

A Platform for Searching Texts for Desired Expressions in a User-editable Pattern Matching Environment for Language Learning

Tatsuya Katsura ^{*}, Koichi Takeuchi [†]

Abstract

In this paper we propose a platform of pattern matching system that can extract required phrases or sentences in texts. Finding certain expressions in texts are often needed in language learning, e.g, examples of case markers between a predicate and an argument, or possible nouns in subject of a verb in a certain meaning. In previous studies, several types of systems, containing concordancers, are proposed. However, it is not easy to apply combined patterns because the pattern matching templates are previously fixed. Thus, we propose a flexible phrase searching system in which the users can create search patterns by combining blocks of basic search templates. One of the characteristics the proposed system is that the user can also specify where to be highlighted in texts with the blocks. To realize the function of combining patterns by the users, the proposed system employs Prolog as the base of the data structure. The platform of the searching system is implemented on an Web server with JavaScript-based interface and database system. In the performance test, we shows that the proposed system can deal with relatively large scale texts (10,000 sentences), and also demonstrate the combined patterns can be applied to the texts. In this paper we discuss the system architecture and the extendability of the pattern matching.

Keywords: Pattern matching, Concordancer, Block-based programming, Prolog

1 Introduction

Extraction of some phrases and expressions in texts is considered essential function in language education. For example, case markers located between a predicate and an object in Japanese language are various rules and alternations¹, then language learners need to search for examples of predicate-argument examples in Japanese texts. As a tool for searching texts, various kinds of concordancer are proposed, however, most concordancers are

^{*} Graduate School of Environmental, Life, Natural Science and Technology, Okayama University, Okayama, Japan

[†] Faculty of Environmental, Life, Natural Science and Technology, Okayama University, Okayama, Japan

¹The verbs *morau* (get) and *eru* (obtain) are almost the same meaning but *morau* has alternation between the case markers *ni/kara* as *kare ni/kara morau* ((I) get (it) from him.), and *eru* can only takes *kara* in this meaning as *kare kara/*ni eru*.

targeted for English. From the results of natural language processing study, several dependency parsers for Japanese are available, thus it is possible to compose a complex text searching system if the language learners can build NLP tools. But language learners are not NLP or software engineer, thus, we need an environment to build patterns by searching for phrases or expressions without needing programming.

when user want to obtain some verb-nouns or predicate preposition, concordancer is available the basic functions character, word matchina are available, but more extendede cases are no method.

, however, more combined search, for example, semantic meaning of predicate buy-sell and linguistic alternation, kara/wo in Japanese dependency

日本語に対するパターンマッチングシステム

2 Related Studies

An IIAI electronic copyright form should accompany your final submission or registration of the conference. The instruction of copyright transfer is announced at the conference webpage. Authors are responsible for obtaining any security clearances for the copyright form submission.

3 Platfrom of Pattern Matching System

user editable pattern with visual handling, Construction of pattern is not easy, the user should apply the assumed patterns, see the results and then tune up the patterns. Thus, editing space and results should be exist.

Patterns are combinations, combining by basic blocks and extract every where matched.

4 References

4.1 Example

4.1.1 *Article in a collection*

4.1.2 *Article in a conference proceedings*

4.1.3 *Article in a journal or magazine*

4.1.4 *Blog*

4.1.5 *Book*

4.1.6 *Book series*

4.1.7 *Electronic publication(Article in a conference proceedings)*

4.1.8 *Electronic publication(Online-only publication)*

4.2 Abbreviations in References

5 Some Common Mistakes

Acknowledgments

Author can include an acknowledgement of this work here.

References

[1] Consolidating the IT Infrastructure, white paper, Oracle Corp., Dec. 2003.