

# Toshiba BRIDJE at NTCIR-4 CLIR: Monolingual/Bilingual IR and Flexible Feedback

Tetsuya Sakai      Makoto Koyama      Akira Kumano  
Toshihiko Manabe  
Knowledge Media Laboratory, Toshiba Corporate R&D Center  
tetsuya.sakai@toshiba.co.jp

## Abstract

*Toshiba participated in the Monolingual/Bilingual tasks at NTCIR-4 CLIR using our CLIR system called BRIDJE. We submitted 24 runs covering three languages (Japanese, English, Chinese) and six language pairs, and achieved the highest performances in the E-J-D, C-J-D, C-J-T, E-E-D, J-E-D, J-E-T subtasks. Based on our formal run results, this paper discusses (a) the feasibility of the MT-based pivot language approach; (b) the effectiveness of our new Flexible Pseudo-Relevance Feedback methods; and (c) the advantages of Q-measure, which is a recently proposed retrieval performance metric based on multigrade relevance.*

**Keywords:** BRIDJE, pivot language, Flexible Pseudo-Relevance Feedback, Q-measure.

## 1 Introduction

Toshiba participated in the Monolingual/Bilingual tasks at NTCIR-4 CLIR using our CLIR system called BRIDJE [13, 14]. The objectives of our participation this year were: (a) To study the feasibility of the *Pivot Language* (or Transitive Translation) approach [1, 3] using Machine Translation (MT) systems; and (b) To devise new methods for *Flexible Pseudo-Relevance Feedback* [8, 9, 10, 11, 12]. In addition, this paper has a third purpose: (c) To show the advantages of *Q-measure*, which is a recently proposed retrieval performance metric based on multigrade relevance [17].

We submitted 24 runs covering three languages (Japanese, English, Chinese) and six language pairs. In addition, there were 12 runs which we generated but could not submit, because we were only allowed to submit up to two runs for each language pair/topic field (i.e. TITLE or DESCRIPTION). (We did not submit a fifth run by mixing different topic fields because we believe that this is not practical.) Table 1 provides a summary of our official and unofficial runs. As indicated in the “Top Performer” rows, we achieved the

highest performances in the E-J-D, C-J-D, C-J-T, E-E-D, J-E-D, J-E-T subtasks, and “silver medal” performances for most of the other tasks, including C-E-D and C-E-T for which we used a pivot language approach. Throughout this paper, we prefer to use the Unofficial Names listed in the third column of this table, as they better reflect the search strategies used.

The remainder of this paper is organised as follows. Section 2 describes the search request translation process of our bilingual runs, including the pivot runs, and briefly discusses their effectiveness. Section 3 introduces two Flexible Pseudo-Relevance Feedback methods and discusses their effectiveness using our monolingual results. It also discusses the advantages of Q-measure as a retrieval performance metric based on multigrade relevance. Finally, Section 4 concludes this paper. We report on our work for the NTCIR-4 QAC2 task in a separate paper [16].

## 2 Search Request Translation

### 2.1 BRIDJE and MT

The BRIDJE Cross-Language Information Access System [13, 14] accepts Japanese or English search requests and retrieves documents from Japanese or English text databases using the Okapi/BM25 algorithm [18]. All of our NTCIR-4 runs used the *default* Okapi parameter values [11]: that is, we did not tune the Okapi parameters at all. Our *traditional* Pseudo-Relevance Feedback (PRF) runs used the *offer weight* (*ow*) for term selection, with  $P = 10$  pseudo relevant documents and  $T = 40$  expansion terms [9, 10, 11, 13]. The algorithms for generating our *flexible* PRF runs will be described in Section 3.

For our E-J and J-E runs, the search requests were simply translated using the Toshiba MT System as in our previous work [7, 11, 13]. For our C-J runs, a new Chinese-Japanese MT system that is currently being developed at Toshiba was used for search request translation. As this new system is not yet complete, its translation quality is not as good as our English-

