

Finding Good Enough: A Task-Based Evaluation of Query Biased Summarization for Cross Language Information Retrieval

Jennifer Williams, Sharon Tam, Wade Shen

MIT Lincoln Laboratory Human Language Technology Group

244 Wood Street, Lexington, MA 02420 USA

jennifer.williams@ll.mit.edu, sharontam@alum.mit.edu
swade@ll.mit.edu

Abstract

In this paper we present our task-based evaluation of query biased summarization for cross-language information retrieval (CLIR) using relevance prediction. We describe our 13 summarization methods each from one of four summarization strategies. We show how well our methods perform using Farsi text from the CLEF 2008 shared-task, which we translated to English automatically. We report precision/recall/F1, accuracy and time-on-task. We found that different summarization methods perform optimally for different evaluation metrics, but overall query biased word clouds are the best summarization strategy. In our analysis, we demonstrate that using the ROUGE metric on our sentence-based summaries cannot make the same kinds of distinctions as our evaluation framework does. Finally, we present our recommendations for creating much-needed evaluation standards and datasets.

1 Introduction

Despite many recent advances in query biased summarization for cross-language information retrieval (CLIR), there are no existing evaluation standards or datasets to make comparisons among different methods, and across different languages (Tombros and Sanderson, 1998; Pingali et al., 2007; McCallum et al., 2012; Bhaskar and Bandyopadhyay, 2012). Consider that creating this kind of summary requires familiarity with techniques from machine translation (MT), summarization, and information retrieval (IR). In this

This work was sponsored by the Federal Bureau of Investigation under Air Force Contract FA8721-05-C-0002. Opinions, interpretations, conclusions, and recommendations are those of the authors and are not necessarily endorsed by the United States Government.

paper, we arrive at the intersection of each of these research areas. Query biased summarization (also known as *query-focused*, *query-relevant*, and *query-dependent*) involves automatically capturing relevant ideas and content from a document with respect to a given query, and presenting it as a condensed version of the original document. This kind of summarization is mostly used in search engines because when search results are tailored to a user's information need, the user can find texts that they are looking for more quickly and more accurately (Tombros and Sanderson, 1998; Mori et al., 2004). Query biased summarization is a valuable research area in natural language processing (NLP), especially for CLIR. Users of CLIR systems meet their information needs by submitting their queries in L_1 to search through documents that have been composed in L_2 , even though they may not be familiar with L_2 (Hovy et al., 1999; Pingali et al., 2007).

There are no standards for objectively evaluating summaries for CLIR – a research gap that we begin to address in this paper. The problem we explore is two-fold: what kinds of summaries are well-suited for CLIR applications, and how should the summaries be evaluated. Our evaluation is *extrinsic*, that is to say we are interested in how summarization affects performance on a different task (Mani et al., 2002; McKeown et al., 2005; Dorr et al., 2005; Murray et al., 2009; McCallum et al., 2012). We use relevance prediction as our *extrinsic* task: a human must decide if a summary for a given document is relevant to a particular information need, or not. Relevance prediction is known to be useful as it correlates with some automatic *intrinsic* methods as well (President and Dorr, 2006; Hobson et al., 2007). To the best of our knowledge, we are the first to apply this evaluation framework to cross language query biased summarization.

Each one of the summarization methods that we

present in this paper belongs to one of the following strategies: (1) unbiased full machine translated text, (2) unbiased word clouds, (3) query biased word clouds, and (4) query biased sentence summaries. The methods and strategies that we present are fast, cheap, and language-independent. All of these strategies are *extractive*, meaning that we used existing parts of a document to create the condensed version, or summary.

We approach our task as an engineering problem: the goal is to decide if summaries are good enough to help CLIR system users find what they are looking for. We have simplified the task by assuming that a set of documents has already been retrieved from a search engine, as CLIR techniques are outside the scope of this paper. We predict that showing the full MT English text as a summarization strategy would not be particularly helpful in our relevance prediction task because the words in the text could be mixed-up, or sentences could be nonsensical, resulting in poor readability. For the same reasons, we expect that showing the full MT English text would take longer to arrive at a relevance decision. Finally, we predict that query biased summaries will result in faster, more accurate decisions from the participants (Tombros and Sanderson, 1998).

We treat the actual CLIR search engine as if it were a black box so that we can focus on evaluating if the summaries themselves are useful. As a starting point, we begin with some principles that we expect to hold true when we evaluate. These principles provide us with the kind of framework that we need for a productive and judicious discussion about how well a summarization method works. We encourage the NLP community to consider the following concepts when developing evaluation standards for this problem:

- End-user intelligibility
- Query-salience
- Retrieval-relevance

Summaries should be presented to the end-user in a way that is both concise and intelligible, even if the machine translated text is difficult to understand. Our notions of *query-salience* and *retrieval-relevance* capture the expectation that good summaries will be efficient enough to help end-users fulfill their information needs. For query-salience, we want users to positively identify relevant documents. Similarly, for retrieval-relevance we want

users to be able to find as many relevant documents as possible.

This paper is structured as follows: Section 2 presents related work; Section 3 describes our data and pre-processing; Section 4 details our summarization methods and strategies; Section 5 describes our experiments; Section 6 shows our results and analysis; and in Section 7, we conclude and discuss some future directions for the NLP community.

2 Related Work

Automatic summarization is generally a well-investigated research area. Summarization is a way of describing the relationships of words in documents to the information content of that document (Luhn, 1958; Edmunson, 1969; Salton and Yang, 1973; Robertson and Walker, 1994; Church and Gale, 1999; Robertson, 2004). Recent work has looked at creating summaries of single and multiple documents (Radev et al., 2004; Erkan and Radev, 2004; Wan et al., 2007; Yin et al., 2012; Chatterjee et al., 2012), as well as summary evaluation (Jing et al., 1998; Tombros and Sanderson 1998; Mani et al., 1998; Mani et al., 1999; Mani, 2001; Lin and Hovy, 2003; Lin, 2004; Nenkova et al., 2007; Hobson et al., 2007; Owczarzak et al., 2012), query and topic biased summarization (Berger and Mittal, 2000; Otterbacher et al., 2005; Daume and Marcu, 2006; Chali and Joty, 2008; Otterbacher et al., 2009; Bando et al., 2010; Bhaskar and Bandyopadhyay, 2012; Harwath and Hazen, 2012; Yin et al., 2012), and summarization across languages (Pingali et al., 2007; Orăsan and Chiorean, 2008; Wan et al., 2010; Azarbondy et al., 2013).

