# ICE-TEA: an Interactive Cross-language Search Engine with Translation Enhancement

Dan Wu
School of Information Management
Wuhan University, Hubei 430072, China

woodan.cn@gmail.com

Daqing He
School of Information Sciences
University of Pittsburgh, PA 15260, USA

dah44@pitt.edu

## Categories and Subject Descriptors

H.3.3 [**Information Storage and Retrieval**]: Information Search and Retrieval—*Relevance Feedback*

## General Terms: Design, Experimentation

## Keywords: CLIR, Translation Enhancement (TE), Query Expansion (QE), Relevance Feedback (RF), ICE-TEA

## EXTENDED ABSTRACT

In cross-language information retrieval (CLIR), relevance feedback (RF) has been demonstrated to be effective in improving retrieval results, especially when reliable RF information can be obtained from users. Though query expansion (QE) is the leading RF approach in CLIR and it can take place before or/and after translating the query, it should not be the only possible RF method. In our demonstration, besides an implementation of post-translation QE, we also implement a novel RF approach called translation enhancement (TE) and the integration of TE and QE.

ICE-TEA is an interactive multilingual information access system with complete relevance feedback techniques. It is developed to utilize the RF information from users to conduct TE, QE, and the integration of them to improve CLIR results. ICE-TEA performs query translation based CLIR. It is built on top of Indri, a leading open-source monolingual search engine. As illustrated in Figure 1, the followings are the main functions of ICE-TEA:

**Translation Enhancement.** As discussed in iCLEF, and as demonstrated by Google's cross-language search engine, returned documents in CLIR need to be translated back to the query language side so that they can be selected and examined by the users. The main idea of TE, therefore, is to extract intended translation relationships of query terms from the relevant document pairs (original returned documents and their translations) and then to enhance the query translation. Through a set of experiments, we design a TE method that acquires translation relationships in relevant document pairs by word alignment obtained using GIZA++. Such translation relationships are used to re-calculate the corresponding translation probabilities of the query terms, so that the translated query obtained from the new translation probabilities will be closer to the true meaning of the user's request.

**Query Expansion.** In current version of ICE-TEA, we implement post-translation QE, which has been proved to perform better than pre-translation QE. Our system utilizes Indri's QE mechanism

which is an adaptation of Lavrenko's relevance models. More interestingly, since QE and TE occur at different stages of the CLIR process, ICE-TEA also combines both QE and TE.

**User-assisted Query Translation with Relevance Judgment.** Following the user-assisted query translation idea [1], ICE-TEA invites users into CLIR process. It lets users type queries and set cumulative probability threshold (CPT) so that initial CLIR search can be initiated. It also displays the original returned documents and their translations (both surrogate and full-text) to help users to make multi-level relevance judgments. After RF, it displays both TE and QE results, and helps users to modify easily the query translations by selecting or deselecting translation alternatives and their probabilities.
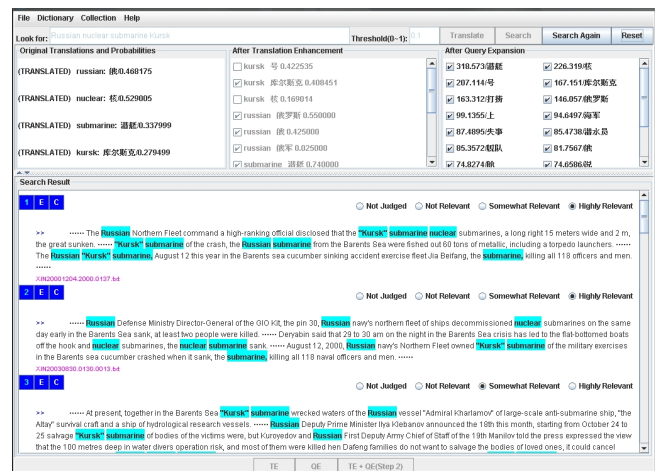


**Figure 1: ICE-TEA's English-to-Chinese Search Interface**

ICE-TEA right now is implemented to conduct CLIR search between English queries and Chinese collections. However, the system's architecture is generic enough to be converted for any major language pairs. The basic required language resource is a bilingual dictionary. However, with a machine translation system or the ability to access a machine translation service between the document and the query language, the translated returned document will be in better quality, so does the TE results.

## REFERENCE

[1] He, D., Oard, D. W., Wang, J., et al. Making MIRACLEs: Interactive Translingual Search for Cebuano and Hindi. *ACM Transactions on Asian Language Information Processing.* 2003, 2(3): 219-244.