**Table 1. TSB Formal Run Results at NTCIR-4 CLIR.**

| Topic Field   | Official Name                        | Unofficial Name      | Relaxed MAP                  | Rigid MAP                    | Description   |
|---|--------------------------------------|----------------------|------------------------------|------------------------------|---|
| (a) Monolingual Japanese runs (55 topics)                                 |                                      |                      |                              |                              |   |
| DESC  | Top Performer at NTCIR-4             |                      | 0.4838                       | 0.3804                       |   |
|   | TSB-J-J-D-01                         | J-J-D-PRF            | 0.4759                       | 0.3667                       | Traditional PRF   |
|   | TSB-J-J-D-03<br><i>not submitted</i> | J-J-D-TE<br>J-J-D-SS | 0.4683<br><b>0.4854</b>      | 0.3578<br><b>0.3677</b>      | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| TITLE   | Top Performer at NTCIR-4             |                      | 0.4864                       | 0.3890                       |   |
|   | TSB-J-J-T-02                         | J-J-T-PRF            | 0.3863                       | 0.2834                       | Traditional PRF   |
|   | TSB-J-J-T-04<br><i>not submitted</i> | J-J-T-TE<br>J-J-T-SS | 0.3829<br><b>0.4538</b> ↑↑↑↑ | 0.2802<br><b>0.3460</b> ↑↑↑↑ | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| (b) English-Japanese runs using E-J MT (55 topics)                        |                                      |                      |                              |                              |   |
| DESC  | Top Performer at NTCIR-4             |                      | 0.3688 (BRIDJE)              | 0.2674                       |   |
|   | TSB-E-J-D-01                         | E-J-D-PRF            | <b>0.3688</b>                | 0.2672↑                      | Traditional PRF   |
|   | TSB-E-J-D-03<br><i>not submitted</i> | E-J-D-TE<br>E-J-D-SS | 0.3620<br>0.3673             | 0.2615<br><b>0.2715</b>      | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| TITLE   | Top Performer at NTCIR-4             |                      | 0.3525                       | 0.2735                       |   |
|   | TSB-E-J-T-02                         | E-J-T-PRF            | 0.3244                       | 0.2388                       | Traditional PRF   |
|   | TSB-E-J-T-04<br><i>not submitted</i> | E-J-T-TE<br>E-J-T-SS | 0.3134<br><b>0.3486</b>      | 0.2284<br><b>0.2557</b>      | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| (c) Chinese-Japanese runs using C-J MT (55 topics)                        |                                      |                      |                              |                              |   |
| DESC  | Top Performer at NTCIR-4             |                      | 0.3008 (BRIDJE)              | 0.2309 (BRIDJE)              |   |
|   | TSB-C-J-D-01                         | C-J-D-PRF            | 0.2986                       | 0.2269                       | Traditional PRF   |
|   | TSB-C-J-D-03<br><i>not submitted</i> | C-J-D-TE<br>C-J-D-SS | <b>0.3008</b><br>0.2997      | <b>0.2309</b><br>0.2282      | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| TITLE   | Top Performer at NTCIR-4             |                      | 0.3193 (BRIDJE)              | 0.2458 (BRIDJE)              |   |
|   | TSB-C-J-T-02                         | C-J-T-PRF            | 0.3193                       | <b>0.2458</b>                | Traditional PRF   |
|   | TSB-C-J-T-04<br><i>not submitted</i> | C-J-T-TE<br>C-J-T-SS | 0.3055<br><b>0.3198</b>      | 0.2324<br>0.2423             | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| (d) Monolingual English runs (58 topics)                                  |                                      |                      |                              |                              |   |
| DESC  | Top Performer at NTCIR-4             |                      | 0.4368 (BRIDJE)              | 0.3469 (BRIDJE)              |   |
|   | TSB-E-E-D-01                         | E-E-D-PRF            | <b>0.4368</b>                | 0.3469                       | Traditional PRF   |
|   | TSB-E-E-D-03<br><i>not submitted</i> | E-E-D-TE<br>E-E-D-SS | 0.4242<br>0.4366             | 0.3381<br><b>0.3510</b>      | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| TITLE   | Top Performer at NTCIR-4             |                      | 0.4512                       | 0.3576                       |   |
|   | TSB-E-E-T-02                         | E-E-T-PRF            | <b>0.4404</b>                | 0.3500                       | Traditional PRF   |
|   | TSB-E-E-T-04<br><i>not submitted</i> | E-E-T-TE<br>E-E-T-SS | 0.4274<br>0.4378             | 0.3367<br><b>0.3522</b>      | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| (e) Japanese-English runs using J-E MT (58 topics)                        |                                      |                      |                              |                              |   |
| DESC  | Top Performer at NTCIR-4             |                      | 0.4227 (BRIDJE)              | 0.3340 (BRIDJE)              |   |
|   | TSB-J-E-D-01                         | J-E-D-PRF            | <b>0.4227*</b>               | <b>0.3340</b>                | Traditional PRF   |
|   | TSB-J-E-D-03<br><i>not submitted</i> | J-E-D-TE<br>J-E-D-SS | 0.4110<br>0.4105             | 0.3253<br>0.3288             | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| TITLE   | Top Performer at NTCIR-4             |                      | 0.4262 (BRIDJE)              | 0.3407 (BRIDJE)              |   |
|   | TSB-J-E-T-02                         | J-E-T-PRF            | <b>0.4262</b>                | <b>0.3407</b>                | Traditional PRF   |
|   | TSB-J-E-T-04<br><i>not submitted</i> | J-E-T-TE<br>J-E-T-SS | 0.4218<br>0.4074             | 0.3369<br>0.3336             | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| (f) Chinese-English <i>pivot</i> runs using C-J MT and J-E MT (58 topics) |                                      |                      |                              |                              |   |
| DESC  | Top Performer at NTCIR-4             |                      | 0.2829                       | 0.2238                       |   |
|   | TSB-C-E-D-01                         | C-E-D-PRF            | 0.2767                       | 0.2183                       | Traditional PRF   |
|   | TSB-C-E-D-03<br><i>not submitted</i> | C-E-D-TE<br>C-E-D-SS | 0.2753<br><b>0.2862</b>      | 0.2169<br><b>0.2303</b>      | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |
| TITLE   | Top Performer at NTCIR-4             |                      | 0.2879                       | 0.2380                       |   |
|   | TSB-C-E-T-02                         | C-E-T-PRF            | 0.2873                       | 0.2207                       | Traditional PRF   |
|   | TSB-C-E-T-04<br><i>not submitted</i> | C-E-T-TE<br>C-E-T-SS | 0.2780<br><b>0.2969</b> ↑    | 0.2114<br><b>0.2370</b> ↑    | Flexible PRF ( Term Exhaustion )<br>Flexible PRF ( Selective Sampling ) |

Based on the Sign Test, TE/SS runs that are significantly better than the corresponding PRF run are indicated by ↑ ( $\alpha = 0.05$ ) and ↑↑ ( $\alpha = 0.01$ ). PRF/SS runs that are significantly better than the corresponding TE run are indicated by ↑ ( $\alpha = 0.05$ ) and ↑↑ ( $\alpha = 0.01$ ). PRF/TE runs that are significantly better than the corresponding SS run are indicated by \* ( $\alpha = 0.05$ ) and \*\* ( $\alpha = 0.01$ ). **Boldface** values indicate the best average performance within each language-pair/topic field.

