, F. & Karttunen, L. Sub-sentential processing of language analysis and understanding. In *Survey of the state of the art in human language technology*, edited by Ron Cole *et al.* Cambridge University Press, 1997, 113-14. <http://www.dfki.de/~hansu/HLT-Survey.pdf> (accessed on 04 June 2020).
18. Goldenstein, J.; Poschmann, P. & Händschke, S.G.M. Linguistic analysis: The study of textual data in management and organisation studies with NLP. *Acad. Manage.* 2007. doi: 10.5465/ambpp.2015.10882abstract.
19. Krippendorff, K. *Content analysis: An introduction to its methodology*. Thousand Oaks, SAGE, 2004, **2**, doi: 10.1177/1094428108324513.
20. Pollach, I. Taming textual data: The contribution of corpus linguistics to computer aided text analysis. *Organ. Res. Methods*, 2012, **15**(2), 263-87. doi: 10.1177/1094428111417451.
21. Bernard, H.R. & Ryan, G. Text analysis: Qualitative and quantitative methods. In *Handbook of methods in cultural anthropology*. 1998. http://www.analytictech.com/mb870/readings/bernard_ryan_text_analysis.pdf (accessed on 11 June 2020).
22. Text mining tools and methods. Illinois Library, <https://guides.library.illinois.edu/textmining> (accessed on 11 June 2020).
23. Murmu, C. *Ol Seched* (in Al Chiki). Pandit Raghunath Murmu Human Welfare Trust, Rairungpur, 2004.
24. Mohanta, B.K. Invention in Al Chiki script for Santali language: An attempt to preserve an endangered tribal dialect. In *Endangered cultures and languages in India: Empirical observations*, edited by G. K. Bera, K. Jose, & North Eastern Institute of Culture & Religion (Gauhati, India), 2015, 211-29.
25. Klessa, K. Language documentation. In *Book of knowledge of languages in danger*, edited by Michael Hornsby. 2014, 106-120. <http://languagesindanger.eu/book-of-knowledge/> (accessed on 03 June 2020)
26. Hundsberger, S. Foreign language learning in second life and the implication for resources provision in academic libraries. *ARCADIA*, 2009. <http://www.arcadiaproject.lib.cam.ac.uk> (accessed on 03 June 2020)
27. Van Esch, D.; Foley, B. & San, N. Future directions in technological support for language documentation. In *Proceedings of the 3rd Workshop on the Use of Computational Methods in the Study of Endangered Languages*, 26-27 February 2019, Honolulu, USA, 2019, 14-22. <https://journals.colorado.edu/index.php/computel/article/view/341> (accessed on 13 June 2020)
28. Acharya, S.; Sanyal, D.K.; Mazumdar, J. & Das, P.P. Archiving endangered mundā languages in a digital library. https://ndlproject.iitkgp.ac.in/sites/default/files/website_images/publications/Paper_Archiving%20Endangered%20Mund%C4%81%20Languages%20in%20a%20Digital%20Library.pdf (accessed on 13 June 2020).
29. Sinclair, S. & Rockwell, G. *Voyant Tools* (v. 2.4). <http://voyant-tools.org/> (accessed on 13 June 2020).
30. Grenoble, L.A. & Singerman, A.R. *Minority languages*. Oxford University Press, 2014. doi: 10.1093/OBO/9780199772810-0176.
31. Pandharipande, R. Minority matters: Issues in minority languages in India. *Int. J. Multicultural Soc.*, 2002, **4**(2), 213-34. www.unesco.org/shs/ijms/vol4/issue2/art4 (accessed on 13 June 2020).
32. Down to Earth. Language census: Many tribal tongues now have fewer takers. 2018 <https://www.downtoearth.org.in/news/environment/language-census-many-tribal-tongues-now-have-fewer-takers-61044> (accessed on 14 June 2020)
33. Office of the Registrar General, India. *Census of India 2011: Paper 1 of 2018, Language India, States and Union territories* (Table C-16). https://censusindia.gov.in/2011Census/C-16_25062018_NEW.pdf (accessed on 14 June 2020).
34. UNESCO atlas of the world's languages in danger: <http://www.unesco.org/languages-atlas/> (accessed on 03 June 2020).
35. Ethnologue: Languages of the world <https://www.ethnologue.com/> (accessed on 03 June 2020).
36. Krauss, Michael E. The World's languages in Crisis. *Language*, 1992, **68**(1), 4-10. doi: 10.1353/lan.1992.0052.

CONTRIBUTORS

Parthasarathi Mukhopadhyay is a Professor in the Department of Library Science, University of Kalyani. His contributions to this paper are - developing the technical architecture including the scripts required for integration of language analysis services in Koha OPAC, building mechanism for text analytics for OI Chiki script and finalizing the objectives, the research questions and the methodology for this research oriented publication.

Anirban Dutta is a Junior Research Fellow (JRF) of the Department of Library Science, University of Kalyani. He has completed his BLIS and MLIS from the same university. He has performed the following activities for this research work – review of literature, overview of the related initiatives, metadata encoding for software framework, and testing & debugging of the language analysis services.