# Spoken Document Retrieval Experiments for SpokenQuery&Doc at Ryukoku University (RYSDT)

Hiroaki NANJO
Faculty of Science and
Technology, Ryukoku
University, Japan
nanjo@rins.ryukoku.ac.jp

Takehiko YOSHIMI
Faculty of Science and
Technology, Ryukoku
University, Japan
yoshimi@rins.ryukoku.ac.jp

Sho MAEDA
Faculty of Science and
Technology, Ryukoku
University, Japan
maeda@nlp.i.ryukoku.ac.jp

Tomohiro NISHIO
Graduate School of Science
and Technology, Ryukoku
University, Japan
nishio@nlp.i.ryukoku.ac.jp

## ABSTRACT

In this paper, we describe spoken document retrieval (SDR) systems in Ryukoku University, which were participated in NTCIR-11 "SpokenQuery&Doc" task. In NTCIR-11 SpokenQuery&Doc task, there are subtasks: "spoken content retrieval (SCR) subtask" and "spoken term detection (STD) subtask". We participated in the SCR and STD subtasks as team RYSDT. In this paper, our SDR and STD systems are described.

## Categories and Subject Descriptors

H3.3 [Information Storage and Retrieval]: Information Search and Retrieval

## Keywords

NTCIR-11, spoken content retrieval, spoken term detection

## 1. INTRODUCTION

With the advance of high-speed networks and mass storage, a great deal of audio contents, such as podcasts, TV news, home videos, and lecture/presentation videos, can be easily stored and published. In some universities, lecture videos are actually published in the open domain through web pages. Since such audio contents continue to increase, we need robust methods that can deal with them. Spoken content retrieval (SCR) and spoken term detection (STD), which process such huge amounts of spoken data for efficient search and browsing, are promising and are the most significant tasks of spoken document retrieval.

We have studied both STD [1] and SCR [2][3][4]. In NTCIR-11, "SpokenQuery&Doc" task is defined, which covers both SCR and STD. In this paper, our SCR and STD approaches and the results of their applications to NTCIR-11 task are described. Our team name is RYSDT in NTCIR-11 SpokenQuery&Doc [5].

## 2. SPOKEN CONTENT RETRIEVAL

### 2.1 SCR Overview and Its Problem

Spoken content retrieval (SCR) is a process finding the spoken document itself or short portions (passages) of spoken document which are relevant to the query. For the SCR task in NTCIR-11, a search target is Japanese oral presentations (lectures) in academic conferences. Since each lecture has a longer duration, just searching each lecture is not enough since we cannot access a specific scene which we want to know even if the suitable presentations are perfectly retrieved. Therefore, NTCIR-11 SCR defined slide group segment (SGS) unit retrieval task and short unit (passage) retrieval task. We tried to search both SGS unit and passage based on an orthodox vector space model (VSM).

We have investigated some VSM-based SCR techniques; query expansion (QE) based on pseudo relevance feedback (PRF)[6] and passage retrieving strategy using local and global document information[3].

PRF is one of the most well-known methods for improving IR accuracy. From IR results obtained by original query q, systems automatically select which documents are relevant or not, and modify a query vector $\mathbf{q}$ to $\mathbf{q_e}$. We have proposed PRFL (PRF for Lectures)[6], which is different from PRF in terms of that first pseudo relevant documents are retrieved with an original query from the index consisting of adequate length document unit and QE is performed, and then, final results are retrieved with the expanded query using an index consisting of retrieval unit. In a short spoken document retrieval, extracting relevant words from relevant documents is not difficult but retrieving relevant documents itself is difficult. The key for PRFL is to find relevant short spoken documents in the initial retrieval.

For the problem, we adopt our another proposed method; an integration method of global information and local information. Specifically, for detecting a short part of an oral presentation ("local document"), we integrate similarities between a given query and a longer unit, for example a whole presentation ("global document"), into a similarity between the given query and a local document contained in the global documents. Moreover, we adopt a search strategy that uses both similarity scores between a target document
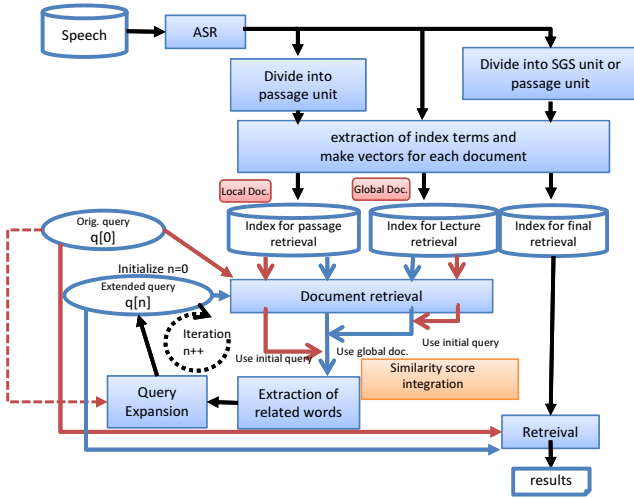
Figure 1: Overview of iterative PRF

