# pre-CODIE – Crosslingual On-Demand Information Extraction

**Kiyoshi Sudo**     **Satoshi Sekine**     **Ralph Grishman**
Computer Science Department
New York University
715 Broadway, 7th Floor,
New York, NY
{sudo,sekine,grishman}@cs.nyu.edu

## 1  Introduction

Our research addresses two central issues of information extraction — portability and multilinguality. We are creating information extraction (IE) systems that take foreign-language input and generate English tables of extracted information, and that can be easily adapted to new extraction tasks. We want to minimize the human intervention required for customization to a new scenario (type of facts or events of interest), and allow the user to interact with the system entirely in English. As a prototype, we have developed the pre-CODIE system, an experimental *C*rosss-lingual *O*n-*D*emand *I*nformation *E*xtraction system that extracts facts or events of interest from Japanese source text without requiring user knowledge of Japanese.

## 2  Overview

To minimize the customization of the IE system across scenarios, the extraction procedure of pre-CODIE is driven by the query from the user. The user starts the procedure by specifying the type of facts or events of interest in the form of a narrative description, and then pre-CODIE customizes itself to the topic by acquiring extraction patterns based on the user's description. Pre-CODIE, as an early attempt at a fully-automated system, still needs user interaction for template definition and slot assignment; automating these steps is left as future work.

Pre-CODIE interacts with its user entirely in English; even for slot assignment of the extraction patterns, the system translates the Japanese extraction patterns, which are based on subtrees of a dependency tree (Sudo et al., 2001), by word-to-word translation of each lexical item in the patterns. For ease of use, the Japanese extraction patterns are not only translated into English, but also shown with translated example sentences which match the pattern.

## 3  System Architecture

Pre-CODIE is implemented as an integration of several modules, as shown in Figure 1: translation, information retrieval, pattern acquisition, and extraction.
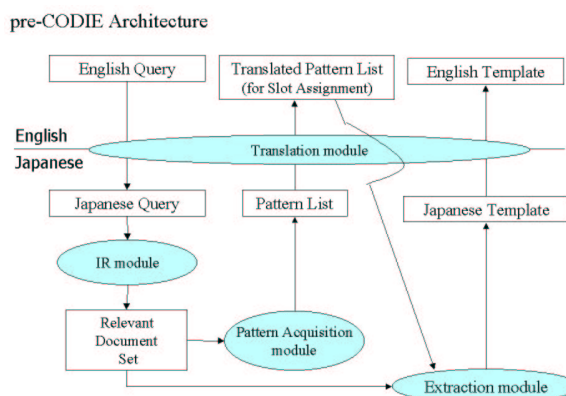


Figure 1: system architecture

First, the system takes the narrative description of the scenario of interest as an English query, and the Translation module (off-the-shelf IBM *King of Translation* system) generates a Japanese query. The IR module retrieves a set of *relevant documents* from Japanese Mainichi Newspaper from 1995. Then, the Pattern Acquisition module produces a list of extraction patterns from the *relevant document set* sorted by their relevance to the scenario (Sudo et al., 2001). Pre-CODIE asks the user to assign each placeholder in the patterns to one of the slots in the template. Finally, the Extraction module performs the pattern matching with slot-assigned patterns to each text in the *relevant document set* and generates a filled Japanese template, which is translated slot-by-slot into English for the user.

## 4 An Example procedure: Management Succession

From the user's point of view, pre-CODIE works as follows with the screenshots in Figure 2.

1. **Query:** The user types in the narrative description of the scenario of interest, one phrase in "description" text-box and more detail optionally given in "narrative" text-box: *"Management Succession: ..."*.

2. **Configuration:** The user adds and/or deletes the slots in the template; Add *"Person-In"*, *"Person-Out"*, *"Post-In"*, *"Post-Out"*, and *"Organization"*.

3. **Slot Assignment:** The user assigns a slot to each placeholder in the pattern by choosing one of the slots defined in step 2; Assign *"(be-promoted (PERSON-SBJ))"* to *"Person-In"*.

   Also, the user can see the example sentences with the match of the pattern highlighted. This will make it easier for the user to understand what each pattern aims to extract.

4. **Extraction:** The user gets the extracted template and repeats this procedure until the user gets the right template by going back to step 3 to change and/or add slot assignments, and by going back to step 2 to delete and/or add slots in the template.

### Acknowledgments

### References

Kiyoshi Sudo, Satoshi Sekine, and Ralph Grishman. 2001. Automatic pattern acquisition for japanese information extraction. In *Proceedings of the Human Language Technology Conference (HLT2001)*, San Diego, California.
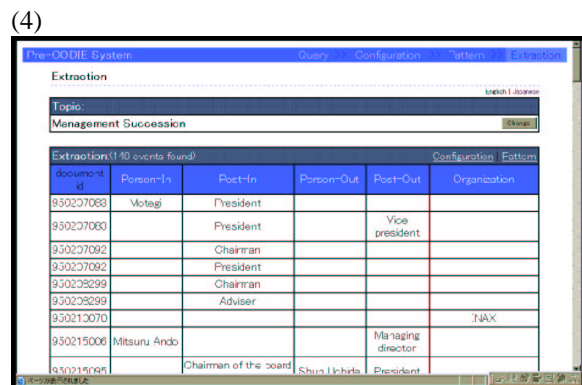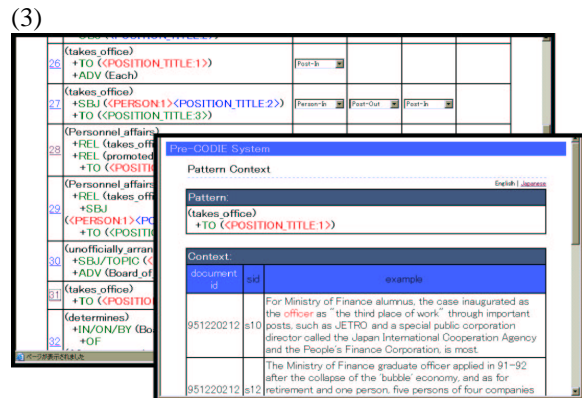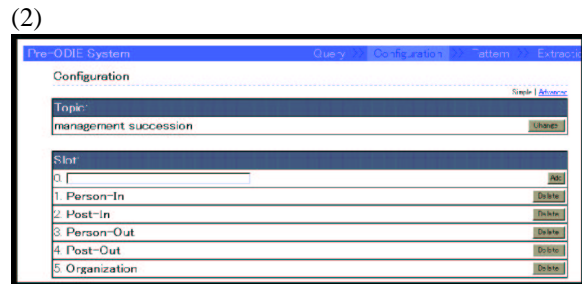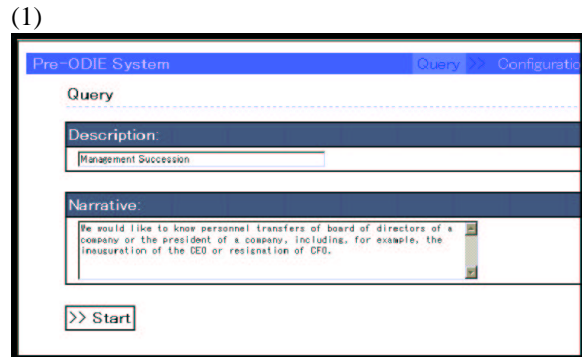


Figure 2: Screenshots of an Example procedure: Each image corresponds to the procedure in Section 4.