Japanese and Japanese-English MT systems. For our C-E runs, we tried a *pivot language* approach instead of using a Chinese-English MT system: The Chinese requests were first translated into Japanese using the new Chinese-Japanese MT system, and the translated requests were further translated into English using our Japanese-English MT system. In short, this is a “Japanese as a pivot language” experiment.

## 2.2 Analysis of Bilingual Runs

Table 2 shows the *relative* performance values of our cross-language runs based on traditional PRF, where, for example, E-J-D-PRF and C-J-D-PRF are compared with the corresponding monolingual baseline J-J-D-PRF. For the E-J and C-J runs, the percentages are considerably higher for the TITLE runs than for the DESCRIPTION runs, due to the fact that the absolute performance of J-J-D-PRF was much better than that of J-J-T-PRF. C-J-D-PRF is considerably less effective than E-J-D-PRF because our Chinese-Japanese MT system is not yet as sophisticated as our English-Japanese one. We expect this difference to disappear eventually as we continue to improve our Chinese-Japanese MT system. (However, note that no such performance difference is visible for the TITLE runs, i.e., C-J-T-PRF vs E-J-T-PRF.)

On the other hand, Table 2 shows that our J-E runs are comparable to the monolingual baselines. That is, our Japanese-English MT did an excellent job. Because of this, our pivoted (i.e. C-E) runs are also reasonably successful: the relative performance of C-E-D-PRF is comparable to that of C-J-D-PRF. Recall that our C-E runs were generated by using Chinese-Japanese MT first, and then Japanese-English MT: As the second MT did not introduce much noise, our Chinese-English translations were almost as good as the Chinese-Japanese ones. Note also that our pivot runs are among the very best C-E runs (Table 1 (f)). Thus, we can conclude that the Pivot Language approach using *good* MT systems is feasible.

## 3 Flexible Feedback

### 3.1 Overview on Flexible Feedback

Traditional PRF relies on at least two parameters:  $P$  (the number of pseudo-relevant documents scooped from the top of the initial ranked output), and  $T$  (the number of expansion terms added to the initial query). Although PRF often improves *average* performance, it typically hurts one-third of a given set of search requests [8]. Various *Flexible PRF* methods have been proposed to enable *per-request adjustment* of these parameters [8, 9, 10, 11, 12], but the results have been inconclusive. Other researchers have also tackled this problem but without success [2, 6].

For NTCIR-4 CLIR, we tried two new Flexible PRF methods for determining  $P$  for each search request, both of which are based on which of the query terms occur in the initially retrieved documents. Sections 3.2 and 3.3 describe these methods.

### 3.2 Term Exhaustion

Our first Flexible PRF method is called *Term Exhaustion*. The idea behind it is simple: Scan the initial ranked output from the top, examining the query terms contained in the retrieved documents. Stop when “novel” query terms (i.e. those that were not in the previous documents) appear to have run out.

Let  $P_{min}$  and  $P_{max}$  denote the minimum/maximum number of pseudo-relevant documents required, respectively. Then, the problem is to automatically determine, for each topic,  $P$  such that  $P_{min} \leq P \leq P_{max}$ . Let  $d(r)$  denote the document at Rank  $r$  in the initial ranked output, and let  $T(d(r))$  denote the set of initial query terms contained in  $d(r)$ . The algorithm shown in Figure 1 determines  $P$  based on Term Exhaustion. Based on our preliminary Japanese monolingual experiments with the NTCIR-3 test collection, we let  $P_{min} = 6$  and  $P_{max} = 20$  for *all* NTCIR-4 Term Exhaustion (TE) runs, including the ones with English documents. As for  $T$ , we simply let  $T = 40$  as in traditional PRF.

### 3.3 Selective Sampling

Our second method, *Selective Sampling*, is unlike any other Flexible PRF method in that it does not necessarily treat the top  $P$  documents as pseudo-relevant. That is, it can *skip* documents. The idea behind it is that there may be similar (and therefore redundant) documents among the top  $P$  documents, and it may be better in such a case to go further down the list to look for more “novel” documents.

In addition to  $P_{min}$  and  $P_{max}$ , we introduce the third parameter called  $P_{scope}$ , so that no more than  $P_{scope}$  documents are examined. The algorithm shown in Figure 2 returns a set of pseudo-relevant documents, namely  $S$ , obtained through Selective Sampling. (Thus, the number of pseudo-relevant documents  $P = |S|$ .) The essence of the algorithm is that it tries to avoid collecting too many documents with the same  $T(d(r))$ . For NTCIR-4, we used  $P_{min} = 3$ ,  $P_{max} = 10$ , and  $P_{scope} = 50$  for all Selective Sampling (SS) runs, again based on our Japanese monolingual experiments with the NTCIR-3 test collection. As with traditional PRF, we let  $T = 40$ . However, as mentioned earlier, these runs were not submitted due to the constraints on the number of runs.