2.1 Query Biased Summarization

Previous work most closely related to our own comes from Pingali et al., (2007). In their work, they present their method for cross-language query biased summarization for Telugu and English. Their work was motivated by the need for people to have access to foreign-language documents from a search engine even though the users were not familiar with the foreign language, in their case English. They used language modeling and translation probability to translate a user's query into L_2 , and then summarized each document in L_2 with respect to the query. In their final step, they translated the summary from L_2 back

to L_1 for the user. They evaluated their method on the DUC 2005 query-focused summarization shared-task with ROUGE scores. We compare our methods to this work also on the DUC 2005 task. Our work demonstrates the first attempt to draw a comparison between user-based studies and intrinsic evaluation with ROUGE. However, one of the limitations with evaluating this way is that the shared-task documents and queries are monolingual.

Bhaskar and Bandyopadhyay (2012) tried a subjective evaluation of extractive cross-language query biased summarization for 7 different languages. They extracted sentences, then scored and ranked the sentences to generate query dependent snippets of documents for their cross lingual information access (CLIA) system. However, the snippet quality was determined subjectively based on scores on a scale of 0 to 1 (with 1 being best). Each score indicated annotator satisfaction for a given snippet. Our evaluation methodology is objective: we ask users to decide if a given document is relevant to an information need, or not.

2.2 Machine Translation Effects

Machine translation quality can affect summarization quality. Wan et al. (2010) researched the effects of MT quality prediction on cross-language document summarization. They generated 5-sentence summaries in Chinese using English source documents. To select sentences, they used predicted translation quality, sentence position, and sentence informativeness. In their evaluation, they employed 4 Chinese-speakers to subjectively rate summaries on a 5-point scale (5 being best) along the dimensions of content, readability, and overall impression. They showed that their approach of using MT quality scores did improve summarization quality on average. While their findings are important, their work did not address query biasing or objective evaluation of the summaries. We attempt to overcome limitations of machine translation quality by using word clouds as one of our summarization strategies.

Knowing when to translate is another challenge for cross-language query biased summarization. Several options exist for when and what to translate during the summarization process: (1) the source documents can be translated, (2) the user's query can be translated, (3) the final summary can be translated, or (4) some combination of these.

An example of translating only the summaries themselves can be found in Wan et al., (2010). On the other hand, Pingali et al. (2007) translated the queries and the summaries. In our work, we used gold-translated queries from the CLEF 2008 dataset, and machine translated source documents. We briefly address this in our work, but note that a full discussion of when and what to translate, and those effects on summarization quality, is outside of the scope of this paper.

2.3 Summarization Evaluation

There has been a lot of work towards developing metrics for understanding what makes a summary good. Evaluation metrics are either *intrinsic* or *extrinsic*. Intrinsic metrics, such as ROUGE, measure the quality of a summary with respect to gold human-generated summaries (Lin, 2004; Lin and Hovy, 2003). Generating gold standard summaries is expensive and time-consuming, a problem that persists with cross-language query biased summarization because those summaries must be query biased as well as in a different language from the source documents.

On the other hand, extrinsic metrics measure the quality of summaries at the system level, by looking at overall system performance on downstream tasks (Jing et al, 1998; Tombros and Sanderson, 1998). One of the most important findings for query biased summarization comes from Tombros and Sanderson (1998). In their monolingual task-based evaluation, they measured user speed and accuracy at identifying relevant documents. They found that query biased summarization improved the user speed and accuracy when the user was asked to make relevance judgements for IR tasks. We also expect that our evaluation will demonstrate that user speed and accuracy is better when summaries are query biased.

3 Data and Pre-Processing

We used data from the Farsi CLEF 2008 ad hoc task (Agirre et al., 2009). Each of the queries included in this dataset consisted of a title, narrative, and description. Figure 1 shows an example of the elements of a CLEF 2008 query. All of the automatic query-biasing in this work was based on the query titles. For our human relevance prediction task on Mechanical Turk, we used the narrative version. The CLEF 2008 dataset included a ground-truth answer key indicating which docu-

ments were relevant to each query. For each query, we randomly selected 5 documents that were relevant as well as 5 documents that were not relevant. The subset of CLEF 2008 data that we used therefore consisted of 500 original Farsi documents and 50 parallel English-Farsi queries. Next we will describe our text pre-processing steps for both languages as well as how we created our parallel English documents.

```

Identifier: 10.2452/552-AH
Title: Tehran's stock market
Description: Find examples of the indexes of Tehran's stock market.
Narrative: Find information on fluctuations in indexes of the stock market, the top stock of the market, probable problems and challenges that the market has dealt with.

```

Figure 1: Full MT English summary and CLEF 2008 English query (title, description, narrative).

3.1 English Documents

All of our English documents were created automatically by translating the original Farsi documents into English (Drexler et al., 2012). The translated documents were sentence-aligned with one sentence per line. For all of our summarization experiments (except unbiased full MT text), we processed the text as follows: removed extra spaces, removed punctuation, folded to lowercase, and removed digits. We also removed common English stopwords² from the texts.

3.2 Farsi Documents

We used the original CLEF 2008 Farsi documents for two of our summarization methods. We stemmed words in each document using automatic morphological analysis with Morfessor CatMAP. We note that within-sentence punctuation was removed during this process (Creutz and Lagus, 2007). We also removed Farsi stopwords and digits.