and original/expanded queries.

The overview of our SCR system is shown in Figure 1. We set up some SCR systems which perform such PRFL iteratively.

## 2.2 SCR Techniques

### 2.2.1 Overview of Pseudo Relevance Feedback (PRF)

In a VSM-based IR system, a query is denoted by a vector whose each element represents a frequency of its corresponding word. In a VSM, a vector $q_e$ for an expanded query $Q_e$ is shown in equation (1).

$$q_e = q + \alpha q_a \qquad (1)$$

Here, q is a vector for the original query. $q_e$ is a vector of the expanded query. $q_a$ is a vector consisting of statistics of new index words to be added. $\alpha$ is a weight. Generally, $\alpha$ less than one is used. In this study, in order to set all elements of $q_e$ to be integer, QE is performed based on the equation (2) using positive integer $\beta$.

$$q_e = \beta q + q_a \qquad (2)$$

We regard top some documents which have higher IR score as pseudo documents, select additional words, and make a vector $q_a$. For $N$th iteration on PRF, top some documents retrieved with $N$-1th expanded query are used.

### 2.2.2 PRF for Lectures (PRFL)

PRFL differs from PRF in terms of that it uses special index for PRF. For a special index, we adopt 10 sequential utterance unit which is just divided oral presentation speech from its beginning into 10 sequential utterances segments. The kind of X-sequential utterance passage unit is introduced in the Japanese SCR test collection [7]. For $N$th iteration on PRFL, top some 10 sequential unit documents retrieved with $N$-1th expanded query are regarded as pseudo relevant document. Extracted words are added to original query and $N$th expanded query is generated.

For the final SGS unit search, we uses similarity scores between a target document and original/all expanded queries

(1st to $N$th). as follows (eq. 3).

$$\begin{aligned} \mathrm{SIM}(Q^0, \cdots, Q^N, D_i) &= (1 - \lambda) \log \mathrm{SMART}(Q^0, D_i) \\ &+ \lambda \mathrm{SIM}(Q^1, \cdots, Q^N, D_i) \end{aligned} \qquad (3)$$

when $M < N$

$$\begin{aligned} \mathrm{SIM}(Q^M, \cdots, Q^N, D_i) &= (1 - \lambda) \log \mathrm{SMART}(Q^M, D_i) \\ &+ \lambda \mathrm{SIM}(Q^{M+1}, \cdots, Q^N, D_i) \end{aligned}$$

and

$$\mathrm{SIM}(Q^N, D_i) = \log \mathrm{SMART}(Q^N, D_i)$$

Here, $Q^n, n = 1, \cdots, N$ are the $n$th expanded query and $Q^0$ is an initial query. $\mathrm{SIM}(Q^{0 \cdots N}, D_i)$ is an integrated score with all. $\lambda$ is a parameter for integration ($= 0.5$).

### 2.2.3 Document Retrieval with Global Information

We have proposed a method for retrieving short part of an oral presentation, that is, an integration method of global information and local information. Specifically, for detecting a short part of an oral presentation ("local document"), we integrate similarities between a given query and a longer unit ("global document"), into a query-local document similarity. In this paper, we denote the method "GDOC".

$$\begin{aligned} \mathrm{SIM}(Q, D_{i,_k}^{k+9}) &= (1 - \lambda') \log \mathrm{SMART}(Q, D_i) \\ &+ \lambda' \log \mathrm{SMART}(Q, D_{i,_k}^{k+9}) \end{aligned} \qquad (4)$$

Here, $Q$ is an original query, $D_i$ is $i$th document, $D_{i,_k}^{k+9}$ is a part of $i$th document form $k$th to $k+9$th utterances ($=10$ sequential unit). $\lambda'$ is a parameter for integration. Then, we find local documents which have higher similarity score $\mathrm{SIM}(Q, D_{i,_k}^{k+9})$.