**Table 2. Relative performance of the cross-language PRF runs.**

| Unofficial name | Relaxed MAP ratio | Rigid MAP ratio | Unofficial name | Relaxed MAP ratio | Rigid MAP ratio |
|-----------------|-------------------|-----------------|-----------------|-------------------|-----------------|
| E-J-D-PRF       | 77%               | 73%             | E-J-T-PRF       | 84%               | 84%             |
| C-J-D-PRF       | 63%               | 62%             | C-J-T-PRF       | 83%               | 87%             |
| J-E-D-PRF       | 97%               | 96%             | J-E-T-PRF       | 97%               | 97%             |
| C-E-D-PRF       | 63%               | 63%             | C-E-T-PRF       | 65%               | 63%             |

```

 $T_O = \phi;$ 
/*  $T_O$  is the set of query terms Observed already. */
 $i = 0;$ 
/*  $i$  is the number of consecutive documents that do not
contain a novel query term. */
for(  $r = 1; r \leq P_{max}; r++$  ){
  if(  $T(d(r)) - T_O == \phi$  ) /* no novel term in  $d(r)$  */
     $i++;$ 
  else /* at least one novel term in  $d(r)$  */
     $i = 0;$  /* start counting from scratch */
  if(  $i + 1 == P_{min}$  )
    return(  $r$  );
   $T_O = T_O \cup T(d(r));$ 
}
return(  $r$  );

```

**Figure 1. Determining  $R$  based on Term Exhaustion.**

### 3.4 New Evaluation Metrics: Q-measure and R-measure

This section briefly describes Average Weighted Precision (AWP), Q-measure and R-measure which we use in Sections 3.5 and 3.6 for analysing our monolingual Flexible PRF results.

At NTCIR, both *Rigid* and *Relaxed* Mean Average Precision are calculated for performance comparison, as Average Precision cannot handle multiple relevance levels. AWP (originally called *weighted average precision* [5]) proposed by Kando *et al.* can handle multigrade relevance *and* are arguably better than the original *Cumulative Gain* [4] as it avoids rank-based averaging [15]. However, Sakai [17] has pointed out a problem with AWP, namely, that it does not give a reliable score if relevant documents are ranked below Rank  $R$ , where  $R$  is the number of known relevant documents. To solve this problem, Sakai has proposed *Q-measure*, which has the reliability of Average Precision *and* the capability of handling multigrade relevance. In addition, Sakai has proposed *R-measure*, which can be used along with Q-measure just like R-Precision is used besides Average Precision.

Formally, let  $gain(X)$  denote the *gain value* for successfully retrieving an  $X$ -relevant document. (We let  $gain(S) = 3$ ,  $gain(A) = 2$ ,  $gain(B) = 1$  throughout this paper.) Let  $L$  denote the size of the ranked output, and let  $X(r)$  denote the relevance level of the

```

 $S = \phi;$ 
/*  $S$  is the set of Sample documents that will be
treated as pseudo-relevant. */
for(  $r = 1; r \leq P_{scope}; r++$  ){
  if( is_good_sample_document(  $r$  ) )
     $S = S \cup d(r);$ 
  if(  $|S| == P_{max}$  )
    return(  $S$  );
}
return(  $S$  );

int is_good_sample_document(  $r$  )
{
   $i = 0;$ 
/*  $i$  is the number of previously seen documents with
the same set of query terms */
for(  $r' = 1; r' \leq r - 1; r'++$  )
  if(  $T(d(r')) == T(d(r))$  )
     $i++;$ 
if(  $i < P_{min}$  )
  return( 1 ); /* a good sample document */
else
  return( 0 ); /* NOT a good sample document */
}

```

**Figure 2. Obtaining the set of pseudo-relevant documents based on Selective Sampling.**

document at Rank  $r$  ( $\leq L$ ). Then, the *gain at Rank  $r$*  is given by  $g(r) = gain(X(r))$  if the document at Rank  $r$  is relevant, and  $g(r) = 0$  if it is nonrelevant. The *cumulative gain at Rank  $r$*  is given by  $cg(r) = g(r) + cg(r - 1)$  for  $r > 1$  and  $cg(1) = g(1)$ .

Let  $cig(r)$  represent the cumulative gain at Rank  $r$  for an *ideal* ranked output. (An ideal ranked output for NTCIR can be obtained by listing up all S-relevant documents, then all A-relevant documents, then all B-relevant documents.) Then, AWP is defined as:

$$AWP = \frac{1}{R} \sum_{1 \leq r \leq L, g(r) > 0} \frac{cg(r)}{cig(r)} \quad (1)$$

The problem with AWP arises from the fact that  $cig(r)$  remains constant for  $r \geq R$ . That is, AWP cannot discriminate between a relevant document at Rank  $R$  and one near the bottom of the ranked list. See [17] for more detailed discussions.

