

# Improve Cross Language Information Retrieval with Pseudo-Relevance Feedback

Lam Tung Giang  
Danang People Committee  
Vietnam

Vo Trung Hung  
University of Danang  
Vietnam

Huynh Cong Phap  
University of Danang  
Vietnam

**Abstract**— In dictionary-based Cross-language Information Retrieval systems, structured query translation has been shown to be an useful method for improving system performance. In this paper, we examine the effects of using pseudo relevance feedback to refine the structured query in the target language. We propose different methods for term weighting based on word distributions and the mutual information between expanded terms and original query terms. Our experimental results in a dictionary-based Vietnamese-English CLIR system show that while changing query terms weights has effects on improving precision, query expansion improves recall rates. The combination of these two techniques helps to improve system performance up to 12%, in terms of Mean Average Precision.

**Keywords**— CLIR, dictionary-based, Vietnamese, structured query, Pseudo-relevance feedback, reweight query terms, query expansion

## I. INTRODUCTION

Cross-Language Information Retrieval (CLIR) has been an important research field with the role to allow users to access documents in languages different from that of query[1][2]. A common approach in CLIR is to translate queries using dictionaries because of the simplicity and the availability of machine readable bilingual dictionaries [3][4].

A major problem in dictionary-based CLIR systems is ambiguity. Given a query in the source language, the translated query in the target language is built by selecting the “correct” translations from a list of candidate translations for each term in the initial query [3]. There are two mutually exclusive techniques to address this problem. The single selection technique tries to find one best translation for each term. The multiple selection technique, on the other hand, builds a structured query in the target language.

This work is motivated by our previous results in query translation by building a structured query in the target language from a given query [5]. In this paper, we apply pseudo-relevance feedback technique to improve the structured query translation by learning weights for query terms from top documents returned by the initial retrieval. Our experimental results in a dictionary-based Vietnamese-English CLIR system show that this method helps to improve precision up to 7%. We also examine the effects of query expansion in the target language. Different methods of calculating weights for expanded terms are examined. Our experiments show that query expansion contributes only a minor improvement in precision, however it helps to retrieve more relevant documents. The combination of using query terms re-weighting and query expansion together is shown as a good solution,

when it helps to improve system performance up to 12%, in term of Mean Average Precision [6].

The article is structured as follows. Section 2 presents several works in pseudo relevance feedback and review our previous work related structured query translation. In section 3, we propose and evaluate effects of query terming re-weighting and query expansion. Section 4 presents and analyses the experimental results. Section 5 presents the conclusions of our study.

## II. RELATED WORKS

### A. Pseudo relevance feedback

Relevance feedback (RF) was introduced in Rocchio's work [7], where the author introduced a formula for forming a new query vector by maximizing its similarity to relevant documents and minimizing its similarity to non-relevant documents in the collection. Initially, this technique was applied for vector space model and uses user feedbacks on the relevance of documents retrieved from the initial ranking and tries to automatically refine the query.

Since real user feedbacks are hard to obtain, Pseudo-Relevance Feedback (PRF) is used as an alternative solution[8]. PRF assumes the top  $n$  documents from initial retrieval as being relevant and uses these pseudo-relevant documents to refine the query for the next retrieval. Due to its automatic manner and effective performance, PRF has been widely applied in different IR frameworks like vector space models, probabilistic IR and language modeling.

Traditional PRF approaches recalculate query term weights based on statistics from retrieved documents and the collection such as term frequency  $tf$ , document frequency  $df$ , or term frequency-inverse document frequency  $tf-idf$ . For example, each term in retrieved documents is assigned an expansion weight  $w(t, D_r)$  as the mean of term weights in each retrieved document [9]:

$$w(t, DR) = \frac{\sum_{d \in D_r} w(t, d)}{R} \quad (1)$$

where  $R$  is the number of retrieved documents,  $w(t, d)$  is the frequency of a term  $t$  in document  $d$ . These term weights then are used to define new query by Rocchio's formula:

$$Q_{new} = \alpha \cdot Q + \beta \cdot \sum_{r \in D_r} \frac{r}{R} \quad (2)$$

Here,  $Q$  and  $Q_{new}$  represent original and new queries,  $D_r$  is the set of pseudo-relevant documents,  $R$  is the number of

retrieved documents,  $r$  is the expansion term weight vector,  $\alpha$  and  $\beta$  are tunable parameters.