### 2.2.4 PRFL with Global Document Information

In PRFL, relevant 10 sequential unit (short passage) is required for QE, however, the short passage search performance is not enough. Irrelevant short passages are wrongly used at QE process. Therefore, we adopt "GDOC" method for the short passage (local document) search. Moreover, at the search, we use an original query. So, in the QE process, we search local and global documents with original and expanded queries, and then, integrate four similarities to get more relevant short passages. In this paper, we denote such kind of PRF as "PRFL+GDOC".

For the final SGS (or passage) unit search, we use similarity scores between a target document and original/all expanded queries (1st to $N$th). The procedure is shown in Figure 1.

## 2.3 Submitted SCR Systems

### 2.3.1 Problem in Vector Space Modeling in Japanese SCR

For a VSM-based SCR, appropriate indexing is significant. Automatic speech recognition (ASR) is performed to make index terms, which essentially contain ASR errors. Therefore, studies of indexing terms that are robust to ASR errors are necessary. In Japanese text, no space is put between words, and word units are ambiguous. Thus,

studies of indexing units are also important. Based on this background, we have investigated several indexing units in Japanese SCR [2] including morpheme unit, character n-gram unit, and phone n-gram unit. We have found that morphemes is suitable for indexing unit and baseforms of nouns and verbs are suitable for index terms. For the NTCIR-11 SCR subtask in SpokenQuery&Doc, we applied above described VSM-based SCR systems.

In this work, the Generic Engine for Transposable Association (GETA)[8] is used for constructing VSM-based SCR systems.

### 2.3.2 SGS Retrieval Systems

First, SGS retrieval systems are described. In a SGS retrieval task, search target is each slide group segment. Transcription or ASR result for each SGS is regarded as a document and VSM-based SCR system is constructed. The submitted systems are listed as follows.

1. PRFL+GDOC iter 3 (set parameters on lecture retrieval task)

2. PRFL+GDOC iter 3

3. PRFL+GDOC iter 2

4. PRFL+GDOC

5. PRF+GDOC

6. GDOC

7. PRFL

8. PRF

9. BASELINE

Here, "BASELINE" system just searches SGS units with the given query, and is the same with NTCIR-11 organizer's baseline system [5].

Integration of original query and expanded query similarities are performed in all systems 1 to 8. In this work, all parameters (number of pseudo relevant documents/words, the weight for additional words on PRF $\beta$ in eq.(2), and combination weight $\lambda$ in eq.(3) and $\lambda'$ in eq.(4) were estimated so that higher 11ptAP was achieved at the past NTCIR SpokenDoc. We assume that each SGS consists of about 60 utterances, and the parameters are selected according to 60 sequential unit document retrieval results.

For constructing indices, we used "REF-WORD-MATCH" transcription which is given by NTCIR-11 organizer. For queries, we used "REF-WORD-MATCH" and "MANUAL" which are given by NTCIR-11 organizer.

### 2.3.3 SGS Retrieval Results

Text query experimet

First, results of SGS retrieval by text query are described. For each query, we tried to retrieve 1000 documents. The results are listed in Figure 2. Mean average precision (MAP) of the BASELINE is 0.159.

Both PRFL/PRF combined with BASELINE (system 7 and 8), which first retrieve 10-utterance/SGS unit and extract additional words from them, achieved comparable IR performances with the baseline. Using global docment information (system 6) improved IR performance, and the PRF+GDOC method (system 5) achieved highest MAP (0.235). PRFL+GDOC systems (system 1 to 4) achieved
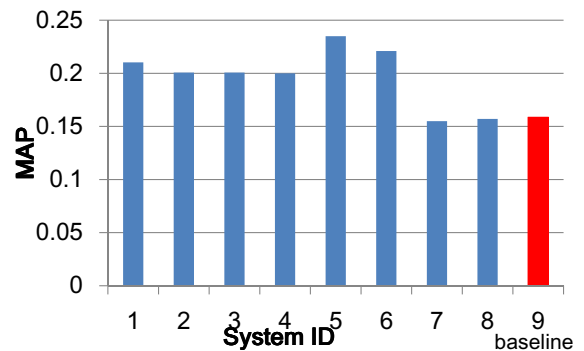