Let the *bonused gain at Rank  $r$*  be given by  $bg(r) =$

$g(r) + 1$  if  $g(r) > 0$  and  $bg(r) = 0$  if  $g(r) = 0$ , and its cumulative version be given by  $cbg(r) = bg(r) + cbg(r - 1)$  for  $r > 1$  and  $cbg(1) = bg(1)$ . Then, Q-measure is defined as:

$$Q\text{-measure} = \frac{1}{R} \sum_{1 \leq r \leq L, g(r) > 0} \frac{cbg(r)}{cig(r) + r} \quad (2)$$

Q-measure is free from the problem of AWP because the denominator  $cig(r) + r$  is guaranteed to increase with  $r$ . Note that this property resembles that of Average Precision, whose denominator is none other than  $r$ :

$$AveP = \frac{1}{R} \sum_{1 \leq r \leq L, g(r) > 0} \frac{count(r)}{r} \quad (3)$$

where  $count(r)$  is the number of relevant documents within Top  $r$ .

Finally, R-measure is defined as:

$$R\text{-measure} = \frac{cbg(R)}{cig(R) + R} \quad (4)$$

### 3.5 Analysis of Monolingual Runs

Table 3 summarises the results of our monolingual runs using the abovementioned metrics based on multigrade relevance. While the Term Exhaustion results are rather disappointing, the Selective Sampling results are very interesting: In particular, J-J-T-SS easily outperforms J-J-T-PRF, and the difference is statistically significant ( $\alpha = 0.01$ ) with the Sign Test as it is actually better than traditional PRF for around 45 topics out of 55 regardless of the performance metric. Unfortunately, however, the *English* Selective Sampling results are not as straightforward as the Japanese ones. In Section 3.6, we shall examine whether this difference is simply due to the fact that we tuned Selective Sampling using the NTCIR-3 *Japanese* test collection or not.

Although it is clear from the definitions that Q-measure is a more reliable performance metric than AWP, we first illustrate its superiority over AWP using actual data. Figure 3 provides a per-topic analysis of J-J-T-SS, which is the most successful Selective Sampling run: Each ‘‘circle’’ represents the value of Q-measure *minus* that of Relaxed Average Precision, while each ‘‘cross’’ represents the value of AWP *minus* that of Relaxed Average Precision. The horizontal axis represents the number of relevant documents  $R$ . From this graph, it is clear that the ‘‘circles’’ are closer to the horizontal axis than the ‘‘crosses’’, and therefore that the property of Q-measure resembles that of Average Precision more than AWP does. Moreover, it is clear that AWP *overestimates* the performance for topics with small  $R$ : This is because, as have been mentioned in Section 3.4, AWP is unreliable when relevant documents are found below Rank  $R$ .

To study the defect of AWP more closely, Table 4 provides some statistics for Topics 009 and 006, which correspond to the two ‘‘crosses’’ at the top left-hand corner of Figure 3. The table shows that the AWP values are over 0.5 even though Relaxed/Rigid Average Precision values are only around 0.1 and the Q-measure ones are around 0.2. Below, we use Topic 009 to illustrate how AWP overestimates performance for topics with small  $R$ .

From Table 4, an ideal ranked output for Topic 009 contains S-relevant documents from Rank 1 to 7, A-relevant documents from Rank 8 to 20, and B-relevant documents from Ranks 21 to 23. Therefore, the cumulative gain at Rank  $r (\geq 23)$  for this ideal list is  $cig(r) = 7 * 3 + 13 * 2 + 3 * 1 = 50$ .

Table 5 shows exactly how AWP and Q-measure are calculated for Topic 009 with J-J-D-SS, by listing up pertinent statistics for all  $r$  such that  $g(r) > 0$  (i.e. for every relevant document retrieved). Thus, AWP is calculated by dividing the sum of values in Column 4 by  $R = 23$ , while Q-measure is calculated by dividing the sum of values in Column 6 by  $R = 23$ . From this table, it is clear that  $cg(r)/cig(r)$  is not suitable for calculating retrieval performance: For example, even though the the twenty-third (i.e. the last) relevant document is at Rank 431,  $cg(431)/cig(431)$  is equal to one, as if to imply Perfect Precision. In contrast, it can be observed that  $bcg(r)/(cig(r)+r)$  imposes a penalty for going down the ranked list, just like Precision does for calculating Average Precision in a binary relevance environment.

### 3.6 Further Analysis of Selective Sampling

Having shown that Q-measure is a reliable evaluation metric, this section examines the Selective Sampling runs more closely using Q-measure.

In Table 3, Selective Sampling is very successful for the Japanese TITLE run, moderately successful for the Japanese DESCRIPTION runs, but not quite so for the English runs. To examine whether this difference arises from the fact that we tuned the Selective Sampling parameter  $P_{max}$  based on *Japanese* NTCIR-3 experiments, we generated some *additional* runs by varying  $P_{max}$  with Selective Sampling as well as varying  $P$  with traditional PRF.

Figures 4 and 5 show the results of the additional experiments, in which the horizontal axis represents  $P$  for the PRF runs and  $P_{max}$  for the Selective Sampling runs. (Our official results correspond to  $P = 10$  and  $P_{max} = 10$ , respectively.) From these graphs, it is clear that the performance is relatively stable with respect to the choice of  $P_{max}$  and  $P$ , and that Selective Sampling is more effective than traditional PRF regardless of the choice of  $P_{max}$  for the Japanese case. Thus, it is not the parameter setting that caused the difference between the Japanese and English results.