More recently, expansion techniques have been introduced within the language modeling framework [10]. In Language Modeling, the probability of a document given a query  $P(d|q)$  is estimated using the Bayes' rule:

$$P(d|q) = \frac{P(q|d) \cdot P(d)}{P(q)} \propto P(q|d) \cdot P(d) \quad (3)$$

Usually,  $P(d)$  is assumed to be uniform in the collection and query terms are assumed independent with each other. The probability  $P(d|q)$  then is equivalent with:

$$P(d|q) \propto P(q|d) = \prod_{t \in q} P(t|d) \quad (4)$$

The values  $P(t|d)$  is traditionally computed using Dirichlet smoothing.

The basic idea of Relevance-Based (RM) approach is to estimate a better query language model by using information given by the pseudo-relevant documents and query terms [11][12]. Following formulas (3) and (4), we have:

$$\begin{aligned} P(w|q) &= \sum_d P(w|d) \cdot P(d|q) \\ &= \sum_d P(w|d) \cdot P(q|d) \\ &= \sum_d P(w|d) \cdot \prod_{t \in q} P(t|d) \end{aligned} \quad (5)$$

Different extensions of Relevance-Based (RM) query model try to eliminate the irrelevant terms, which appear frequently in both returned documents and in the whole collection. The following formula interpolates the probability of a term in the collection [12][13]:

$$P(w|q') = \lambda \cdot P(w|q) + (1 - \lambda) \cdot P(w|C) \quad (6)$$

Here,  $\lambda$  is a tunable parameter for controlling irrelevant terms. Top  $n$  terms with highest  $P(w|q')$  are selected to add into the new query.

Other developments for improving query expansion in the literature include approaches using external information, such as Wikipedia [14], Wordnet [15] or query logs of web search[16]; using machine learning to find good feedback documents [17] and good expansion terms [18][19].

### B. Pseudo relevance feedback for CLIR

For cross-lingual information retrieval, PRF can be applied in different retrieval stages of pre-translation, post-translation or the combination of both with the aim of increasing retrieval performance [3][20]. The PRF strategy gives an average improvement across query topics. It works well if there are many relevant documents retrieved in the initial top  $n$ , but is less successful when the initial retrieval effectiveness is poor[21].

Hiemstra [22] builds a structured query in the target language by grouping translations for each query term in the source language with the same weight and then use relevance feedback to learn translation probabilities. Daqing and Dan[23] proposes a Translation Enhancement method to expand queries

and to enhance query translation by adjusting translation probabilities.

### C. Building structured query in the target language

In this part, we review our previous results in query translation by building a structured query in the target language from a given query [5].

Given a Vietnamese query, we apply a heuristic method combining the use of dictionaries and a tagger tool for extracting keywords from a given query in Vietnamese, then propose methods based on Mutual Information to select  $n$  best translation candidates for each Vietnamese keyword from the dictionaries (we set  $n=5$ ) and then to build a structured query in English. In the English structured query, the translation of each query term in the source language is a group of English words, containing the best candidate assigned weight 1 and other candidates assigned weight 0.5. Each group also is assigned a weight depending on the tag assigned to the term in the source language by a tagger tool.

Formally, given a Vietnamese query  $q_v = (v_1, \dots, v_n)$ , the translated query in English has the following form:

$$\begin{aligned} q_e &= ((e_{1,1}, w_{1,1}) \text{ OR} \dots \text{OR} (e_{1,m_1}, w_{1,m_1}); w_1) \text{ AND} \\ &\dots \\ &\text{AND} ((e_{n,1}, w_{n,1}) \text{ OR} \dots \text{OR} (e_{n,m_n}, w_{n,m_n}); w_n) \end{aligned} \quad (7)$$

Here  $v_1, \dots, v_n$  are query terms in the given query, values  $m_1, \dots, m_n$  are numbers of translation candidates of  $v_1, \dots, v_n$ . Each pair  $(e_{j,k}, w_{j,k})$  contains a translation candidate and its assigned weight for a word  $v_j$ . Each value  $w_j$  is assigned weight for the group containing translation candidates of the term  $v_j$ .