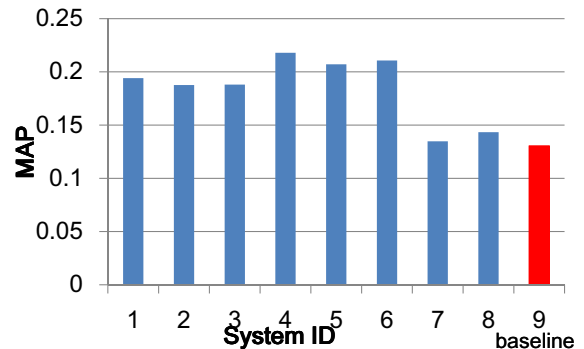Figure 2: SGS SCR result by text query



Figure 3: SGS SCR result by spoken query

higher performances than baseline system. Iterative PRF does not work well in this work.

Spoken query experiment

Next, results of SGS retrieval by spoken query are described. For each query, we also tried to retrieve 1000 documents. The results are listed in Figure 3. Mean average precision (MAP) of the BASELINE is 0.131.

As same as the text query experiment, Both PRFL/PRF combined with BASELINE (system 7 and 8) achieved comparable IR performance with the baseline. Using global document information (system 6) is also effective. PRF/PRFL +GDOC method (system 4 and 5) are effective and system 4 achieved highest MAP (0.218). Iterative PRF (system 1 to 3) does not work well in this work.

Averaged precisions for each query of system 4 and 9 are also shown in Figure 4. For 74% queries (26 of 35 queries), system 4 outperforms baseline and underperforms for 6% queries.

In the task, PRFL effect seems to be smaller than in our past NTCIR SpokenDoc2 task. One possible reason is that the queries in the task are long and include a lot of words. For short queries QE is promising since significant words for retrieval are often omitted in a short query. On the other hand, long query may already include significant words, and it may weaken QE effects.

In spoken query experiment, PRF/PRFL seems to be more effective than in text query experiment. Spoken queries are transcribed by an ASR system and contain ASR errors. When a significant words are not recognized correctly in spoken query, we cannot find the collect document by the
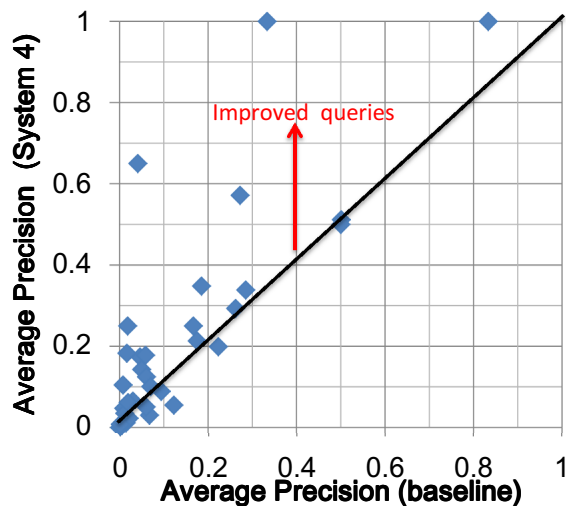
Figure 4: Comparison of proposed system and baseline in SGS SCR spoken query task

Table 1: passage retrieval result (text query)

| System ID | uMAP | pwMAP | fMAP |
|---|---|---|---|
| 1 | 0.028 | 0.114 | 0.043 |
| 2 | 0.029 | 0.117 | 0.044 |
| 3 | 0.029 | 0.115 | 0.044 |
| 4 | 0.029 | 0.116 | 0.044 |
| 5 | 0.029 | 0.116 | 0.044 |
| 6 | 0.032 | 0.125 | 0.045 |
| 7 | 0.018 | 0.085 | 0.032 |
| 8 | 0.018 | 0.085 | 0.032 |
| baseline | 0.021 | 0.090 | 0.034 |

Table 2: passage retrieval result (spoken query)

| System ID | uMAP | pwMAP | fMAP |
|---|---|---|---|
| 1 | 0.023 | 0.098 | 0.037 |
| 2 | 0.023 | 0.095 | 0.037 |
| 3 | 0.023 | 0.097 | 0.038 |
| 4 | 0.024 | 0.098 | 0.038 |
| 5 | 0.024 | 0.098 | 0.038 |
| 6 | 0.030 | 0.098 | 0.041 |
| 7 | 0.017 | 0.070 | 0.030 |
| 8 | 0.017 | 0.070 | 0.030 |
| baseline | 0.016 | 0.090 | 0.024 |

query. It can be one reason that the QE seems to be more effective in spoken query.