4 Summarization Strategies

All of our summarization methods were extractive except for unbiased full machine translated text. In this section, we describe each of our 13 summarization methods which we have organized into one of the following strategies: (1) unbiased full machine translated text, (2) unbiased

²English and Farsi stopword lists from: <http://members.unine.ch/jacques.savoy/clef/index.html>

word cloud summaries, (3) query biased word cloud summaries, and (4) query biased sentence summaries. Regardless of which summarization method used, we highlighted words in yellow that also appeared in the query. Let t be a term in document d where $d \in D_L$ and D_L is a collection of documents in a particular language. Note that for our summarization methods, term weightings were calculated separately for each language. While $|D| = 1000$, we calculated term weightings based on $|D_E| = 500$ and $|D_F| = 500$. Finally, let q be a query where $q \in Q$ and Q is our set of 50 parallel English-Farsi CLEF queries. Assume that \log refers to \log_{10} .

```

H-770622-42472S8 (document):
the price changes in the stock market tehran
group economic: yesterday indicator of the
entire stock price than he made the remarks
337 increased and 151664... rate of also 30
increased and the last price yesterday 2102
each dollars. yesterday business in the stock
market 693 and a hundred thousand and one part
at a cost of 2 billion and 989 million and
452 and 839 by 361 of people like. yesterday
stock price changes in the following.
the prices of is:

10.2452/552-AH (query title):
Tehran's stock market

```

Figure 2: Full MT English summary and CLEF 2008 English query.

4.1 Unbiased Full Machine Translated English

Our first baseline approach was to use all of the raw machine translation output (no subsets of the sentences were used). Each summary therefore consisted of the full text of an entire document automatically translated from Farsi to English (Drexler et al., 2012). Figure 2 shows an example full text document translated from Farsi to English and a gold-standard English CLEF query. Note that we use this particular document-query pair as an example throughout this paper (document: H-770622-42472S8, query: 10.2452/552-AH). According to the CLEF answer key, the sample document is relevant to the sample query.

4.2 Unbiased Word Clouds

For our second baseline approach, we ranked terms in a document and displayed them as *word clouds*. Word clouds are one a way to arrange a collection of words where each word can vary

in size. We used word clouds as a summarization strategy to overcome any potential disfluencies from the machine translation output and also to see if they are feasible at all for summarization. All of our methods for word clouds used words from machine translated English text. Each term-ranking method below generates different ranked lists of terms, which we used to create different word clouds. We created one word cloud per document using the top 12 ranked words. We used the raw term scores to scale text font size, so that words with a higher score appeared larger and more prominent in a word cloud. Words were shuffled such that the exact ordering of words was at random.

I: Term Frequency (TF) Term frequency is very commonly used for finding important terms in a document. Given a term t in a document d , the number of times that term occurs is:

$$tf_{t,d} = |t \in d|$$

II: Inverse Document Frequency (IDF) The *idf* term weighting is typically used in IR and other text categorization tasks to make distinctions between documents. The version of *idf* that we used throughout our work came from Erkan and Radev (2004) and Otterbacher et al. (2009), in keeping consistent with theirs. Let N be the number of documents in the collection, such that $N = |D|$ and n_t is the number of documents that contain term t , such that $n_t = |\{d \in D : t \in d\}|$, then:

$$idf_t = \log \frac{N + 1}{0.5 \times n_t}$$

While *idf* is usually thought of as a type of heuristic, there have been some discussions about its theoretical basis (Robertson, 2004; Robertson and Walker, 1994; Church and Gale, 1999; Salton and Yang, 1973). An example of this summary is shown in Figure 3.

III: Term Frequency Inverse Document Frequency (TFIDF) We use $tfidf_{t,d}$ term weighting to find terms which are both rare and important for a document, with respect to terms across all other documents in the collection:

$$tfidf_{t,d} = tf_{t,d} \times idf_t$$

4.3 Query Biased Word Clouds

We generated query biased word clouds following the same principles as our unbiased word clouds,

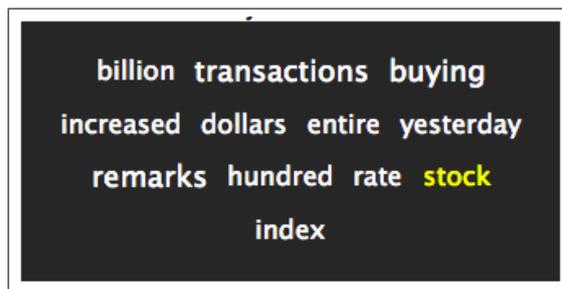


Figure 3: Word cloud summary for inverse document frequency (IDF), for query “Tehran’s stock market”.

namely the text font scaling and highlighting remained the same.

IV. Query Biased Term Frequency (TFQ) In Figure 4 we show a sample word cloud summary based on query biased term frequency. We define query biased term frequency tfQ at the document level, as:

$$tfQ_{t,d,q} = \begin{cases} 2tf_{t,d}, & \text{if } t \in q \\ tf_{t,d}, & \text{otherwise} \end{cases}$$



Figure 4: Word cloud summary for query biased term frequency (TFQ), for query “Tehran’s stock market”.

V. Query Biased Inverse Document Frequency (IDFQ) Since *idf* helps with identifying terms that discriminate documents in a collection, we would expect that query biased *idf* would help to identify documents that are relevant to a query:

$$idfQ_{t,q} = \begin{cases} 2idf_t, & \text{if } t \in q \\ idf_t, & \text{otherwise} \end{cases}$$

VI. Query Biased TFIDF (TFIDFQ) We define query biased $tf \times idf$ similarly to our TFQ and IDFQ, at the document level:

$$tfidfQ_{t,d,q} = \begin{cases} 2tf_{t,d} \times idf_t, & \text{if } t \in q \\ tf_{t,d} \times idf_t, & \text{otherwise} \end{cases}$$



Figure 5: Word cloud summary for scaled query biased term frequency (SFQ) for query “Tehran’s stock market”.