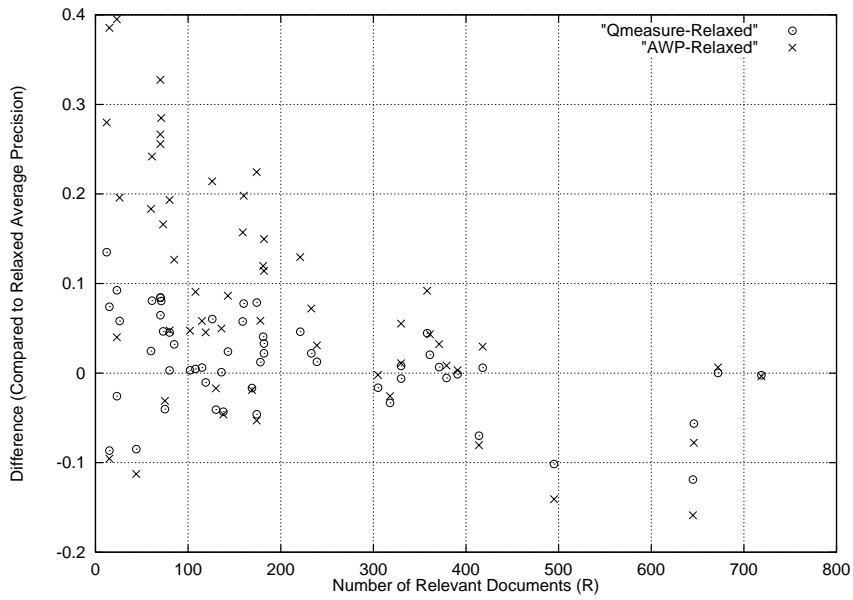


Figure 3.  $R$  vs Q-measure (AWP) minus Relaxed Average Precision.

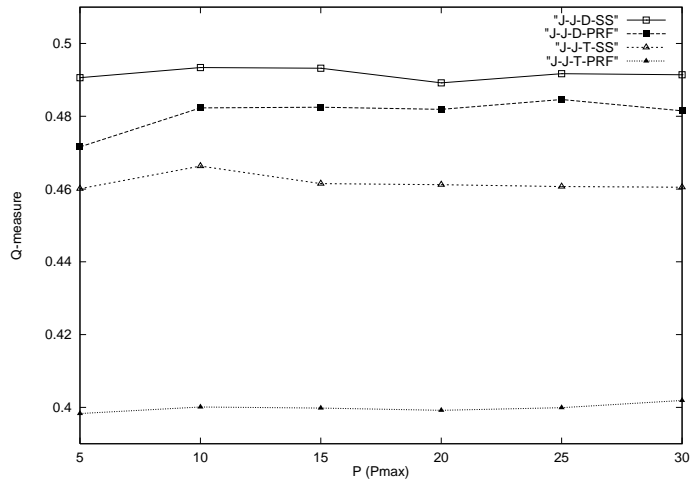


Figure 4. The effect of varying  $P$  ( $P_{max}$ ) for the J-J task.

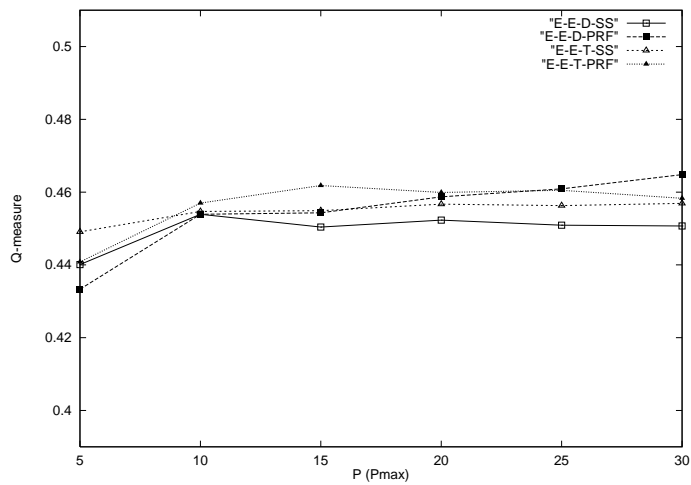


Figure 5. The effect of varying  $P$  ( $P_{max}$ ) for the E-E task.

**Table 3. The monolingual results in terms of Q-measure, R-measure, AWP and R-WP.**

| Official Name | Relaxed MAP        | Rigid MAP          | Q-measure          | AWP                | R-measure          |
|---------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| J-J-D-PRF     | 0.4759             | 0.3667             | 0.4823             | 0.5466             | 0.4997             |
| J-J-D-TE      | 0.4683             | 0.3578             | 0.4738             | 0.5360             | 0.4906             |
| J-J-D-SS      | <b>0.4854</b>      | <b>0.3677</b>      | <b>0.4934</b>      | <b>0.5597</b> ↑    | <b>0.5086</b>      |
| J-J-T-PRF     | 0.3863             | 0.2834             | 0.4001             | 0.4725             | 0.4350             |
| J-J-T-TE      | 0.3829             | 0.2802             | 0.3976             | 0.4718             | 0.4309             |
| J-J-T-SS      | <b>0.4538</b> ↑↑↑↑ | <b>0.3460</b> ↑↑↑↑ | <b>0.4663</b> ↑↑↑↑ | <b>0.5385</b> ↑↑↑↑ | <b>0.4816</b> ↑↑↑↑ |
| E-E-D-PRF     | <b>0.4368</b>      | 0.3469             | <b>0.4539</b>      | <b>0.5471</b>      | 0.4652             |
| E-E-D-TE      | 0.4242             | 0.3381             | 0.4430             | 0.5367             | 0.4532             |
| E-E-D-SS      | 0.4366             | <b>0.3510</b>      | <b>0.4539</b>      | 0.5461             | <b>0.4654</b>      |
| E-E-T-PRF     | <b>0.4404</b>      | 0.3500             | <b>0.4570</b> *    | <b>0.5449</b>      | <b>0.4717</b>      |
| E-E-T-TE      | 0.4274             | 0.3367             | 0.4423             | 0.5275             | 0.4612             |
| E-E-T-SS      | 0.4378             | <b>0.3522</b>      | 0.4547             | 0.5378             | 0.4696             |

The significance test results are given in the same way as in Table 1.