### 2.3.4 Passage Retrieval Systems

Next, passage retrieval systems are described. In a passage retrieval task, search target is a short part (utterance sequences of arbitrary length) in lectures. Although participants are requested to detect start and end points of such parts, determining such arbitrary length passages is time consuming. Here, we used a uniformly automatically segmented unit as a passage unit instead of such arbitrary length unit. Actually, we just divided oral presentation speech from its beginning into several segments which consist of N sequential utterances, which is introduced in the Japanese SCR test collection [7]. We regarded each 10 sequential utterance as a passage/document, and VSM-based SCR system is constructed.

The submitted systems are listed as follows.

1. PRFL+GDOC iter 3 (set $\lambda'$ to 0.6 (other systems: 0.5))
2. PRFL+GDOC iter 3
3. PRFL+GDOC iter 2
4. PRFL+GDOC
5. PRF+GDOC
6. GDOC
7. PRFL
8. PRF
9. BASELINE

Here, "BASELINE" system just searches passages with given query, and is the same with organizer's baseline system [5].

Integration of original query and expanded query similarities are performed in all systems 1 to 8. Also, all parameters were estimated so that higher 11ptAP was achieved at the past NTCIR SpokenDoc task.

For constructing indices, we used "REF-WORD-MATCH" transcription which is given by NTCIR-11 organizer. For queries, we used "REF-WORD-MATCH" and "MANUAL" which are given by NTCIR-11 organizer.

### 2.3.5 Passage Retrieval Results

Text query experiment

First, results of passage retrieval by text query are described. 10-utterance-based unit is regarded as a document, and for each text query, we tried to retrieve 1000 documents.

The results are listed in Table 1. uMAP, pwMAP, and fMAP of the baseline system are 0.021, 0.090, and 0.034, respectively. Our systems and baseline system always try to output 1000 of 10-sequential-utterances (total max 10000 utterances), therefore, it is difficult to achieve higher uMAP and fMAP. In our system, retrieved document is a uniformly divided segment, and we did not consider that the center of each segment is relevant to the query. Therefore, it is difficult to achieve higher pwMAP.

PRF and PRFL (system 8 and 7) are the same since they use 10-utterance-based document index for QE, and seemed not be effective. Using global document information (system 6) worked well and improved IR performance, and combinations with PRF techniques (system 1 to 5) do not work well in this task.

Spoken query experiment

Next, results of passage retrieval by spoken query are described. Also, 10-utterance-based unit is regarded as a document, and for each spoken query, we tried to retrieve 1000 documents. The results are listed in Table 2.

Almost the same results with text query results. System 6 shows the highest IR performance among the 9 systems, and achieves comparable IR performance with text query case in terms of uMAP and fMAP.

Here, we confirm that the effect of PRF seems to be small in long query and text query cases.

# 3. SPOKEN TERM DETECTION (STD)

## 3.1 STD Overview

Spoken term detection (STD) is the process finding the positions of a query term from the set of spoken documents. The common STD method is that ASR is performed first and perform text matching. We adopt the kind of STD system. In text matching process, phoneme is selected as a unit, that is, STD is performed based on a distance between document phoneme sequence and query phoneme sequence. As a distance measure, edit distance is adopted, which is easily calculated by continuous dynamic programing (CDP). In NTCIR-11 STD subtask, we conducted such CDP-based STD systems.

## 3.2 STD Techniques

### 3.2.1 Problems of CDP-based STD

CDP-based STD algorithm is described below.

— CDP-based STD algorithm —

1. Accept query $Q$ (length $L$)

2. Perfome STD with $Q$, and get STD result. Here, STD result consists of the file (lecture) ID-utterance ID list with edit distance information.

3. Calculate the score ($1 - \frac{\text{Edit Distance}}{L}$), which reflects a similarity between $Q$ and each spoken document. Then, list the results according to the score (decending order).

In CDP-based STD, there exist some problems. First one is homonym. Words which are pronounced in the same way are homonyms, and the all homonyms are falsely detected by CDP-based STD when phoneme sequence matching is performed. Another problem is substring matching. Especially, short queries face on the problem. Query phoneme sequence are unexpectedly matched with a part of other words, which is falsely detected. Also, due to ASR errors, queries are accidentally matched with other words, which is

falsely detected.