VII. Query Biased Scaled Frequency (SFQ)

This term weighting scheme, which we call scaled query biased term frequency or sfQ , is a variant of the traditional $tf \times idf$ weighting. First, we project the usual term frequency into log-space, for a term t in document d with:

$$tfS_{t,d} = \log(tf_{t,d})$$

We let $tfS_{t,d} \approx 0$ when $tf_{t,d} = 1$. We believe that singleton terms in a document provide no indication that a document is query-relevant, and treatment of singleton terms in this way would have the potential to reduce false-positives in our relevance prediction task. Note that scaled term frequency differs from Robertson’s (2004) *inverse total term frequency* in the sense that our method involves no consideration of term position within a document. Scaled query biased term frequency, shown in Figure 5, is defined as:

$$sfQ_{t,d,q} = \begin{cases} 2tfS_{t,d} \times idf_t, & \text{if } t \in q \\ tfS_{t,d} \times idf_t, & \text{otherwise} \end{cases}$$

VIII. Word Relevance (W) We adapted an existing relevance weighting from Allan et al., (2003), that was originally formulated for ranking sentences with respect to a query. However, we modified their original ranking method so that we could rank individual terms in a document instead of sentences. Our method for word relevance, W is defined as:

$$W_{t,d,q} = \log(tf_{t,d} + 1) \times \log(tf_{t,q} + 1) \times idf_t$$

In W , term frequency values are *smoothed* by adding 1. The smoothing could especially affect rare terms and singletons, when $tf_{t,d}$ is very

low. All terms in a query or a document will be weighted and each term could potentially contribute to summary.

4.4 Query Biased Sentence Summaries

Sentences are a canonical unit to use in extractive summaries. In this section we describe four different sentence scoring methods that we used. These methods show how to calculate sentence scores for a given document with respect to a given query. Sentences for a document were always ranked using the raw score value output generated from a scoring method. Each document summary contained the top 3 ranked sentences where the sentences were simply listed out. Each of these methods used sentence-aligned English machine translated documents, and two of them also used the original Farsi text.

IX. Sentence Relevance (REL) Our sentence relevance scoring method comes from Allan et al. (2003). The sentence weight is a summation over words that appear in the query. We provide their sentence scoring formula here. This calculates the relevance score for a sentence s from document d , to a query q :

$$rel_{(s|q)} = \sum_{t \in s} \log(tf_{t,s} + 1) \times \log(tf_{t,q} + 1) \times idf_t$$

Terms will occur in either the sentence or the query, or both. We applied this method to machine translated English text. The output of this method is a relevance score for each sentence in a given document. We used those scores to rank sentences in each document from our English machine translated text.

X. Query Biased Lexrank (LQ) We implemented query biased LexRank, a well-known graph-based summarization method (Otterbacher et al., 2009). It is a modified version of the original LexRank algorithm (Erkan and Radev, 2004; Page et al., 1998). The similarity metric, $sim_{x,y}$, also known as *idf-modified cosine similarity*, measures the distance between two sentences x and y in a document d , defined as:

$$sim_{x,y} = \frac{\sum_{t \in x,y} tf_{t,x} \times tf_{t,y} \times (idf_t)^2}{\sqrt{\sum_{t \in x} tf_{t,x} idf_{t,x}^2} \sqrt{\sum_{t \in y} tf_{t,y} idf_{t,y}^2}}$$

We used $sim_{x,y}$ to score the similarity of sentence-to-sentence, resulting in a similarity

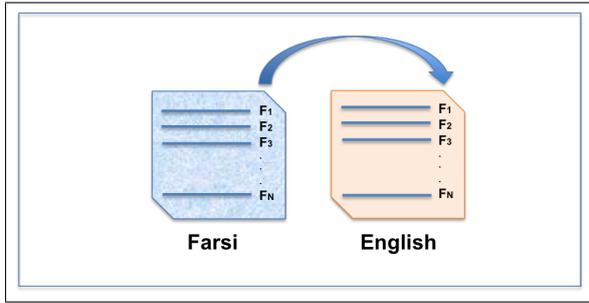


Figure 6: LQP - projecting Farsi sentence scores onto parallel English sentences.

graph where each vertex was a sentence and each edge was the cosine similarity between sentences. We normalized the cosine matrix with a similarity threshold ($t = 0.05$), so that sentences above this threshold were given similarity 1, and 0 otherwise. We used $rel_{s|q}$ to score sentence-to-query. The LexRank score for each sentence was then calculated as:

$$LQ_{s|q} = \frac{d \times rel_{s|q}}{\sum_{z \in C} rel_{z|q}} + (1 - d) \times \sum_{v \in adj[s]} \frac{sim_{s,v}}{\sum_{r \in adj[v]} sim_{v,r}} LQ_{v|q}$$

where C is the set of all sentences in a given document. Here the parameter d is just a damper to designate a probability of randomly jumping to one of the sentences in the graph ($d = 0.7$). We found the stationary distribution by applying the power method ($\epsilon = 5$), which is guaranteed to converge to a stationary distribution (Otterbacher et al., 2009). The output of LQ is a score for each sentence from a given document with respect to a query. We used that score to rank sentences in each document from our English machine translated text.

XI. Projected Cross-Language Query Biased Lexrank (LQP) We introduce LQP to describe a way of scoring and ranking sentences such that the L_1 (English) summaries are biased from the L_2 (Farsi) query and source document. Our gold-standard Farsi queries were included with our CLEF 2008 data, making them more reliable than what we could get from automatic translation. First, sentences from each Farsi document were scored with Farsi queries using LQ , described above. Then each LQ score was projected onto sentence-aligned English. We demonstrate LQP

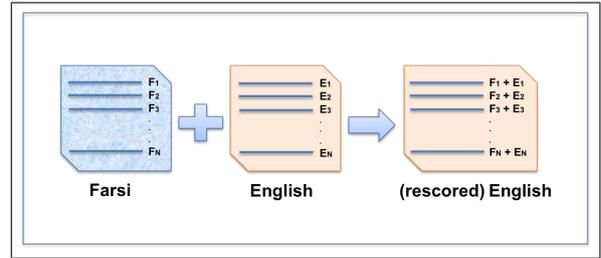


Figure 7: LQC - Farsi sentence scores are combined with parallel English sentence scores to obtain sentence re-ranking.

in Figure 6. By doing this, we simulated translating the user’s English query into Farsi with the best possible query translation, before proceeding with summarization. This approach to cross-language summarization could be of interest for CLIR systems that do query translation on-the-fly. It is also of interest for summarization systems that need to utilize previously translated source documents the capability is lacking to translate summaries from L_2 to L_1 .