For instance, from the Vietnamese query *quản lý quản lý quản lý quản lý* (*process* *sản xuất* *production*), we get the query translation (*management*<sup>1</sup> *OR* *regulate*<sup>0.5</sup> *OR* *control*<sup>0.5</sup>)<sup>2</sup> (*method*<sup>1</sup> *OR* *process*<sup>0.5</sup> *OR* *instruction*<sup>0.5</sup>)<sup>4</sup> (*production* <sup>1</sup> *OR* *manufacture*<sup>0.5</sup> *OR* *fabricate*<sup>0.5</sup>)<sup>2</sup>, which is used in Solr search engine.

In our previous work, the terms weights in the structured query are assigned heuristically. It is obviously not the best method.

## III. OUR PROPOSED APPROACH

In this article, we examine the effects of using pseudo-relevance feedback for refining the structured query in the target language. Given a query  $q$  in the form of formula (7) and the collection of pseudo relevant documents returned from the initial retrieval, we propose an algorithm for reweighting query terms in the target language and adding new terms to build a new query. The algorithm consists of 5 separated steps:

Step 1: Calculate query terms weights based on term distribution in pseudo relevant documents.

Step 2: Change the query terms weight, using the new query to retrieve a new list of pseudo relevant documents.

Step 3: Select top  $m$  popular terms contained in new retrieved documents ( $m=100$ ).

Step 4: Calculate terms weights and select  $n$  terms with highest weights ( $n=5, 10, 15, 20, 25$ ).

Step 5: Build the expand query by adding these  $n$  terms.

### A. Calculate query term weights

Denoting  $R$  as the set of returned documents in the initial retrieval. The query term weight for a term contained in query  $q$  is calculated by:

$$w(t) = \sum_{d \in R} \text{score}(d) \frac{\text{count}(t, d)}{\text{length}(d)} \quad (8)$$

Where  $\text{count}(t, d)$  is the number of times a term  $t$  appears in a document  $d$  contained in  $R$ ,  $\text{length}(d)$  is the length of the document  $d$ ,  $\text{score}(d)$  is the score assigned to the document  $d$  by the search engine.

The term weights calculated by formula (8) are used to reformulate a new query. At this step, the new query has the following form:

$$q' = (e_{1,1}w_{1,1} \text{ OR } \dots \text{ OR } e_{1,n1}w_{1,n1}) \text{ AND } \dots (e_{n,1}w_{n,1} \text{ OR } \dots \text{ OR } e_{n,nm}w_{1,nm}) \quad (9)$$

Here  $e_{ij}$  and  $w_{ij}$  are similar with those in formula (7). Please note that the weights assigned for translation groups in formula (7) are removed in the new query  $q'$ .

### B. Select top popular terms

With the set  $R$  of returned documents in the second retrieval, all documents  $d_i$  in  $R$  are vectorized. In the result, a dictionary  $D = \{t_1, \dots, t_{|D|}\}$  is created, containing all terms belonging documents in  $R$ . Each document is represented as a vector  $d_i = \{w_{i,1}, \dots, w_{i,|D|}\}$ , where  $w_{i,j}$  is the  $tf$ - $idf$  weight of the term  $t_j$  of the document  $d_i$  in the set  $R$ . For each term not contained in the query  $q$ , the term weight is calculated by the next formula:

$$w(t_j) = \frac{\lambda}{|R|} \sum_{d_i \in R} w_{i,j} \quad (10)$$

with  $\lambda$  is a turnable parameter (we set  $\lambda=1$  here). The formula (10) is used to select top  $m$  popular terms ( $m=100$  in our experiments) in the retrieved documents. The expanded terms will be chosen among these terms.