**Table 4. J-J-T-SS performance values for Topics 006, 009, 044 and 045.**

| Topic ID | $R$ | $R_S$ | $R_A$ | $R_B$ | Relaxed | Rigid  | Q-measure | AWP    |
|----------|-----|-------|-------|-------|---------|--------|-----------|--------|
| 006      | 15  | 0     | 11    | 4     | 0.1759  | 0.1168 | 0.2500    | 0.5615 |
| 009      | 23  | 7     | 13    | 3     | 0.1092  | 0.0868 | 0.2017    | 0.5043 |

One possible explanation for the above inconsistent behaviour of Selective Sampling would be that, as Selective Sampling tries to skip *redundant* documents in the initial ranked output, it works better with homogeneous document collection than with a heterogeneous one: the NTCIR-4 Japanese collection is composed of Mainichi and Yomiuri newspapers only, while the NTCIR-4 English collection is composed of Taiwan News, China Times English News, Mainichi Daily News, Korea Times, Xinhua, and Hong Kong Standard. We plan to test this hypothesis in the near future, by evaluating the effectiveness of Selective Sampling for other homogeneous/heterogeneous test collections.

We have tried to investigate the “degree of redundancy” in the initial ranked output, at least to some extent, by examining the average number of *skipped* documents for each Selective Sampling run. Table 6 summarises the results. Recall that our Selective Sampling runs picked up  $P = |S|$  pseudo-relevant documents from top  $P_{scope} = 50$  documents for each topic, such that  $P_{min} = 3 \leq P \leq P_{max} = 10$ . Let  $r_{last} (\leq P_{scope})$  be the rank of the  $P$ -th pseudo-relevant document that has been selected. Then, clearly, the number of skipped documents is given by  $r_{last} - P$ . For example, if the documents at Ranks 1,3,5 have been selected, then the number of skipped documents is  $5 - 3 = 2$ .

From Table 6, it appears that skipping more documents does not necessarily lead to more successful Selective Sampling, as E-E-T-SS skipped many documents but was not as effective as J-J-T-SS and J-J-D-SS. Thus our analysis is not sufficient for explaining when Selective Sampling works. On the other

hand, it appears that document skipping occurs more frequently with TITLE runs than with DESCRIPTION runs, probably because fewer query terms imply larger groups of similar documents. This observation is in agreement with the fact that Selective Sampling was more successful with Japanese TITLES than with Japanese DESCRIPTIONS.

## 4 Conclusions

Toshiba participated in the Monolingual/Bilingual tasks at NTCIR-4 CLIR. Our main findings are as follows:

1. The “Japanese as a pivot language” approach using *two* MT systems is feasible;
2. Flexible Feedback based on Selective Sampling is effective for the NTCIR-4 *Japanese* test collection, especially with the TITLE fields; and
3. Q-measure is a useful metric for evaluation with multigrade relevance.

As our Selective Sampling results for the NTCIR-4 *English* test collection were inconclusive, we plan to examine Selective Sampling using other homogeneous/heterogeneous test collections to clarify when it works and when it does not.

## References

- [1] Ballesteros, L. A.: Cross-Language Retrieval via Transitive Translation, In Croft, W. B. (ed.) *Advances in Information Retrieval: Recent Research*

**Table 5. AWP/Q-measure calculation for Topic 009 (J-J-D-SS).**

| $r$ | $cig(r)$ | $cg(r)$ | $\frac{cg(r)}{cig(r)}$ | $bcg(r)$ | $\frac{bcg(r)}{cig(r)+r}$ |
|-----|----------|---------|------------------------|----------|---------------------------|
| 12  | 31       | 1       | 0.0323                 | 2        | 0.0465                    |
| 19  | 45       | 3       | 0.0667                 | 5        | 0.0781                    |
| 37  | 50       | 4       | 0.0800                 | 7        | 0.0805                    |
| 41  | 50       | 6       | 0.1200                 | 10       | 0.1099                    |
| 43  | 50       | 9       | 0.1800                 | 14       | 0.1505                    |
| 46  | 50       | 11      | 0.2200                 | 17       | 0.1771                    |
| 48  | 50       | 14      | 0.2800                 | 21       | 0.2143                    |
| 52  | 50       | 16      | 0.3200                 | 24       | 0.2353                    |
| 56  | 50       | 19      | 0.3800                 | 28       | 0.2642                    |
| 69  | 50       | 21      | 0.4200                 | 31       | 0.2605                    |
| 88  | 50       | 23      | 0.4600                 | 34       | 0.2464                    |
| 91  | 50       | 26      | 0.5200                 | 38       | 0.2695                    |
| 103 | 50       | 28      | 0.5600                 | 41       | 0.2680                    |
| 117 | 50       | 30      | 0.6000                 | 44       | 0.2635                    |
| 126 | 50       | 32      | 0.6400                 | 47       | 0.2670                    |
| 141 | 50       | 34      | 0.6800                 | 50       | 0.2618                    |
| 168 | 50       | 37      | 0.7400                 | 54       | 0.2477                    |
| 179 | 50       | 38      | 0.7600                 | 56       | 0.2445                    |
| 196 | 50       | 41      | 0.8200                 | 60       | 0.2439                    |
| 276 | 50       | 43      | 0.8600                 | 63       | 0.1933                    |
| 309 | 50       | 45      | 0.9000                 | 66       | 0.1838                    |
| 338 | 50       | 48      | 0.9600                 | 70       | 0.1804                    |
| 431 | 50       | 50      | 1.0000                 | 73       | 0.1518                    |