We try to reduce false detections. Specifically, we perform query expansion and suppress STD outputs which are likely misdetected ones according to STD results obtained by expanded queries.

### 3.2.2 Query Expansion for STD

Algorithm of STD with original and expanded queries "QE-STD" is described bellow. The overview is illustrated in Figure5.

— QE-STD algorithm —

1. Accept query $Q$ (length $L$)

2. Perform QE and generate expanded queries $QE_i (i = 1 \dots N)$.

3. Perform STD with original query $Q$, and get STD result $R_Q$. Here, STD result consists of the file (lecture) ID-utterance ID list with edit distance information.

4. Perform STD with expanded queries $QE_i (i = 1 \dots N)$, and get STD result $R_{QE}$. Here, STD result consists of file ID list.

5. For each element of $R_Q$, check that the same file ID is included in $R_{QE}$. If not included, that is, no expanded query and original query co-occurs in the same lecture, add penalty $p$ to its edit distance (modified edit distane).

6. Calculate the score ($1 - \frac{\text{Modified Edit Distance}}{L}$), which reflects a similarity between $Q$ and each spoken document. Then, list the results according to the score (decending order).

QE-STD reduces scores for spoken documents in which no expanded query and original query co-occurs, therefore, such spoken documents drop ranks in the STD output list. Otherwise, that if original query and at least one expanded query are both found in the same lecture, the QE-STD score is maintained.

### 3.2.3 Generation of Expanded Queries

For QE-STD, it is significant how to get appropriate expanded queries. Here, we propose that to make expanded queries by adding words to the original query. Specifically, words which likely occur in the preceding/following to the query word are added to the original query, thus, expanded queries are longer than the original query. Using such longer expanded queries, it is expected to reduce false detection for short queries.

We can select words which likely appear at the preceding/following to the target word in various ways. Here, we use 10 case-marker particles[1] in this work based on the assumption that a query word is often nouns, especially proper nouns, and at the preceding/following to nouns case-marker particles (kaku-joshi in Japanese) likely appear in Japanese.

Specifically, adding each case-marker particles to a head or tail of the original query word, we get 20 expanded queries



Figure 5: QE-STD overview
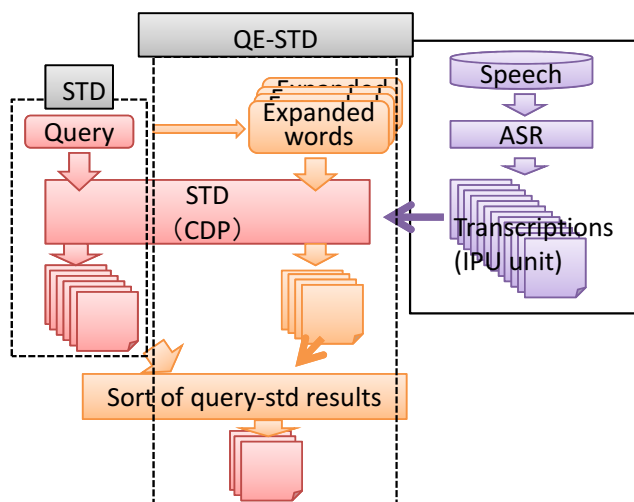
---

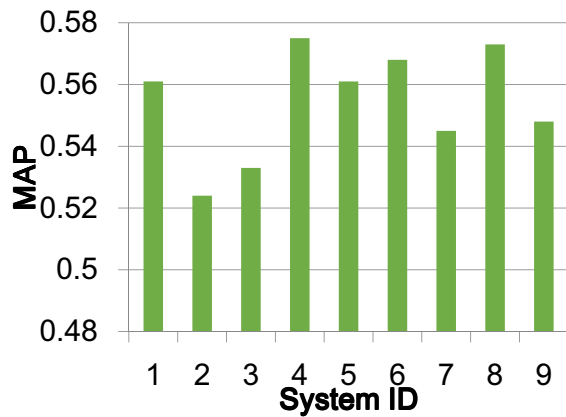[1] "ga", "no", "ni", "o", "e", "to", "de", "yori", "kara", and "ya"

Figure 6: STD results (text query)

for each query. For example, for query word "w e b u", tail-added expanded queries "w e b u g a", "w e b u n o", "w e b u n i", $\cdots$ and head-added expanded queries "g a w e b u", "n o w e b u", "n i w e b u", $\cdots$ are used.