### C. Calculate new term weights

This part introduces 4 ways for calculating expanded term weights. The first way of calculating term weights simply uses the formula (10), denoted as FW1. The  $n$  terms with highest weights then are used to add into the new query.

The second way combines the local  $tf$ - $idf$  weight and the global  $idf$  weight of terms. For each term  $t_j$ , term weight is calculated by the formula FW2:

$$w(t_j) = \frac{\lambda}{|R|} \sum_{d_i \in R} w_{i,j} \cdot \log\left(\frac{N+1}{N_{t_i}+1}\right) \quad (11)$$

Here  $N$  is the total number of documents in the collection.  $N_{t_i}$  is the number of document containing the term  $t_i$ ,  $\lambda$  is a tunable parameter.

Based on the assumption that terms closer to the query terms are more likely to be relevant to the query topic, we propose the third way of calculating term weights uses mutual information of expanded terms and the initial query terms. First, we learn a "local" co-occurrence model of term pairs from the pseudo-relevant documents. For each term  $t_j$ , and a query term  $q_k$ , we denote  $mi(t_j, q_k)$  as the number of time two

these terms are in a distance of 3. The term weight is then calculated by formula FW3 as follows:

$$w(t_j) = \lambda \cdot \sum_{q_k \in q} mi(t_j, q_k) \quad (12)$$

Another way of using Mutual Information is building a "global" co-occurrence model of term pairs in the whole collection, then the term weight is calculated by formula FW4:

$$w(t_j) = \lambda \cdot \sum_{q_k \in q} mi(t_j, q_k) \cdot \log\left(\frac{N+1}{N_{q_k}+1}\right) \quad (13)$$

By adding top  $n$  terms with highest term weights, the final query has the following form:

$$q_{final} = q' \text{ AND (expanded terms)} \\ = (e_{1,1}w_{1,1} \text{ OR } \dots \text{ OR } e_{1,n1}w_{1,n1}) \text{ AND } \dots \\ \dots (e_{n,1}w_{n,1} \text{ OR } \dots \text{ OR } e_{n,nm}w_{1,nm}) \text{ AND } t_1w_1 \dots t_nw_n \quad (14)$$

Here  $e_{ij}$  and  $w_{ij}$  are similar with those in formula (7),  $t_i$  is expanded term and  $w_i$  is the weight assigned for  $t_i$ .

## IV. EXPERIMENTAL RESULTS

### A. Test configuration

To evaluate presented methods, we conduct the following experiment: first, we collect 24000 English documents from Web and build an English monolingual IR system on top of the open source search tool Solr<sup>1</sup>. We use 50 Vietnamese queries with an average length of 8,73 words for our experiment. At first, we apply query translation methods in [5] to build structured queries in English. After that, we follow the algorithm presented in section 3 to change query term weights. With top  $m$  100 popular terms ( $m=100$ ) from 50 retrieved documents, we use formulas FW1, FW2, FW3, FW4 defined in section III to calculate term weights. Top  $n$  terms ( $n=5, 10, 15, 20$  or  $25$ ) with highest weights are used for query expansion.

### B. Results

With each formula, we examine the system performance with different values of  $\lambda$ : 0.0001, 0.0002, 0.0005, 0.001, 0.002, 0.005, 0.01, 0.02, 0.05, 0.1, 0.2 and 0.5. It turns out that the tunable parameter  $\lambda$  is an important factor for selecting expanded terms. We conduct experiments with different values and select the final values for  $\lambda$  as 0.1, 0.01, 0.001 and 0.001 when applying formulas FW1, FW2, FW3 and FW4 respectively.

The next two tables show the MAP scores and the number of relevant documents being retrieved for different test configurations. In each table, the first row (used as the *baseline*) and the second row (marked as the *CW* configuration) show the performance when using original translated query and after changing query term weights. The next 4 rows show the performance when we expand query with 5, 10, 15, 20 or 25 top terms by applying different term weighting formulas FW1, FW2, FW3, FW4.