**Table 6. Average number of skipped documents.**

| Unofficial name | #docs skipped |
|-----------------|---------------|
| J-J-D-SS        | 8.8           |
| J-J-T-SS        | 11.3          |
| E-E-D-SS        | 6.7           |
| E-E-T-SS        | 15.4          |

from the CIIR, pp. 203-234, Kluwer Academic Publishers (2000).

- [2] Billerdeck, B. and Zobel, J.: When Query Expansion Fails, *ACM SIGIR 2003 Proceedings*, pp. 387-388 (2003).
- [3] Gey, F. C., Jiang, H., Chen, A. and Larson, R. R.: Manual Queries and Machine Translation in Cross-language Retrieval and Interactive Retrieval with Cheshire II at TREC-7, *TREC-7 Proceedings*, pp. 527-540 (1999).
- [4] Järvelin, K. and Kekäläinen, J.: IR Evaluation Methods for Retrieving Highly Relevant Documents, *ACM SIGIR 2000 Proceedings*, pp. 41-48 (2000).
- [5] Kando, N., Kuriyama, K. and Yoshioka, M.: Information Retrieval System Evaluation using Multi-Grade Relevance Judgments - Discussion on Averageable Single-Numbered Measures (in Japanese), *IPSJ SIG Notes*, FI-63-12, pp. 105-112 (2001).
- [6] Murata, M. *et al.*: CRL at NTCIR2, *NTCIR-2 Proceedings*, pp. 119-129 (2001).

- [7] Sakai, T. *et al.*: Sakai, T., Kajiura, M., Sumita, K., Jones, G. and Collier, N.: A Study on English-to-Japanese / Japanese-to-English Cross-Language Information Retrieval using Machine Translation (in Japanese), *IPSJ Journal*, Vol. 40, No. 11, pp. 4075-4086 (1999).
- [8] Sakai, T., Kajiura, M. and Sumita K: A First Step towards Flexible Local Feedback for Ad hoc Retrieval, *IRAL 2000 Proceedings*, pp.95-102 (2000).
- [9] Sakai, T., Robertson, S. E. and Walker, S.: Flexible Pseudo-Relevance Feedback for NTCIR-2, *NTCIR-2 Proceedings*, pp.5-(59-66) (2001).  
<http://research.nii.ac.jp/ntcir/workshop/OnlineProceedings2/sakai.pdf>
- [10] Sakai, T., Robertson, S.E. and Walker, S.: Flexible Pseudo-Relevance Feedback via Direct Mapping and Categorization of Search Requests, *BCS-IRSG ECIR 2001 Proceedings*, pp.3-14 (2001).
- [11] Sakai, T.: Japanese-English Cross-Language Information Retrieval Using Machine Translation and Pseudo-Relevance Feedback, *IJCPOL*, Vol. 14, No. 2, pp. 83-107 (2001).
- [12] Sakai, T. and Robertson, S.E.: Flexible Pseudo-Relevance Feedback Using Optimization Tables, *ACM SIGIR 2001 Proceedings*, pp.396-397 (2001).
- [13] Sakai, T., Koyama, M., Suzuki, M. and Manabe, T.: Toshiba KIDS at NTCIR-3: Japanese and English-Japanese IR, *NTCIR-3 Proceedings* (2003).  
<http://research.nii.ac.jp/ntcir/workshop/OnlineProceedings3/NTCIR3-CLIR-SakaiT>
- [14] Sakai, T. *et al.*: BRIDGE over a Language Barrier: Cross-Language Information Access by Integrating Translation and Retrieval, *IRAL 2003 Proceedings*, pp.65-76 (2003). <http://acl.ldc.upenn.edu/W/W03/W03-1109.pdf>
- [15] Sakai, T.: Average Gain Ratio: A Simple Retrieval Performance Measure for Evaluation with Multiple Relevance Levels, *ACM SIGIR 2003 Proceedings*, pp. 417-418 (2003).
- [16] Sakai, T. *et al.*: Toshiba ASKMi at NTCIR-4 QAC2 *NTCIR-4 QAC2 Working Notes*, to appear (2004).
- [17] Sakai, T: New Performance Metrics based on Multi-grade Relevance: Their Application to Question Answering, *ACM SIGIR 2004 Proceedings*, submitted. (2004).
- [18] Sparck Jones, K., Walker, S. and Robertson, S. E.: A Probabilistic Model of Information Retrieval: Development and Comparative Experiments, *Information Processing and Management* 36, Part I (pp. 779-808) and Part II (pp. 809-840), (2000).