XII. Combinatory Query Biased Lexrank (LQC) Another variation of LexRank that we introduce in this work is LQC , which combines LexRank scores from both languages to re-rank sentences. A visual summary of this method is shown in Figure 7. We accomplished our re-ranking by first running LQ on Farsi and English separately, then adding the two scores together. This combination of Farsi and English scores provided us with a different way to score and rank sentences, compared with LQ and LQP . The idea behind combinatory query biased LexRank is to take advantage of sentences which are high-ranking in Farsi but not in English. The LQC method exploits all available resources in our dataset: L_1 and L_2 queries as well as L_1 and L_2 documents.

5 Experiments

We tested each of our summarization methods and overall strategies in a task-based evaluation framework using relevance prediction. We used Mechanical Turk for our experiments since it has been shown to be useful for evaluating NLP systems (Callison-Burch 2009; Gillick and Liu, 2010). We obtained human judgments for whether or not a document was considered relevant to a query, or information need. We measured the relevance

judgements by precision/recall/F1, accuracy, and also time-on-task based on the average response time per Human Intelligence Task (HIT).

5.1 Mechanical Turk

In our Mechanical Turk experiment, we used terminology from CLEF 2008 to describe a query as an “information need”. All of the Mechanical Turk workers were presented with the following for their individual HIT: instructions, an information need and one summary for a document. Workers were asked to indicate if the given summary for a document was relevant to the given information need (Hobson et al., 2007). Workers were not shown the original Farsi source documents. We paid workers \$0.01 per HIT. We obtained 5 HITs for each information need and summary pair. We used a built-in approval rate qualification provided by Mechanical Turk to restrict which workers could work on our tasks. Each worker had an approval rate of at least 95

Instructions: Each image below consists of a statement summarizing the information you are trying to find from a set of documents followed by a summary of one of the documents returned when you query the documents. Based on the summary, choose whether you think the document returned is relevant to the information need. NOTE: It may be difficult to distinguish whether the document is relevant as the text may be difficult to understand. Just use your best judgment.

6 Results and Analysis

We present our experiment results and additional analysis. First, we report the results of our relevance prediction task, showing performance for individual summarization methods as well as performance for the overall strategies. Then we show analysis of our results from the monolingual question-biased shared-task for DUC 2005, as well as a comparison to previous work.

6.1 Results for Individual Methods

Our results are shown in Table 1. We report performance for 13 individual methods as well as overall performance on the 4 different summarization strategies. To calculate the performance for each

strategy, we used the arithmetic mean of the corresponding individual methods. We measured precision, recall and F1 to give us a sense of our summaries might influence document retrieval in an actual CLIR system. We also measured accuracy and time-on-task. For these latter two metrics, we distinguish between summaries that were relevant (R) and non-relevant (NR).

All of the summarization-based methods favored recall over precision: documents were marked ‘relevant’ more often than ‘non-relevant’. For many of the methods shown in Table 1, workers spent more time correctly deciding ‘relevant’ than correctly deciding ‘non-relevant’. This suggests some workers participated in our Mechanical Turk task purposefully. For many of the summarization methods, workers were able to positively identify relevant documents.

From Table 1 we see that Full MT performed better on precision than all of the other methods and strategies, but we note that performance on precision was generally very low. This might be due to Mechanical Turk workers overgeneralizing by marking summaries as relevant when they were not. Some individual methods preserve our principle of *retrieval-relevance*, as indicated by the higher recall scores for SQF, LQEF, and TFQ. That is to say, these particular query biased summarization methods can be used to assist users with identifying more relevant documents. The accuracy on relevant documents addresses our principle of *query-salience*, and it is especially high for our query-biased methods: LQEF, SQF, LQ, and TFQ. The results also seem to fit our intuition that the summary in Figure 3 seems less relevant to the summaries shown in Figures 4 & 5 even though these are the same documents biased on the same query “Tehran stock market”.

Overall, query biased word clouds outperform the other summarization strategies for 5 out of 7 metrics. This could be due to the fact that word clouds provide a very concise and overview of a document, which is one of the main goals for automatic summarization. Along these lines, word clouds are probably not subject to the effects of MT quality and we believe it is possible that MT quality could have had a negative impact on our query biased extracted sentence summaries, as well as our full MT English texts.

Table 1: Individual method results: precision/recall/F1, time-on-task, and accuracy. Note that results for time-on-task and accuracy scores are distinguished for relevant (R) and non-relevant (NR) documents.

Summarization Strategy	Precision, Recall, F1			Time-on-Task		Accuracy	
	Prec.	Rec.	F1	R	NR	R	NR
Unbiased Full MT English	0.653	0.636	0.644	219.5	77.6	0.696	0.712
TF	0.615	0.777	0.686	33.5	34.6	0.840	0.508
IDF	0.537	0.470	0.501	84.7	45.8	0.444	0.700
TFIDF	0.647	0.710	0.677	33.2	38.2	0.772	0.656
Unbiased Word Clouds	0.599	0.652	0.621	50.5	39.5	0.685	0.621
TFQ	0.605	0.809	0.692	55.3	82.4	0.864	0.436
IDFQ	0.582	0.793	0.671	23.6	31.6	0.844	0.436
TFIDFQ	0.599	0.738	0.661	37.9	26.9	0.804	0.500
SFQ	0.591	0.813	0.685	55.7	49.4	0.876	0.504
W	0.611	0.738	0.669	28.2	28.9	0.840	0.564
Query Biased Word Clouds	0.597	0.778	0.675	36.4	34.2	0.846	0.488
REL	0.582	0.746	0.654	30.6	44.3	0.832	0.548
LQ	0.549	0.783	0.646	64.4	54.8	0.868	0.292
LQP	0.578	0.734	0.647	28.2	28.0	0.768	0.472
LQC	0.557	0.810	0.660	33.9	38.8	0.896	0.292
Query Biased Sentences	0.566	0.768	0.651	39.2	41.5	0.841	0.401

Table 2: Comparison of peer systems on DUC 2005 shared-task for monolingual question-biased summarization, f-scores from ROUGE-2 and ROUGE-SU4.