### 3.2.4 STD systems

We used "REF-WORD-MATCH" transcription which is given by NTCIR-11 organizer. For queries, we used text query "MANUAL" which are given by NTCIR-11 organizer.

We set up following 9 systems. Here, baseline system marked "OrigQuery" just performs CDP-match between original query and each spoken document (IPU unit). System 1 to 8, which are marked "ExpQuery", perform STD with original and expanded queries. "penalty" is the parameter $p$ defined in QE-STD algorithm procedure 5, which is added to edit distance where no expanded query and original query co-occur in the same lecture.

1. Orig+ExpQuery (add to head and tail), penalty 2.5
2. Orig+ExpQuery (add to tail), penalty 2.5
3. Orig+ExpQuery (add to head), penalty 2.5
4. Orig+ExpQuery (add to head and tail), penalty 1.5
5. Orig+ExpQuery (add to tail), penalty 1.5
6. Orig+ExpQuery (add to head), penalty 1.5
7. Orig+ExpQuery (add to head and tail), penalty 2
8. Orig+ExpQuery (add to head and tail), penalty 1
9. OrigQuery (BASELINE)

### 3.2.5 STD result

For each query, we tried to detect 1000 utterances (IPU: Inter-Pausal Unit). In STD with expanded queries, we output results whose edit distance is equal to the minimum edit distance given by original query. Specifically, when M is the minimum edit distance for a given query in expanded queries search, only IPUs which have edit distance M are output.

The STD results are listed in Figure 6. QE-STD effect is confirmed. Especially, using both kind of expanded queries (head added and tail added expanded queries), higher MAPs are achieved than the case using one kind of the expanded queries. System 1 outperforms system 2 and 3, also system 4 outperforms system 5 and 6. Focusing on the penalty $p$, we got the highest MAP when $p$ is set to 1.5 (system 4).

## 4. CONCLUSIONS

We participated in NTCIR-11 "SpokenQuery&Doc" task as a team "RYSDT". In this paper, our SCR and STD systems were described. Vector space model based SCR systems with pseudo relevance feedback and integration of global document information are evaluated. CDP-based STD systems with query expansion are evaluated. We confirmed that our proposed methods are effective.

## 5. ACKNOWLEDGMENTS

## 6. REFERENCES

[1] Kazuyuki Noritake, Hiroaki Nanjo, and Takehiko Yoshimi. Image processing filters for line detection-based spoken term detection. In INTERSPEECH 2011, pages 2125–2128, 2011.

[2] Koji Shigeyasu, Hiroaki Nanjo, and Takehiko Yoshimi. A study of indexing units for japanese spoken document retrieval. In 10th Western Pacific Acoustics Conference (WESPAC X), 2009.

[3] Hiroaki Nanjo, Yusuke Iyonaga, and Takehiko Yoshimi. Spoken Document Retrieval for Oral Presentations Integrating Global Document Similarities into Local Document Similarities. In INTERSPEECH 2010, pages 1285–1288, 2010.

[4] Hiroaki Nanjo, Kazuyuki Noritake, and Takehiko Yoshimi. Spoken Document Retrieval Experiments for SpokenDocat Ryukoku University (RYSDT). In NTCIR-9, 2011.

[5] Tomoyosi Akiba, Hiromitsu Nishizaki, Hiroaki Nanjo, and Gareth J.F.Jones. Overview of the NTCIR-11 SpokenQuery&Doc task. In NTCIR-11, 2014.

[6] Hiroaki Nanjo, Tomohiro Nishio, and Takehiko Yoshimi. Spoken Document Retrieval Experiments for SpokenDoc-2 at Ryukoku University (RYSDT). In NTCIR-10, 2013.

[7] Tomoyosi Akiba, Kiyoaki Aikawa, Yoshiaki Itoh, Tatsuya Kawahara, Hiroaki Nanjo, Hiromitsu Nishizaki, Norihito Yasuda, Yoichi Yamashita, and Katunobu Itou. Construction of a test collection for spoken document retrieval from lecture audio data. Journal of Information Processing, 17:82–94, 2009.

[8] Akihiko Takano, Yoshiki Niwa, Shingo Nishioka, Makoto Iwayama, Toru Hisamitsu, Osamu Imaichi, and Hirofumi Sakurai. Information Access Based on Associative Calculation. In SOFSEM 2000: Theory and Practice of Informatics, Lecture Notes in Computer Science, pages 15–35, 2000.