By changing query term weights using formula (8), the MAP score is improved 7%. With  $n=10$ , the query expansion method using the formula FW2 gives the best MAP score of

<sup>1</sup> <http://lucene.apache.org/solr/>

0.425, which is 112% of the *baseline* and 104% of *CW* configurations. The formula FW1 also gives a high MAP score of 0.421, which is 111% of the *baseline* and 103% of the *CW* configurations. It can be seen that query term reweighting is the main factor to contribute to the system precision.

TABLE I. MAP SCORES

		n=5	n=10	n=15	n=20	n=25
<i>Baseline</i>	0.380					
<i>CW</i>	0.407					
FW1		0.416	0.421	0.417	0.416	0.410
FW2		0.416	<b>0.425</b>	0.418	0.415	0.411
FW3		0.414	0.411	0.413	0.411	0.412
FW4		0.404	0.400	0.388	0.386	0.367

The results in the table 2 show that the number of retrieved relevant documents is declined in the *CW* configuration. However, this number is increased when we apply query expansion. The best result of 5179 retrieved documents from total 6109 relevant documents is reached with  $n=15$  and the formula FW4 based on global mutual information is used.

TABLE II. NUMBER OF RETRIEVED RELEVANT DOCUMENTS

		n=5	n=10	n=15	n=20	n=25
<i>Baseline</i>	4999					
<i>CW</i>	4961					
FW1		5044	5047	5075	5075	5071
FW2		5010	5067	5061	5082	5099
FW3		5081	5075	5095	5070	5072
FW4		5019	5004	5179	5098	5127

The figure 1 presents the interpolated 11-point average precision for 4 configurations: *baseline*, *CW*, FW1 and FW2 with  $n=10$ . It shows a clear advantages of the combination of proposed algorithms with term ranking formulas FW1 and FW2 over the baseline.

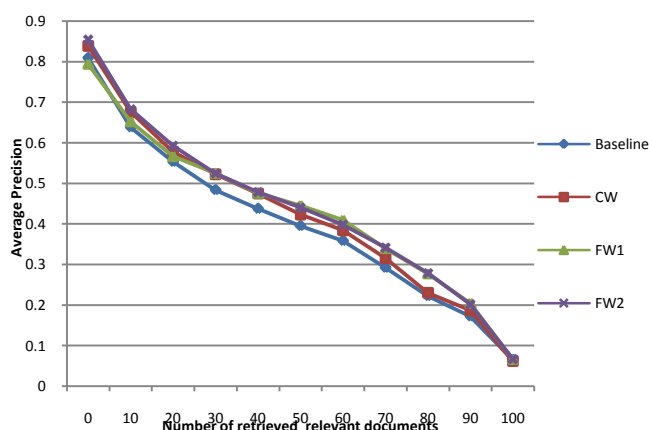


Fig. 1. Interpolated 11-point average precision (n=10)

## V. CONCLUSIONS

Pseudo-relevance feedback is proven to improve system performance in IR and CLIR. In this article, we propose a two-step PRF algorithm to separate query term re-weighting and query expansion. We examine different variants for terms weighting using returned documents from the initial retrieval, including using local *tf-idf* term weight, combining local *tf-idf* term weight and global *idf* weight, and using mutual information of terms and initial query terms. Our experimental results show that the combination of the proposed algorithm and weighting functions helps to improve system precision and recall rates.