Peer ID	ROUGE-2	ROUGE-SU4
17	0.07170	0.12970
8	0.06960	0.12790
4	0.06850	0.12770
Tel-Eng-Sum	0.06048	0.12058
LQ	0.05124	0.09343
REL	0.04914	0.09081

6.2 Analysis with DUC 2005

We analysed our summarization methods by comparing two of our sentence-based methods (LQ and REL) with peers from the monolingual question-biased summarization shared-task for DUC 2005. Even though DUC 2005 is a monolingual task, we decided to use it as part of our analysis for two reasons: (1) to see how well we could do with query/question biasing while ignoring the variables introduced by MT and cross-language text, and (2) to make a comparison to previous work. Pingali et al., (2007) also used this the same DUC task to assess their cross-language query biased summarization system. Systems

from the DUC 2005 question-biased summarization task were evaluated automatically against human gold-standard summaries using ROUGE (Lin and Hovy, 2003). Our results from the DUC 2005 shared-task are shown in Table 2, reported as ROUGE-2 and ROUGE-SU4 f-scores, as these two variations of ROUGE are the most helpful (Dang, 2005; Pingali et al., 2007).

Table 2 shows scores for several top peer systems, as well as results for the Tel-Eng-Sum method from Pingali et al., (2007). While we have reported f-scores in our analysis, we also note that our implementations of LQ and REL outperform all of the DUC 2005 peer systems for precision, as shown in Table 3. We also know that ROUGE cannot be used for comparing sentence summaries to ranked lists of words and there are no existing intrinsic methods to make that kind of comparison. Therefore we were able to successfully compare just 2 of our sentence-based methods to previous work using ROUGE.

7 Discussion and Future Work

Cross-language query biased summarization is an important part of CLIR, because it helps the user decide which foreign-language documents they might want to read. But, how do we know if

Table 3: Top 3 system precision scores for ROUGE-2 and ROUGE-SU4.

Peer ID	ROUGE-2	ROUGE-SU4
LQ	0.08272	0.15197
REL	0.0809	0.15049
15	0.07249	0.13129

a query biased summary is “good enough” to be used in a real-world CLIR system? We want to be able to say that we can do query biased summarization just as well for both monolingual and cross-language IR systems. From previous work, there has been some variability with regard to when and what to translate - variables which have no impact on monolingual summarization. We attempted to address this issue with two of our methods: LQP and LQC. To fully exploit the MT variable, we would need many more relevance prediction experiments using humans who know L_1 and others who know L_2 . Unfortunately in our case, we were not able to find Farsi speakers on Mechanical Turk. Access to these speakers would have allowed us to try further experiments as well as other kinds of analysis.

Our results on the relevance prediction task tell us that query biased summarization strategies help users identify relevant documents faster and with better accuracy than unbiased summaries. Our findings support the findings of Tombros and Sanderson (1998). Another important finding is that now we can weigh tradeoffs so that different summarization methods could be used to optimize over different metrics. For example, if we want to optimize for retrieval-relevance we might select a summarization method that tends to have higher recall, such as scaled query biased term frequency (SFQ). Similarly, we could optimize over accuracy on relevant documents, and use Combinatory LexRank (LQC) with Farsi and English together.

We have shown that the relevance prediction tasks can be crowdsourced on Mechanical Turk with reasonable results. The data we used from the Farsi CLEF 2008 ad-hoc task included an answer key, but there were no parallel English documents. However, in order for the NLP community to make strides in evaluating cross-language query biased summarization for CLIR, we will need standards and data. Optimal data would be parallel datasets consisting of documents in L_1 and L_2 with queries in L_1 and L_2 along with an answer

key specifying which documents are relevant to the queries. Further we would also need sets of human gold-standard query biased summaries in L_1 and L_2 . These standards and data would allow us to compare method-to-method across different languages, while simultaneously allowing us to tease apart other variables such as: when and what to translate, translation quality, methods for biasing, and type of summarization strategy (sentences, words, etc). And of course it would be better if this standard dataset was multilingual instead of bilingual, for obvious reasons.

We have approached cross-language query biased summarization as a stand-alone problem, treating the CLIR system and document retrieval as a black box. However, summaries need to preserve query-salience: summaries should not make it more difficult to positively identify relevant documents. And they should also preserve retrieval-relevance: summaries should help users identify as many relevant documents as possible.

Acknowledgments

We would like to express thanks to David Harwath at MIT Computer Science and Artificial Intelligence Laboratory (CSAIL), who helped us develop and implement ideas in this paper. We also want to thank Terry Gleason from MIT Lincoln Laboratory for providing machine translations.

References

- Eneko Agirre, Giorgio Maria Di Nunzio, Nicola Ferro, Thomas Mandl, and Carol Peters. CLEF 2008: Ad hoc track overview. In *Evaluating Systems for Multilingual and Multimodal Information Access*, pp 15–37. Springer Berlin Heidelberg, 2009.
- James Allan, Courtney Wade, and Alvaro Bolivar. Retrieval and Novelty Detection at the Sentence Level. In *Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Informaion Retrieval*, (SIGIR '03). ACM, New York, NY, USA, 314-321.
- Hosein Azarbyonad, Azadeh Shakery, and Hesham Faili. Exploiting Multiple Translation Resources for English-Persian Cross Language Information Retrieval. In P. Forner, H. Müller, R. Paredes, P. Rosso, and B. Stein, editors, *Information Access Evaluation. Multilinguality, Multimodality, and Visualization*, volume 8138 of *Lecture Notes in Computer Science*, pp 93–99. Springer Berlin Heidelberg, 2013.
- Lorena Leal Bando, Falk Scholer, Andrew Turpin. Constructing Query-biased Summaries: A Comparison of Human and System Generated Snippets. In