## REFERENCES

- [1] Gerard Salton, Automatic processing of foreign language documents, J. Am. Soc. Inf. Sci., vol. 21, no. 3, pp. 187–194, 1970.
- [2] Jian-Yun Nie, Cross-Language Information Retrieval. 2010.
- [3] Dong Zhou, Mark Truran, Tim Brailsford, Vincent Wade, and Helen Ashman, Translation techniques in cross-language information retrieval, ACM Comput. Surv., vol. 45, no. 1, pp. 1–44, 2012.
- [4] Ari Pirkola, Turid Hedlund, Heikki Keskestalo, and Kalervo Järvelin, Dictionary-Based Cross-Language Information Retrieval: Problems, Methods, and Research Findings, Inf. Retr. Boston., vol. 4, no. 3, pp. 209–230, 2001.
- [5] Lam Tung Giang, Vo Trung Hung, and Huynh Cong Phap, Building Structured Query in Target Language for Vietnamese – English Cross Language Information Retrieval Systems, Int. J. Eng. Res. Technol., vol. 4, no. 04, pp. 146–151, 2015.
- [6] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze, Introduction to Information Retrieval. Cambridge: Cambridge University Press, 2008.
- [7] J. J. Rocchio, Relevance Feedback in Information Retrieval, SMART Retr. Syst. Experiments Autom. Doc. Process., 1971.
- [8] Gerard Salton and Chris Buckley, Improving retrieval performance by relevance feedback, J. Am. Soc. Inf. Sci., vol. 41, no. 4, pp. 288–297, 1990.
- [9] Zheng Ye, Ben He, Xiangji Huang, and Hongfei Lin, Revisiting Rocchio's relevance feedback algorithm for probabilistic models, Lect. Notes Comput. Sci., vol. 6458 LNCS, pp. 151–161, 2010.
- [10] John Lafferty and Chengxiang Zhai, Document language models, query models, and risk minimization for information retrieval, Proc. 24th Annu. Int. ACM SIGIR Conf. Res. Dev. Inf. Retr. SIGIR 01, vol. 01, pp. 111–119, 2001.
- [11] Victor Lavrenko and W. Bruce Croft, Relevance-Based Language Models, Sigir'01, p. 8, 2001.
- [12] Chengxiang Zhai and John D. Lafferty, Model-based Feedback In The Language Modeling Approach To Information Retrieval, Cikm, pp. 403–410, 2001.
- [13] Stéphane Clinchant and Eric Gaussier, A Theoretical Analysis of Pseudo-Relevance Feedback Models, Proc. 2013 Conf. Theory Inf. Retr. - ICTIR '13, pp. 6–13, 2013.
- [14] Yang Xu, Gareth J. F. Jones, and Bin Wang, Query Dependent Pseudo-Relevance Feedback based on Wikipedia, Sigir, 2009.
- [15] M. Voorhees and Road Nj, Expansion using lexical-semantic relations, SIGIR'94 Proc. 17th Annu. Int. ACM SIGIR Conf. Res. Dev. Inf. Retr., pp. 61–69, 1994.
- [16] Ben He and Iadh Ounis, Combining fields for query expansion and adaptive query expansion, Inf. Process. Manag., vol. 43, no. 5, pp. 1294–1307, 2007.
- [17] B. He and I. Ounis, Finding good feedback documents, pp. 2–5, 2009.

- [18] Guihong Cao, Jian-Yun Nie, Jianfeng Gao, and Stephen Robertson, Selecting good expansion terms for pseudo-relevance feedback, Proc. 31st Annu. Int. ACM SIGIR Conf. Res. Dev. Inf. Retr. - SIGIR '08, p. 243, 2008.
- [19] Saúl Vargas, Rlt Santos, Craig Macdonald, and I. Ounis, Selecting Effective Expansion Terms for Diversity, OAIR 2013 10th Int. Conf. RIAO Ser., 2013.
- [20] Kazuaki Kishida, Technical issues of cross-language information retrieval: a review, *Information Processing & Management*, vol. 41, no. 3, pp. 433–455, 2005.
- [21] M. Sanderson and P. Clough, Measuring pseudo relevance feedback & CLIR, pp. 484–485, 2004.
- [22] Djoerd Hiemstra, Wessel Kraaij, Ren´ee Pohlmann, and Thijs Westerveld, Translation Resources , Merging Strategies , and Relevance Feedback for Cross-Language, *CrossLanguage Inf. Retr. Eval. Work. CrossLanguage Eval. Forum CLEF 2000*, pp. 102–115, 2000.
- [23] Daqing He and Dan Wu, Translation enhancement: a new relevance feedback method for cross-language information retrieval, Proc. 17th ACM Conf. Inf. Knowl. Manag., pp. 729–738, 2008.