- Proceedings of the Third Symposium on Information Interaction in Context (IliX '10)*, ACM 2010, New York, NY, USA, 195–204.
- Adam Berger and Vibhu O Mittal. Query-Relevant Summarization Using FAQs. In *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, ACL 2000.
- Pinaki Bhaskar and Sivaji Bandyopadhyay. Cross-Lingual Query Dependent Snippet Generation. *International Journal of Computer Science and Information Technology (IJCSIT)*, 3(4), 2012.
- Pinaki Bhaskar and Sivaji Bandyopadhyay. Language Independent Query Focused Snippet Generation. In T. Catarci, P. Forner, D. Hiemstra, A. Peñas, and G. Santucci, editors, *Information Access Evaluation. Multilinguality, Multimodality, and Visual Analytics*, volume 7488 of *Lecture Notes in Computer Science*, pp 138–140. Springer Berlin Heidelberg, 2012.
- Stephen P. Borgatti, Kathleen M. Carley, David Krackhardt. On the Robustness of Centrality Measures Under Conditions of Imperfect Data. *Social Networks*, (28):124–136, 2006.
- Florian Boudin, Stéphane Huet, and Juan-Manuel Torres-Moreno. A Graph-Based Approach to Cross-Language Multi-Document Summarization. *Polibits*, (43):113–118, 2011.
- Chris Callison-Burch. Fast, Cheap, and Creative: Evaluating Translation Quality Using Amazon’s Mechanical Turk. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp 286–295, Singapore, ACL 2009.
- Yllias Chali and Shafiq R. Joty. Unsupervised Approach for Selecting Sentences in Query-Based Summarization. In *Proceedings of the Twenty-First International FLAIRS Conference*, 2008.
- Niladri Chatterjee, Amol Mittal, and Shubham Goyal. Single Document Extractive Text Summarization Using Genetic Algorithms. In *Emerging Applications of Information Technology (EAIT), 2012 Third International Conference*, pp 19–23, 2012.
- Kenneth W. Church and William A. Gale. Inverse Document Frequency (IDF): A Measure of Deviations From Poisson. In *Natural language processing using very large corpora*, pages 283–295. Springer, 1999.
- Mathias Creutz and Krista Lagus. Unsupervised Models for Morpheme Segmentation and Morphology Learning. *ACM Transactions on Speech and Language Processing*, 4(1):3:1–3:34, February 2007.
- Hoa Trang Dang. Overview of DUC 2005. In *Proceedings of the Document Understanding Conference*, 2005.
- Hal Daumé III, Daniel Marcu. Bayesian Query-Focused Summarization. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, ACL 2006.
- Jennifer Drexler, Wade Shen, Terry P. Gleason, Timothy R. Anderson, Raymond E. Slyh, Brian M. Ore, and Eric G. Hansen. The MIT-LL/AFRL IWSLT-2012 MT System. In *Proceedings of the International Workshop on Spoken Language Translation (IWSLT)*, Hong Kong, December 2012.
- Bonnie J. Dorr, Christof Monz, Stacy President, Richard Schwartz, and David Zajic. A Methodology for Extrinsic Evaluation of Text Summarization: Does ROUGE Correlate? In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pp 1-8. Ann Arbor, ACL 2005.
- H. P. Edmundson. New Methods in Automatic Extracting. In *Journal of the ACM*, 16(2):264–285, April 1969.
- Güneş Erkan and Dragomir R. Radev. Lexrank: Graph-based lexical centrality as salience in text summarization. *Journal of Artificial Intelligence Research*, 22(1):457–479, December 2004.
- Dan Gillick and Yang Liu. Non-Expert Evaluation of Summarization Systems is Risky. In *Proceedings of NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk*, pp 148-151, Los Angeles, California, USA, June, 2010.
- David Harwath and Timothy J. Hazen. Topic Identification Based Extrinsic Evaluation of Summarization Techniques Applied to Conversational Speech. In *Proceedings of ICASSP*, 2012: 5073-5076.
- Stacy P. Hobson, Bonnie J. Dorr, Christof Monz, and Richard Schwartz. Task-Eased Evaluation of Text Summarization Using Relevance Prediction. In *Information Processing Management*, 43(6): 1482-1499, 2007.
- Hongyan Jing, Regina Barzilay, Kathleen McKeown, and Michael Elhadad. Summarization Evaluation Methods: Experiments and Analysis. In *Proceedings of American Association for Artificial Intelligence (AAAI)*, 1998.
- Reza Karimpour, Amineh Ghorbani, Azadeh Pishdad, Mitra Mohtarami, Abolfazl AleAhmad, Hadi Amiri, and Farhad Oroumchian. Improving Persian Information Retrieval Systems Using Stemming and Part of Speech Tagging. In *Proceedings of the 9th Cross-language Evaluation Forum Conference on Evaluating Systems for Multilingual and Multimodal Information Access*, CLEF 2008, pp 89–96, Berlin, Heidelberg, 2009. Springer-Verlag.

- Chin-Yew Lin. Looking For A Few Good Metrics: Automatic Summarization Evaluation - How Many Samples Are Enough? In *Proceedings of NTCIR Workshop 4*, Tokyo, Japan, June 2004.
- Annie Louis and Ani Nenkova. Automatic Summary Evaluation without Human Models. In *Proceedings of Empirical Methods in Natural Language Processing*, EMNLP 2009.
- H. P. Luhn. The Automatic Creation of Literature Abstracts. *IBM Journal of Research and Development*, 2(2):159–165, April 1958.
- Inderjeet Mani, Eric Bloedorn, and Barbara Gates. Using Cohesion and Coherence Models for Text Summarization. In *AAAI Symposium Technical Report SS-989-06*, AAAI Press, 69–76, 1998.
- Inderjeet Mani, David House, Gary Klein, Lynette Hirschman, Therese Firmin, and Beth Sundheim. The TIPSTER SUMMAC Text Summarization Evaluation. In *Proceedings of European Association for Computational Linguistics*, EACL 1999.
- Inderjeet Mani. Summarization Evaluation: An Overview. In *Proceedings of the NTCIR Workshop*, Vol. 2, 2001.
- Inderjeet Mani, Gary Klein, David House, Lynette Hirschman, Therese Firmin, and Beth Sundheim. SUMMAC: A Text Summarization Evaluation. *Natural Language Engineering*, 8(1) 43-68. March 2002.
- Kathleen McKeown, Rebecca J. Passonneau, David K. Elson, Ani Nenkova, and Julia Hirschberg. Do Summaries Help? A Task-Based Evaluation of Multi-Document Summarization. In *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp 210-217. ACM 2005.
- Anthony McCallum, Gerald Penn, Cosmin Munteanu, and Xiaodan Zhu. Ecological Validity and the Evaluation of Speech Summarization Quality. In *Proceedings of Workshop on Evaluation Metrics and System Comparison for Automatic Summarization*. 2012 Association for Computational Linguistics, Stroudsburg, PA, USA, 28-35.
- Tatsunori Mori, Masanori Nozawa, and Yoshiaki Asada. Multi-Answer Focused Multi-Document Summarization Using a Question-Answering Engine. In *Proceedings of the 20th International Conference on Computational Linguistics*, COLING '04, Stroudsburg, PA, USA, ACL 2004.
- Gabriel Murray, Thomas Kleinbauer, Peter Poller, Tilman Becker, Steve Renals, and Jonathan Kilgour. Extrinsic Summarization Evaluation: A Decision Audit Task. *ACM Transactions on Speech and Language Processing*, 6(2) Article 2, October 2009.
- Ani Nenkova and Kathleen McKeown. A Survey of Text Summarization Techniques. In C. C. Aggarwal and C. Zhai, editors, *Mining Text Data*, pp 43–76. Springer US, 2012.
- Constantin Orăsan and Oana Andreea Chiorean. Evaluation of a Cross-Lingual Romanian-English Multi-Document Summariser. In *Proceedings of Language Resources and Evaluation Conference*, LREC 2008.
- Jahna Otterbacher, Güneş Erkan, and Dragomir R. Radev. Using Random Walks for Question-focused Sentence Retrieval. In *Proceedings of Human Language Technology Conference on Empirical Methods in Natural Language Processing*, Vancouver, Canada, pp 915-922, EMNLP 2005.
- Jahna Otterbacher, Güneş Erkan, and Dragomir R. Radev. Biased LexRank: Passage Retrieval Using Random Walks With Question-Based Priors. In *Information Processing Management*, 45(1), January 2009, pp 42-54.
- Karolina Owczarzak, John M. Conroy, Hoa Trang Dang, and Ani Nenkova. An Assessment of the Accuracy of Automatic Evaluation in Summarization. In *Proceedings of the Workshop on Evaluation Metrics and System Comparison for Automatic Summarization*, pp 1-9, Montréal, Canada, ACL 2012.
- Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The Pagerank Citation Ranking: Bringing Order to the Web. Technical report, Stanford Digital Library Technologies Project, 1998.
- Prasad Pingali, Jagadeesh Jagarlamudi, and Vasudeva Varma. Experiments in Cross Language Query Focused Multi-Document Summarization. In *Workshop on Cross Lingual Information Access Addressing the Information Need of Multilingual Societies*, IJCAI 2007.
- Stacy F. President and Bonnie J. Dorr. Text Summarization Evaluation: Correlating Human Performance on an Extrinsic Task With Automatic Intrinsic Metrics. No. LAMP-TR-133. University of Maryland College Park Language and Media Processing Laboratory Institute for Advanced Computer Studies (UMIACS), 2006.
- Dragomir R. Radev, Hongyan Jing, Malgorzata Styś, and Daniel Tam. Centroid-Based Summarization of Multiple Documents. In *Proceedings of Information Processing Management*, 40(6):919–938, Nov. 2004.
- Stephen Robertson. Understanding Inverse Document Frequency: on Theoretical Arguments for IDF. *Journal of Documentation*, 60(5):503–520, 2004.
- Stephen E. Robertson and Steve Walker. Some Simple Effective Approximations to the 2-Poisson Model for Probabilistic Weighted Retrieval. In *Proceedings of the 17th annual international ACM SIGIR*

conference on Research and development in information retrieval, pp 232–241. Springer-Verlag New York, Inc., 1994.

Gerard Salton and Chung-Shu Yang. On the Specification of Term Values in Automatic Indexing. *Journal of Documentation*, 29(4):351–372, 1973.

Anastasios Tombros and Mark Sanderson. Advantages of Query Biased Summaries in Information Retrieval. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pp 2–10. ACM, 1998.

Xiaojun Wan and Jianguo Xiao. Graph-Based Multi-Modality Learning for Topic-Focused Multi-Document Summarization. In *Proceedings of the 21st international joint conference on Artificial intelligence (IJCAI'09)*, San Francisco, CA, USA, 1586–1591.

Xiaojun Wan, Huiying Li, and Jianguo Xiao. Cross-Language Document Summarization Based on Machine Translation Quality Prediction. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics (ACL '10)*. Association for Computational Linguistics, Stroudsburg, PA, USA, 917–926.

Xiaojun Wan, Houping Jia, Shanshan Huang, and Jianguo Xiao. Summarizing the Differences in Multilingual News. In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '11*, pp 735–744, New York, NY, USA, 2011. ACM.

Wenpeng Yin, Yulong Pei, Fan Zhang, and Lian'en Huang. SentTopic-MultiRank: A Novel Ranking Model for Multi-Document Summarization. In *Proceedings of COLING*, pages 2977–2992, 2012.

Junlin Zhang, Le Sun, and Jinming Min. Using the Web Corpus to Translate the Queries in Cross-Lingual Information Retrieval. In *Proceedings in 2005 IEEE International Conference on Natural Language Processing and Knowledge Engineering, 2005, IEEE NLP-KE '05*, pp 493–498, 2005.